

---

# Infinite Task Learning in RKHSs

---

Romain Brault<sup>1†</sup> Alex Lambert<sup>2†</sup> Zoltán Szabó<sup>3</sup> Maxime Sangnier<sup>4</sup> Florence d’Alché-Buc<sup>2</sup>

<sup>1</sup>CentraleSupélec; <sup>2</sup>Télécom ParisTech, IP Paris; <sup>3</sup>École Polytechnique, IP Paris; <sup>4</sup>Sorbonne Université.

## Abstract

Machine learning has witnessed tremendous success in solving tasks depending on a single hyperparameter. When considering simultaneously a finite number of tasks, multi-task learning enables one to account for the similarities of the tasks via appropriate regularizers. A step further consists of learning a continuum of tasks for various loss functions. A promising approach, called *Parametric Task Learning*, has paved the way in the continuum setting for affine models and piecewise-linear loss functions. In this work, we introduce a novel approach called *Infinite Task Learning*: its goal is to learn a function whose output is a function over the hyperparameter space. We leverage tools from operator-valued kernels and the associated Vector-Valued Reproducing Kernel Hilbert Space that provide an explicit control over the role of the hyperparameters, and also allows us to consider new type of constraints. We provide generalization guarantees to the suggested scheme and illustrate its efficiency in cost-sensitive classification, quantile regression and density level set estimation.

## 1 INTRODUCTION

Several fundamental problems in machine learning and statistics can be phrased as the minimization of a loss function described by a hyperparameter. The hyperparameter might capture numerous aspects of the problem: (i) the tolerance w.r.t. outliers as the  $\epsilon$ -insensitivity in Support Vector Regression (Vapnik et al., 1997), (ii) importance of smoothness or sparsity such as the weight of the  $\ell_2$ -norm in

Proceedings of the 22<sup>nd</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2019, Naha, Okinawa, Japan. PMLR: Volume 89. Copyright 2019 by the author(s).

Tikhonov regularization (Tikhonov et al., 1977),  $\ell_1$ -norm in LASSO (Tibshirani, 1996), or more general structured-sparsity inducing norms (Bach et al., 2012), (iii) Density Level-Set Estimation (DLSE), see for example one-class support vector machines One-Class Support Vector Machine (OCSVM, Schölkopf et al., 2000), (iv) confidence as exemplified by Quantile Regression (QR, Koenker et al., 1978), or (v) importance of different decisions as implemented by Cost-Sensitive Classification (CSC, Zadrozny et al., 2001).

In various cases including QR, CSC or DLSE, one is interested in solving the parameterized task for several hyperparameter values. Multi-Task Learning (Evgeniou and Pontil, 2004) provides a principled way of benefiting from the relationship between similar tasks while preserving local properties of the algorithms:  $\nu$ -property in DLSE (Glazer et al., 2013) or quantile property in QR (Takeuchi, Le, et al., 2006).

A natural extension from the traditional multi-task setting is to provide a prediction tool being able to deal with *any* value of the hyperparameter. In their seminal work, (Takeuchi, Hongo, et al., 2013) extended multi-task learning by considering an infinite number of parametrized tasks in a framework called Parametric Task Learning (PTL). Assuming that the loss is piecewise affine in the hyperparameter, the authors are able to get the whole solution path through parametric programming, relying on techniques developed by Hastie et al. (2004).<sup>1</sup>

In this paper, we relax the affine model assumption on the tasks as well as the piecewise-linear assumption on the loss, and take a different angle. We propose Infinite Task Learning (ITL) within the framework of function-valued function learning to handle a continuum number of parameterized tasks. For that purpose we leverage tools from operator-valued kernels and the associated Vector-Valued Reproducing Kernel Hilbert Space (vv-RKHS, Pedrick, 1957). The idea is that

---

<sup>†</sup>Both authors contributed equally to this work.

<sup>1</sup>Alternative optimization techniques to deal with countable or continuous hyperparameter spaces could include semi-infinite (Stein, 2012) or bi-level programming (Wen et al., 1991).

the output is a function on the hyperparameters — modelled as scalar-valued Reproducing Kernel Hilbert Space (RKHS)—, which provides an explicit control over the role of the hyperparameters, and also enables us to consider new type of constraints. In the studied framework each task is described by a (scalar-valued) RKHS over the input space which is capable of dealing with nonlinearities. The resulting ITL formulation relying on vv-RKHS specifically encompasses existing multi-task approaches including joint quantile regression (Sangnier et al., 2016) or multi-task variants of density level set estimation (Glazer et al., 2013) by encoding a continuum of tasks.

Our **contributions** can be summarized as follows:

- We propose ITL, a novel vv-RKHS-based scheme to learn a continuum of tasks parametrized by a hyperparameter and design new regularizers.
- We prove excess risk bounds on ITL and illustrate its efficiency in quantile regression, cost-sensitive classification, and density level set estimation.

The paper is structured as follows. The ITL problem is defined in Section 2. In Section 3 we detail how the resulting learning problem can be tackled in vv-RKHSs. Excess risk bounds is the focus of Section 4. Numerical results are presented in Section 5. Conclusions are drawn in Section 6. Details of proofs are given in the supplement.

## 2 FROM PARAMETERIZED TO INFINITE TASK LEARNING

After introducing a few notations, we gradually define our goal by moving from single parameterized tasks (Section 2.1) to ITL (Section 2.3) through multi-task learning (Section 2.2).

**Notations:**  $\mathbb{1}_S$  is the indicator function of set  $S$ . We use the  $\sum_{i,j=1}^{n,m}$  shorthand for  $\sum_{i=1}^n \sum_{j=1}^m$ .  $|x|_+$  =  $\max(x, 0)$  denotes positive part.  $\mathcal{F}(\mathcal{X}; \mathcal{Y})$  stands for the set of  $\mathcal{X} \rightarrow \mathcal{Y}$  functions. Let  $\mathcal{Z}$  be Hilbert space and  $\mathcal{L}(\mathcal{Z})$  be the space of  $\mathcal{Z} \rightarrow \mathcal{Z}$  bounded linear operators. Let  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Z})$  be an operator-valued kernel, i.e.  $\sum_{i,j=1}^n \langle z_i, K(x_i, x_j) z_j \rangle_{\mathcal{Z}} \geq 0$  for all  $n \in \mathbb{N}^*$  and  $x_1, \dots, x_n \in \mathcal{X}$  and  $z_1, \dots, z_n \in \mathcal{Z}$  and  $K(x, z) = K(z, x)^*$  for all  $x, z \in \mathcal{X}$ .  $K$  gives rise to the Vector-Valued Reproducing Kernel Hilbert Space  $\mathcal{H}_K = \overline{\text{span}} \{K(\cdot, x)z \mid x \in \mathcal{X}, z \in \mathcal{Z}\} \subset \mathcal{F}(\mathcal{X}; \mathcal{Z})$ , where  $\overline{\text{span}} \{\cdot\}$  denotes the closure of the linear span of its argument. For further details on vv-RKHS the reader is referred to (Carmeli et al., 2010).

### 2.1 Learning Parameterized Tasks

A *supervised parametrized task* is defined as follows. Let  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$  be a random variable with joint distribution  $\mathbf{P}_{X,Y}$  which is assumed to be fixed but unknown; we also assume that  $\mathcal{Y} \subset \mathbb{R}$ . We have access to  $n$  independent identically distributed (i.i.d.) observations called training samples:  $\mathcal{S} := ((x_i, y_i))_{i=1}^n \sim \mathbf{P}_{X,Y}^{\otimes n}$ . Let  $\Theta$  be the domain of hyperparameters, and  $v_\theta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  be a loss function associated to  $\theta \in \Theta$ . Let  $\mathcal{H} \subset \mathcal{F}(\mathcal{X}; \mathcal{Y})$  denote our hypothesis class; throughout the paper  $\mathcal{H}$  is assumed to be a Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ . For a given  $\theta$ , the goal is to estimate the minimizer of the expected risk

$$\mathbf{R}^\theta(\mathbf{h}) := \mathbf{E}_{X,Y}[v_\theta(Y, \mathbf{h}(X))] \quad (1)$$

over  $\mathcal{H}$ , using the training sample  $\mathcal{S}$ . This task can be addressed by solving the regularized empirical risk minimization problem

$$\min_{\mathbf{h} \in \mathcal{H}} \mathbf{R}_S^\theta(\mathbf{h}) + \Omega(\mathbf{h}), \quad (2)$$

where  $\mathbf{R}_S^\theta(\mathbf{h}) := \frac{1}{n} \sum_{i=1}^n v_\theta(y_i, \mathbf{h}(x_i))$  is the empirical risk and  $\Omega : \mathcal{H} \rightarrow \mathbb{R}$  is a regularizer. Below we give three examples.

**Quantile Regression:** In this setting  $\theta \in (0, 1)$ . For a given hyperparameter  $\theta$ , in Quantile Regression the goal is to predict the  $\theta$ -quantile of the real-valued output conditional distribution  $\mathbf{P}_{Y|X}$ . The task can be tackled using the pinball loss (Koenker et al., 1978) defined in Eq. (3) and illustrated in Fig. S.3.

$$v_\theta(\mathbf{y}, \mathbf{h}(x)) = |\theta - \mathbb{1}_{\mathbb{R}_-}(\mathbf{y} - \mathbf{h}(x))| |\mathbf{y} - \mathbf{h}(x)|, \quad (3)$$

$$\Omega(\mathbf{h}) = \frac{\lambda}{2} \|\mathbf{h}\|_{\mathcal{H}}^2, \quad \lambda > 0.$$

**Cost-Sensitive Classification:** Our next example considers binary classification ( $\mathcal{Y} = \{-1, 1\}$ ) where a (possibly) different cost is associated with each class, as it is often the case in medical diagnosis. The sign of  $\mathbf{h} \in \mathcal{H}$  yields the estimated class and in cost-sensitive classification one takes

$$v_\theta(\mathbf{y}, \mathbf{h}(x)) = \left| \frac{1}{2}(\theta + 1) - \mathbb{1}_{\{-1\}}(\mathbf{y}) \right| |1 - \mathbf{y}\mathbf{h}(x)|_+, \quad (4)$$

$$\Omega(\mathbf{h}) = \frac{\lambda}{2} \|\mathbf{h}\|_{\mathcal{H}}^2, \quad \lambda > 0.$$

The  $\theta \in [-1, 1]$  hyperparameter captures the trade-off between the importance of correctly classifying the samples having  $-1$  and  $+1$  labels. When  $\theta$  is close to  $-1$ , the obtained  $\mathbf{h}$  focuses on classifying well class  $-1$ , and vice-versa. Typically, it is desirable for a physician to choose *a posteriori* the value of the hyperparameter at which he wants to predict. Since this cost can rarely be considered to be fixed, this motivates to learn one model giving access to all hyperparameter values.

**Density Level-Set Estimation:** Examples of parameterized tasks can also be found in the unsupervised setting. For instance in outlier detection, the goal is to separate outliers from inliers. A classical technique to tackle this task is OCSVM (Schölkopf et al., 2000). OCSVM has a free parameter  $\theta \in (0, 1]$ , which can be proven to be an upper bound on the fraction of outliers. When using a Gaussian kernel with a bandwidth tending towards zero, OCSVM consistently estimates density level sets (Vert et al., 2006). This unsupervised learning problem can be empirically described by the minimization of a regularized empirical risk  $\mathcal{R}_s^\theta(\mathbf{h}, \mathbf{t}) + \Omega(\mathbf{h})$ , solved *jointly* over  $\mathbf{h} \in \mathcal{H}$  and  $\mathbf{t} \in \mathbb{R}$  with

$$v_\theta(\mathbf{t}, \mathbf{h}(\mathbf{x})) = -\mathbf{t} + \frac{1}{\theta}|\mathbf{t} - \mathbf{h}(\mathbf{x})|_+, \quad \Omega(\mathbf{h}) = \frac{1}{2}\|\mathbf{h}\|_{\mathcal{H}}^2.$$

## 2.2 Solving a Finite Number of Tasks as Multi-Task Learning

In all the aforementioned problems, one is rarely interested in the choice of a single hyperparameter value ( $\theta$ ) and associated risk ( $\mathcal{R}_s^\theta$ ), but rather in the joint solution of multiple tasks. The naive approach of solving the different tasks independently can easily lead to inconsistencies. A principled way of solving many parameterized tasks has been cast as a MTL problem (Evgeniou, Micchelli, et al., 2005) which takes into account the similarities between tasks and helps providing consistent solutions. For example it is possible to encode the similarities of the different tasks in MTL through an explicit constraint function (Ciliberto et al., 2017). In the current work, the similarity between tasks is designed in an implicit way through the use of a kernel on the hyperparameters. Moreover, in contrast to MTL, in our case the input space and the training samples are the same for each task; a task is specified by a value of the hyperparameter. This setting is sometimes referred to as multi-output learning (Álvarez et al., 2012).

Formally, assume that we have  $p$  tasks described by parameters  $(\theta_j)_{j=1}^p$ . The idea of multi-task learning is to minimize the sum of the local loss functions  $\mathcal{R}_s^{\theta_j}$ , i. e.

$$\arg \min_{\mathbf{h}} \sum_{j=1}^p \mathcal{R}_s^{\theta_j}(\mathbf{h}_j) + \Omega(\mathbf{h}),$$

where the individual tasks are modelled by the real-valued  $\mathbf{h}_j$  functions, the overall  $\mathbb{R}^p$ -valued model is the vector-valued function  $\mathbf{h}: \mathbf{x} \mapsto (\mathbf{h}_1(\mathbf{x}), \dots, \mathbf{h}_p(\mathbf{x}))$ , and  $\Omega$  is a regularization term encoding similarities between tasks.

It is instructive to consider two concrete examples:

- In joint quantile regression one can use the regularizer to encourage that the predicted conditional

quantile estimates for two similar quantile values are similar. This idea forms the basis of the approach proposed by Sangnier et al. (2016) who formulates the joint quantile regression problem in a vector-valued Reproducing Kernel Hilbert Space with an appropriate decomposable kernel that encodes the links between the tasks. The obtained solution shows less quantile curve crossings compared to estimators not exploiting the dependencies of the tasks as well as an improved accuracy.

- A multi-task version of DLSE has recently been presented by Glazer et al. (2013) with the goal of obtaining nested density level sets as  $\theta$  grows. Similarly to joint quantile regression, it is crucial to take into account the similarities of the tasks in the joint model to efficiently solve this problem.

## 2.3 Towards Infinite Task Learning

In the following, we propose a novel framework called Infinite Task Learning in which we learn a function-valued function  $\mathbf{h} \in \mathcal{F}(\mathcal{X}; \mathcal{F}(\Theta; \mathcal{Y}))$ . Our goal is to be able to handle new tasks after the learning phase and thus, not to be limited to given predefined values of the hyperparameter. Regarding this goal, our framework generalizes the Parametric Task Learning approach introduced by Takeuchi, Hongo, et al. (2013), by allowing a wider class of models and relaxing the hypothesis of piece-wise linearity of the loss function. Moreover a nice byproduct of this vv-RKHS based approach is that one can benefit from the functional point of view, design new regularizers and impose various constraints on the whole continuum of tasks, e. g.,

- The continuity of the  $\theta \mapsto \mathbf{h}(\mathbf{x})(\theta)$  function is a natural desirable property: for a given input  $\mathbf{x}$ , the predictions on similar tasks should also be similar.
- Another example is to impose a shape constraint in QR: the conditional quantile should be increasing w. r. t. the hyperparameter  $\theta$ . This requirement can be imposed through the functional view of the problem.
- In DLSE, to get nested level sets, one would want that for all  $\mathbf{x} \in \mathcal{X}$ , the decision function  $\theta \mapsto \mathbb{1}_{\mathbb{R}_+}(\mathbf{h}(\mathbf{x})(\theta) - \mathbf{t}(\theta))$  changes its sign only once.

To keep the presentation simple, in the sequel we are going to focus on ITL in the supervised setting; unsupervised tasks can be handled similarly.

Assume that  $\mathbf{h}$  belongs to some space  $\mathcal{H} \subseteq \mathcal{F}(\mathcal{X}; \mathcal{F}(\Theta; \mathcal{Y}))$  and introduce an integrated loss function

$$V(\mathbf{y}, \mathbf{h}(\mathbf{x})) := \int_{\Theta} v(\theta, \mathbf{y}, \mathbf{h}(\mathbf{x})(\theta)) d\mu(\theta), \quad (5)$$

where the local loss  $v: \Theta \times \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  denotes  $v_\theta$  seen as a function of three variables including the hyperparameter and  $\mu$  is a probability measure on  $\Theta$  which encodes the importance of the prediction at different hyperparameter values. Without prior information and for compact  $\Theta$ , one may consider  $\mu$  to be uniform. The true risk reads then

$$\mathbf{R}(\mathbf{h}) := \mathbf{E}_{\mathcal{X}, \mathcal{Y}} [V(\mathbf{Y}, \mathbf{h}(\mathbf{X}))]. \quad (6)$$

Intuitively, minimizing the expectation of the integral over  $\theta$  in a rich enough space corresponds to searching for a pointwise minimizer  $\mathbf{x} \mapsto \mathbf{h}^*(\mathbf{x})(\theta)$  of the parametrized tasks introduced in Eq. (1) with, for instance, the implicit space constraint that  $\theta \mapsto \mathbf{h}^*(\mathbf{x})(\theta)$  is a continuous function for each input  $\mathbf{x}$ . We show in Proposition S.7.1 that this is precisely the case in QR.

Interestingly, the empirical counterpart of the true risk minimization can now be considered with a much richer family of penalty terms:

$$\min_{\mathbf{h} \in \mathcal{H}} \mathbf{R}_S(\mathbf{h}) + \Omega(\mathbf{h}), \quad \mathbf{R}_S(\mathbf{h}) := \frac{1}{n} \sum_{i=1}^n V(\mathbf{y}_i, \mathbf{h}(\mathbf{x}_i)). \quad (7)$$

Here,  $\Omega(\mathbf{h})$  can be a weighted sum of various penalties

- imposed directly on  $(\theta, \mathbf{x}) \mapsto \mathbf{h}(\mathbf{x})(\theta)$ , or
- integrated constraints on either  $\theta \mapsto \mathbf{h}(\mathbf{x})(\theta)$  or  $\mathbf{x} \mapsto \mathbf{h}(\mathbf{x})(\theta)$  such as

$$\int_{\mathcal{X}} \Omega_1(\mathbf{h}(\mathbf{x})(\cdot)) d\mathbf{P}(\mathbf{x}) \text{ or } \int_{\Theta} \Omega_2(\mathbf{h}(\cdot)(\theta)) d\mu(\theta)$$

which allow the property enforced by  $\Omega_1$  or  $\Omega_2$  to hold pointwise on  $\mathcal{X}$  or  $\Theta$  respectively.

It is worthwhile to see a concrete example before turning to the numerical solution (Section 3): in quantile regression, the monotonicity assumption of the  $\theta \mapsto \mathbf{h}(\mathbf{x})(\theta)$  function can be encoded by choosing  $\Omega_1$  as

$$\Omega_1(f) = \lambda_{\text{nc}} \int_{\Theta} |-(\partial f)(\theta)|_+ d\mu(\theta).$$

Many different models ( $\mathcal{H}$ ) could be applied to solve this problem. In our work we consider Reproducing Kernel Hilbert Spaces as they offer a simple and principled way to define regularizers by the appropriate choice of kernels and exhibit a significant flexibility.

### 3 SOLVING THE PROBLEM IN RKHSs

This section is dedicated to solving the ITL problem defined in Eq. (7). In Section 3.1 we focus on the objective ( $\tilde{\mathbf{V}}$ ). The applied vv-RKHS model family is detailed in Section 3.2 with various penalty examples followed by representer theorems, giving rise to computational tractability.

#### 3.1 Sampled Empirical Risk

In practice solving Eq. (7) can be rather challenging due to the integral over  $\theta$ . One might consider different numerical integration techniques to handle this issue. We focus here on Quasi Monte Carlo (QMC) methods<sup>2</sup> as they allow (i) efficient optimization over vv-RKHSs which we will use for modelling  $\mathcal{H}$  (Proposition 3.1), and (ii) enable us to derive generalization guarantees (Proposition 4.1). Indeed, let

$$\tilde{\mathbf{V}}(\mathbf{y}, \mathbf{h}(\mathbf{x})) := \sum_{j=1}^m w_j v(\theta_j, \mathbf{y}, \mathbf{h}(\mathbf{x})(\theta_j)) \quad (8)$$

be the QMC approximation of Eq. (5). Let  $w_j = m^{-1} F^{-1}(\theta_j)$ , and  $(\theta_j)_{j=1}^m$  be a sequence with values in  $[0, 1]^d$  such as the Sobol or Halton sequence where  $\mu$  is assumed to be absolutely continuous w. r. t. the Lebesgue measure and  $F$  is the associated cdf. Using this notation and the training samples  $\mathcal{S} = ((\mathbf{x}_i, \mathbf{y}_i))_{i=1}^n$ , the empirical risk takes the form

$$\tilde{\mathbf{R}}_S(\mathbf{h}) := \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{V}}(\mathbf{y}_i, \mathbf{h}(\mathbf{x}_i)) \quad (9)$$

and the problem to solve is

$$\min_{\mathbf{h} \in \mathcal{H}} \tilde{\mathbf{R}}_S(\mathbf{h}) + \Omega(\mathbf{h}). \quad (10)$$

#### 3.2 Hypothesis class ( $\mathcal{H}$ )

Recall that  $\mathcal{H} \subseteq \mathcal{F}(\mathcal{X}; \mathcal{F}(\Theta; \mathcal{Y}))$ , in other words  $\mathbf{h}(\mathbf{x})$  is a  $\Theta \mapsto \mathcal{Y}$  function for all  $\mathbf{x} \in \mathcal{X}$ . In this work we assume that the  $\Theta \mapsto \mathcal{Y}$  mapping can be described by an RKHS  $\mathcal{H}_{k_\Theta}$  associated to a  $k_\Theta: \Theta \times \Theta \rightarrow \mathbb{R}$  scalar-valued kernel defined on the hyperparameters. Let  $k_{\mathcal{X}}: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a scalar-valued kernel on the input space. The  $\mathbf{x} \mapsto$  (hyperparameter  $\mapsto$  output) relation, i.e.  $\mathbf{h}: \mathcal{X} \rightarrow \mathcal{H}_{k_\Theta}$  is then modelled by the Vector-Valued Reproducing Kernel Hilbert Space  $\mathcal{H}_{\mathbf{K}} = \overline{\text{span}} \{ \mathbf{K}(\cdot, \mathbf{x}) \mathbf{f} \mid \mathbf{x} \in \mathcal{X}, \mathbf{f} \in \mathcal{H}_{k_\Theta} \}$ , where the operator-valued kernel  $\mathbf{K}$  is defined as  $\mathbf{K}(\mathbf{x}, \mathbf{z}) = k_{\mathcal{X}}(\mathbf{x}, \mathbf{z}) \mathbf{I}$ , and  $\mathbf{I} = \mathbf{I}_{\mathcal{H}_{k_\Theta}}$  is the identity operator on  $\mathcal{H}_{k_\Theta}$ .

This so-called decomposable Operator-Valued Kernel has several benefits and gives rise to a function space with a well-known structure. One can consider elements  $\mathbf{h} \in \mathcal{H}_{\mathbf{K}}$  as mappings from  $\mathcal{X}$  to  $\mathcal{H}_{k_\Theta}$ , and also as functions from  $(\mathcal{X} \times \Theta)$  to  $\mathbb{R}$ . It is indeed known that there is an isometry between  $\mathcal{H}_{\mathbf{K}}$  and  $\mathcal{H}_{k_{\mathcal{X}}} \otimes \mathcal{H}_{k_\Theta}$ , the RKHS associated to the product kernel  $k_{\mathcal{X}} \otimes k_\Theta$ . The equivalence between these views allows a great flexibility and enables one to follow a functional point of view (to analyse statistical aspects) or to leverage the

<sup>2</sup>See Section S.10.1 of the supplement for a discussion on other integration techniques.

tensor product point of view (to design new kind of penalization schemes). Below we detail various regularizers before focusing on the representer theorems.

- **Ridge Penalty:** For QR and CSC, a natural regularization is the squared vv-RKHS norm

$$\Omega^{\text{RIDGE}}(\mathbf{h}) = \frac{\lambda}{2} \|\mathbf{h}\|_{\mathcal{H}_{\mathcal{K}}}^2, \quad \lambda > 0. \quad (11)$$

This choice is amenable to excess risk analysis (see Proposition 4.1). It can be also seen as the counterpart of the classical (multi-task regularization term introduced by Sangnier et al. (2016), compatible with an infinite number of tasks.  $\|\cdot\|_{\mathcal{H}_{\mathcal{K}}}^2$  acts by constraining the solution to a ball of a finite radius within the vv-RKHS, whose shape is controlled by both  $k_{\mathcal{X}}$  and  $k_{\Theta}$ .

- **$L^{2,1}$ -penalty:** For DLSE, it is more adequate to apply an  $L^{2,1}$ -RKHS mixed regularizer:

$$\Omega^{\text{DLSE}}(\mathbf{h}) = \frac{1}{2} \int_{\Theta} \|\mathbf{h}(\cdot)(\theta)\|_{\mathcal{H}_{k_{\mathcal{X}}}}^2 d\mu(\theta) \quad (12)$$

which is an example of a  $\Theta$ -integrated penalty. This  $\Omega$  choice allows the preservation of the  $\theta$ -property (see Fig. 2), i. e. that the proportion of the outliers is  $\theta$ .

- **Shape Constraints:** Taking the example of QR it is advantageous to ensure the monotonicity of the estimated quantile function. Let  $\partial_{\Theta} \mathbf{h}$  denotes the derivative of  $\mathbf{h}(x)(\theta)$  with respect to  $\theta$ . Then one should solve

$$\begin{aligned} & \arg \min_{\mathbf{h} \in \mathcal{H}_{\mathcal{K}}} \tilde{\mathcal{R}}_{\mathcal{S}}(\mathbf{h}) + \Omega^{\text{RIDGE}}(\mathbf{h}) \\ & \text{s. t. } \forall (x, \theta) \in \mathcal{X} \times \Theta, (\partial_{\Theta} \mathbf{h})(x)(\theta) \geq 0. \end{aligned}$$

However, the functional constraint prevents a tractable optimization scheme. To mitigate this bottleneck, we penalize if the derivative of  $\mathbf{h}$  w. r. t.  $\theta$  is negative:

$$\Omega_{\text{nc}}(\mathbf{h}) := \lambda_{\text{nc}} \int_{\mathcal{X}} \int_{\Theta} |-(\partial_{\Theta} \mathbf{h})(x)(\theta)|_+ d\mu(\theta) d\mathbf{P}(x). \quad (13)$$

When  $\mathbf{P} := \mathbf{P}_{\mathcal{X}}$  this penalization can rely on the same anchors and weights as the ones used to approximate the integrated loss function:

$$\tilde{\Omega}_{\text{nc}}(\mathbf{h}) = \lambda_{\text{nc}} \sum_{i,j=1}^{n,m} w_j |-(\partial_{\mathcal{X}} \mathbf{h})(x_i)(\theta_j)|_+. \quad (14)$$

Thus, one can modify the overall regularizer in QR to be

$$\Omega(\mathbf{h}) := \Omega^{\text{RIDGE}}(\mathbf{h}) + \tilde{\Omega}_{\text{nc}}(\mathbf{h}). \quad (15)$$

### 3.3 Representer Theorems

Apart from the flexibility of regularizer design, the other advantage of applying vv-RKHS as hypothesis

class is that it gives rise to finite-dimensional representation of the ITL solution under mild conditions. The representer theorem Proposition 3.1 applies to CSC when  $\lambda_{\text{nc}} = 0$  and to QR when  $\lambda_{\text{nc}} > 0$ .

**Proposition 3.1** (Representer). *Assume that for  $\forall \theta \in \Theta, v_{\theta}$  is a proper lower semicontinuous convex function with respect to its second argument. Then*

$$\arg \min_{\mathbf{h} \in \mathcal{H}_{\mathcal{K}}} \tilde{\mathcal{R}}_{\mathcal{S}}(\mathbf{h}) + \Omega(\mathbf{h}), \quad \lambda > 0$$

with  $\Omega(\mathbf{h})$  defined as in Eq. (15), has a unique solution  $\mathbf{h}^*$ , and  $\exists (\alpha_{ij})_{i,j=1}^{n,m}, (\beta_{ij})_{i,j=1}^{n,m} \in \mathbb{R}^{2nm}$  such that  $\forall x \in \mathcal{X}$

$$\mathbf{h}^*(x) = \sum_{i=1}^n k_{\mathcal{X}}(x, x_i) \left( \sum_{j=1}^m \alpha_{ij} k_{\Theta}(\cdot, \theta_j) + \beta_{ij} (\partial_2 k_{\Theta})(\cdot, \theta_j) \right).$$

**Sketch of the proof.** *First, we prove that the function to minimize is coercive, convex, lower semicontinuous, hence it has a unique minimum. Then  $\mathcal{H}_{\mathcal{K}}$  is decomposed into two orthogonal subspaces and we use the reproducing property to get the finite representation.*

For DLSE, we similarly get a representer theorem with the following modelling choice. Let  $k_{\mathcal{B}} : \Theta \times \Theta \rightarrow \mathbb{R}$  be a scalar-valued kernel (possibly different from  $k_{\Theta}$ ),  $\mathcal{H}_{k_{\mathcal{B}}}$  the associated RKHS and  $\mathbf{t} \in \mathcal{H}_{k_{\mathcal{B}}}$ . Assume also that  $\Theta \subseteq [\epsilon, 1]$  where  $\epsilon > 0$ .<sup>3</sup> Then, learning a continuum of level sets boils down to the minimization problem

$$\arg \min_{\mathbf{h} \in \mathcal{H}_{\mathcal{K}}, \mathbf{t} \in \mathcal{H}_{k_{\mathcal{B}}}} \tilde{\mathcal{R}}_{\mathcal{S}}(\mathbf{h}, \mathbf{t}) + \tilde{\Omega}(\mathbf{h}, \mathbf{t}), \quad \lambda > 0, \quad (16)$$

where  $\tilde{\Omega}(\mathbf{h}, \mathbf{t}) = \frac{1}{2} \sum_{j=1}^m w_j \|\mathbf{h}(\cdot)(\theta_j)\|_{\mathcal{H}_{k_{\mathcal{X}}}}^2 + \frac{\lambda}{2} \|\mathbf{t}\|_{\mathcal{H}_{k_{\mathcal{B}}}}^2$ ,  $\tilde{\mathcal{R}}_{\mathcal{S}}(\mathbf{h}, \mathbf{t}) = \frac{1}{n} \sum_{i,j=1}^{n,m} \frac{w_i}{\theta_j} (|\mathbf{t}(\theta_j) - \mathbf{h}(x_i)(\theta_j)|_+ - \mathbf{t}(\theta_j))$ .

**Proposition 3.2** (Representer). *Assume that  $k_{\Theta}$  is bounded:  $\sup_{\theta \in \Theta} k_{\Theta}(\theta, \theta) < +\infty$ . Then the minimization problem described in Eq. (16) has a unique solution  $(\mathbf{h}^*, \mathbf{t}^*)$  and there exist  $(\alpha_{ij})_{i,j=1}^{n,m} \in \mathbb{R}^{n \times m}$  and  $(\beta_j)_{j=1}^m \in \mathbb{R}^m$  such that for  $\forall (x, \theta) \in \mathcal{X} \times [\epsilon, 1]$ ,*

$$\begin{aligned} \mathbf{h}^*(x)(\theta) &= \sum_{i,j=1}^{n,m} \alpha_{ij} k_{\mathcal{X}}(x, x_i) k_{\Theta}(\theta, \theta_j), \\ \mathbf{t}^*(\theta) &= \sum_{j=1}^m \beta_j k_{\mathcal{B}}(\theta, \theta_j). \end{aligned}$$

**Sketch of the proof.** *First we show that the infimum exists, and that it must be attained in some subspace of  $\mathcal{H}_{\mathcal{K}} \times \mathcal{H}_{k_{\mathcal{B}}}$  over which the objective function is coercive. By the reproducing property, we get the claimed finite decomposition.*

#### Remarks:

<sup>3</sup>We choose  $\Theta \subseteq [\epsilon, 1]$ ,  $\epsilon > 0$  rather than  $\Theta \subseteq [0, 1]$  because the loss might not be integrable on  $[0, 1]$ .

Table 1: Quantile Regression on 20 UCI datasets. Reported:  $100 \times$  value of the pinball loss,  $100 \times$  crossing loss (smaller is better). p-val.: outcome of the Mann-Whitney-Wilcoxon test of JQR vs.  $\infty$ -QR and Independent vs.  $\infty$ -QR. Boldface: significant values w. r. t.  $\infty$ -QR.

DATASET	JQR				IND-QR				$\infty$ -QR	
	(PINBALL)	P.-VAL.)	(CROSS	P.-VAL.)	(PINBALL	P.-VAL.)	(CROSS	P.-VAL.)	PINBALL	CROSS
COBARORE	159 ± 24	$9 \cdot 10^{-01}$	0.1 ± 0.4	$6 \cdot 10^{-01}$	150 ± 21	$2 \cdot 10^{-01}$	0.3 ± 0.8	$7 \cdot 10^{-01}$	165 ± 36	2.0 ± 6.0
ENGEL	175 ± 555	$6 \cdot 10^{-01}$	0.0 ± 0.2	$1 \cdot 10^{+00}$	63 ± 53	$8 \cdot 10^{-01}$	4.0 ± 12.8	$8 \cdot 10^{-01}$	47 ± 6	0.0 ± 0.1
BOSTONHOUSING	49 ± 4	$8 \cdot 10^{-01}$	0.7 ± 0.7	$2 \cdot 10^{-01}$	49 ± 4	$8 \cdot 10^{-01}$	<b>1.3 ± 1.2</b>	$1 \cdot 10^{-05}$	49 ± 4	0.3 ± 0.5
CAUTION	88 ± 17	$6 \cdot 10^{-01}$	0.1 ± 0.2	$6 \cdot 10^{-01}$	89 ± 19	$4 \cdot 10^{-01}$	<b>0.3 ± 0.4</b>	$2 \cdot 10^{-04}$	85 ± 16	0.0 ± 0.1
FTCOLLINSNOW	154 ± 16	$8 \cdot 10^{-01}$	0.0 ± 0.0	$6 \cdot 10^{-01}$	155 ± 13	$9 \cdot 10^{-01}$	0.2 ± 0.9	$8 \cdot 10^{-01}$	156 ± 17	0.1 ± 0.6
HIGHWAY	103 ± 19	$4 \cdot 10^{-01}$	0.8 ± 1.4	$2 \cdot 10^{-02}$	99 ± 20	$9 \cdot 10^{-01}$	<b>6.2 ± 4.1</b>	$1 \cdot 10^{-07}$	105 ± 36	0.1 ± 0.4
HEIGHTS	127 ± 3	$1 \cdot 10^{+00}$	0.0 ± 0.0	$1 \cdot 10^{+00}$	127 ± 3	$9 \cdot 10^{-01}$	0.0 ± 0.0	$1 \cdot 10^{+00}$	127 ± 3	0.0 ± 0.0
SNIFFER	43 ± 6	$8 \cdot 10^{-01}$	0.1 ± 0.3	$2 \cdot 10^{-01}$	44 ± 5	$7 \cdot 10^{-01}$	<b>1.4 ± 1.2</b>	$6 \cdot 10^{-07}$	44 ± 7	0.1 ± 0.1
SNOWGEESE	55 ± 20	$7 \cdot 10^{-01}$	0.3 ± 0.8	$3 \cdot 10^{-01}$	53 ± 18	$6 \cdot 10^{-01}$	0.4 ± 1.0	$5 \cdot 10^{-02}$	57 ± 20	0.2 ± 0.6
UFC	81 ± 5	$6 \cdot 10^{-01}$	<b>0.0 ± 0.0</b>	$4 \cdot 10^{-04}$	82 ± 5	$7 \cdot 10^{-01}$	<b>1.0 ± 1.4</b>	$2 \cdot 10^{-04}$	82 ± 4	0.1 ± 0.3
BIGMAC2003	80 ± 21	$7 \cdot 10^{-01}$	<b>1.4 ± 2.1</b>	$4 \cdot 10^{-04}$	74 ± 24	$9 \cdot 10^{-02}$	<b>0.9 ± 1.1</b>	$7 \cdot 10^{-05}$	84 ± 24	0.2 ± 0.4
UN3	98 ± 9	$8 \cdot 10^{-01}$	0.0 ± 0.0	$1 \cdot 10^{-01}$	99 ± 9	$1 \cdot 10^{+00}$	<b>1.2 ± 1.0</b>	$1 \cdot 10^{-05}$	99 ± 10	0.1 ± 0.4
BIRTHWT	141 ± 13	$1 \cdot 10^{+00}$	0.0 ± 0.0	$6 \cdot 10^{-01}$	140 ± 12	$9 \cdot 10^{-01}$	0.1 ± 0.2	$7 \cdot 10^{-02}$	141 ± 12	0.0 ± 0.0
CRABS	11 ± 1	$4 \cdot 10^{-05}$	0.0 ± 0.0	$8 \cdot 10^{-01}$	11 ± 1	$2 \cdot 10^{-04}$	<b>0.0 ± 0.0</b>	$2 \cdot 10^{-05}$	13 ± 3	0.0 ± 0.0
GAGURINE	61 ± 7	$4 \cdot 10^{-01}$	0.0 ± 0.1	$3 \cdot 10^{-03}$	62 ± 7	$5 \cdot 10^{-01}$	<b>0.1 ± 0.2</b>	$4 \cdot 10^{-04}$	62 ± 7	0.0 ± 0.0
GEYSER	105 ± 7	$9 \cdot 10^{-01}$	0.1 ± 0.3	$9 \cdot 10^{-01}$	105 ± 6	$9 \cdot 10^{-01}$	0.2 ± 0.3	$6 \cdot 10^{-01}$	104 ± 6	0.1 ± 0.2
GILGAIS	51 ± 6	$5 \cdot 10^{-01}$	0.1 ± 0.1	$1 \cdot 10^{-01}$	49 ± 6	$6 \cdot 10^{-01}$	<b>1.1 ± 0.7</b>	$2 \cdot 10^{-05}$	49 ± 7	0.3 ± 0.3
TOPO	69 ± 18	$1 \cdot 10^{+00}$	0.1 ± 0.5	$1 \cdot 10^{+00}$	71 ± 20	$1 \cdot 10^{+00}$	<b>1.7 ± 1.4</b>	$3 \cdot 10^{-07}$	70 ± 17	0.0 ± 0.0
MCYCLE	66 ± 9	$9 \cdot 10^{-01}$	0.2 ± 0.3	$7 \cdot 10^{-03}$	66 ± 8	$9 \cdot 10^{-01}$	<b>0.3 ± 0.3</b>	$7 \cdot 10^{-06}$	65 ± 9	0.0 ± 0.1
CPUS	<b>7 ± 4</b>	$2 \cdot 10^{-04}$	<b>0.7 ± 1.0</b>	$5 \cdot 10^{-04}$	<b>7 ± 5</b>	$3 \cdot 10^{-04}$	<b>1.2 ± 0.8</b>	$6 \cdot 10^{-08}$	16 ± 10	0.0 ± 0.0

- Models with bias: it can be advantageous to add a bias to the model, which is here a function of the hyperparameter  $\theta$ :  $h(x)(\theta) = f(x)(\theta) + b(\theta)$ ,  $f \in \mathcal{H}_{\mathcal{K}}$ ,  $b \in \mathcal{H}_{\mathcal{K}_b}$ , where  $\mathcal{K}_b : \Theta \times \Theta \rightarrow \mathbb{R}$  is a scalar-valued kernel. This can be the case for example if the kernel on the hyperparameters is the constant kernel, i.e.  $k_{\Theta}(\theta, \theta') = 1$  ( $\forall \theta, \theta' \in \Theta$ ), hence the model  $f(x)(\theta)$  would not depend on  $\theta$ . An analogous statement to Proposition 3.1 still holds for the biased model if one adds a regularization  $\lambda_b \|b\|_{\mathcal{H}_{\mathcal{K}_b}}^2$ ,  $\lambda_b > 0$  to the risk.
- Relation to JQR: In  $\infty$ -QR, by choosing  $k_{\Theta}$  to be the Gaussian kernel,  $\mathcal{K}_b(x, z) = \mathbb{1}_{\{x\}}(z)$ ,  $\mu = \frac{1}{m} \sum_{j=1}^m \delta_{\theta_j}$ , where  $\delta_{\theta}$  is the Dirac measure concentrated on  $\theta$ , one gets back Sangnier et al. (2016)’s Joint Quantile Regression (JQR) framework as a special case of our approach. In contrast to the JQR, however, in  $\infty$ -QR one can predict the quantile value at any  $\theta \in (0, 1)$ , even outside the  $(\theta_j)_{j=1}^m$  used for learning.
- Relation to q-OCSVM: In DLSE, by choosing  $k_{\Theta}(\theta, \theta') = 1$  (for all  $\theta, \theta' \in \Theta$ ) to be the constant kernel,  $\mathcal{K}_b(\theta, \theta') = \mathbb{1}_{\{\theta\}}(\theta')$ ,  $\mu = \frac{1}{m} \sum_{j=1}^m \delta_{\theta_j}$ , our approach specializes to q-OCSVM (Glazer et al., 2013).
- Relation to Kadri et al. (2016): Note that Operator-Valued Kernels for functional outputs have also been used in (Kadri et al., 2016), under the form of

integral operators acting on  $L^2$  spaces. Both kernels give rise to the same space of functions; the benefit of our approach being to provide an *exact* finite representation of the solution (see Proposition 3.1).

- Efficiency of the decomposable kernel: this kernel choice allows to rewrite the expansions in Propositions 3.1 and 3.2 as a Kronecker products and the complexity of the prediction of  $n'$  points for  $m'$  quantile becomes  $\mathcal{O}(m'mn + n'n m)$  instead of  $\mathcal{O}(m'mn'n)$ .

## 4 Excess Risk Bounds

Below we provide a generalization error analysis to the solution of Eq. (10) for QR and CSC (with Ridge regularization and without shape constraints) by stability argument (Bousquet et al., 2002), extending the work of Audiffren et al. (2013) to Infinite-Task Learning. The proposition (finite sample bounds are given in Corollary S.9.6) instantiates the guarantee for the QMC scheme.

**Proposition 4.1 (Generalization).** *Let  $h^* \in \mathcal{H}_{\mathcal{K}}$  be the solution of Eq. (10) for the QR or CSC problem with QMC approximation. Under mild conditions on the kernels  $\mathcal{K}_X, \mathcal{K}_{\Theta}$  and  $\mathbf{P}_{X,Y}$ , stated in the supplement, one has*

$$\mathbb{R}(h^*) \leq \tilde{\mathbb{R}}_s(h^*) + \mathcal{O}_{\mathbf{P}_{X,Y}} \left( \frac{1}{\sqrt{\lambda n}} \right) + \mathcal{O} \left( \frac{\log(m)}{\sqrt{\lambda m}} \right). \quad (17)$$

**Sketch of the proof.** *The error resulting from sam-*

pling  $\mathbf{P}_{X,Y}$  and the inexact integration is respectively bounded by  $\beta$ -stability (Kadri et al., 2016) and QMC results.<sup>4</sup>

**( $\mathbf{n}, \mathbf{m}$ ) Trade-off:** The proposition reveals the interplay between the two approximations,  $\mathbf{n}$  (the number of training samples) and  $\mathbf{m}$  (the number of locations taken in the integral approximation), and allows to identify the regime in  $\lambda = \lambda(\mathbf{n}, \mathbf{m})$  driving the excess risk to zero. Indeed by choosing  $\mathbf{m} = \sqrt{\mathbf{n}}$  and discarding logarithmic factors for simplicity,  $\lambda \gg \mathbf{n}^{-1}$  is sufficient. The mild assumptions imposed are: boundedness on both kernels and the random variable  $Y$ , as well as some smoothness of the kernels.

## 5 Numerical Examples

In this section we provide numerical examples illustrating the efficiency of the proposed ITL approach.<sup>5</sup> We used the following datasets in our experiments:

- **Quantile Regression:** we used (i) a sine synthetic benchmark (Sangnier et al., 2016): a sine curve at 1Hz modulated by a sine envelope at 1/3Hz and mean 1, distorted with a Gaussian noise of mean 0 and a linearly decreasing standard deviation from 1.2 at  $x = 0$  to 0.2 at  $x = 1.5$ . (ii) 20 standard regression datasets from UCI. The number of samples varied between 38 (CobarOre) and 1375 (Height). The observations were standardised to have unit variance and zero mean for each attribute.
- **Density Level-Set Estimation:** The Wilt database from the UCI repository with 4839 samples and 5 attributes, and the Spambase UCI dataset with 4601 samples and 57 attributes served as benchmarks.

Additional experiments related to the CSC problem are provided in Section S.10.5.

**Note on Optimization:** There are several ways to solve the non-smooth optimization problems associated to the QR, DLSE and CSC tasks. One could proceed for example by duality—as it was done in JQR Sangnier et al. (2016)—, or apply sub-gradient descent techniques (which often converge quite slowly). In order to allow unified treatment and efficient solution in our experiments we used the L-BFGS-B (Zhu et al., 1997) optimization scheme which is widely popular in large-scale learning, with non-smooth extensions (Keskar et al., 2017; Skajaa, 2010). The technique requires only evaluation of objective function along with its gradient, which can be computed automatically using reverse mode automatic differentiation (as in Abadi et al. (2016)). To benefit from the available fast smooth implementations (Fei et al., 2014; Jones et al., 2001), we applied an infimal convolution

(see Section S.10.3 of the supplementary material) on the non-differentiable terms of the objective. Under the assumption that  $\mathbf{m} = \mathcal{O}(\sqrt{\mathbf{n}})$  (see Proposition 4.1), the complexity per L-BFGS-B iteration is  $\mathcal{O}(\mathbf{n}^2\sqrt{\mathbf{n}})$ . An experiment showing the impact of increasing  $\mathbf{m}$  on a synthetic dataset is provided in Fig. S.4.

**QR:** The efficiency of the non-crossing penalty is illustrated in Fig. 1 on the synthetic sine wave dataset described in Section 5 where  $\mathbf{n} = 40$  and  $\mathbf{m} = 20$  points have been generated. Many crossings are visible on the right plot, while they are almost not noticeable on the left plot, using the non-crossing penalty. Concerning our real-world examples, to study the efficiency of the proposed scheme in quantile regression the following experimental protocol was applied. Each dataset (Section 5) was splitted randomly into a training set (70%) and a test set (30%). We optimized the hyperparameters by minimizing a 5-folds cross validation with a Bayesian optimizer<sup>6</sup> (For further details see Section S.10.4). Once the hyperparameters were obtained, a new regressor was learned on the whole training set using the optimized hyperparameters. We report the value of the pinball loss and the crossing loss on the test set for three methods: our technique is called  $\infty$ -QR, we refer to Sangnier et al. (2016)’s approach as JQR, and independent learning (abbreviated as IND-QR) represents a further baseline.

We repeated 20 simulations (different random training-test splits); the results are also compared using a Mann-Whitney-Wilcoxon test. A summary is provided in Table 1. Notice that while JQR is tailored to predict finite many quantiles, our  $\infty$ -QR method estimates the *whole quantile function* hence solves a more challenging task. Despite the more difficult problem solved, as Table 1 suggest that the performance in terms of pinball loss of  $\infty$ -QR is comparable to that of the state-of-the-art JQR on all the twenty studied benchmarks, except for the ‘crabs’ and ‘cpus’ datasets (p.-val.  $< 0.25\%$ ). In addition, when considering the non-crossing penalty one can observe that  $\infty$ -QR outperforms the IND-QR baseline on eleven datasets (p.-val.  $< 0.25\%$ ) and JQR on two datasets. This illustrates the efficiency of the constraint based on the continuum scheme.

**DLSE:** To assess the quality of the estimated model by  $\infty$ -OCSVM, we illustrate the  $\theta$ -property

<sup>4</sup>The QMC approximation may involve the Sobol sequence with discrepancy  $m^{-1} \log(m)^s$  ( $s = \dim(\Theta)$ ).

<sup>5</sup>The code is available at <https://bitbucket.org/RomainBrault/itl>.

<sup>6</sup>We used a Gaussian Process model and minimized the Expected improvement. The optimizer was initialized using 27 samples from a Sobol sequence and ran for 50 iterations.

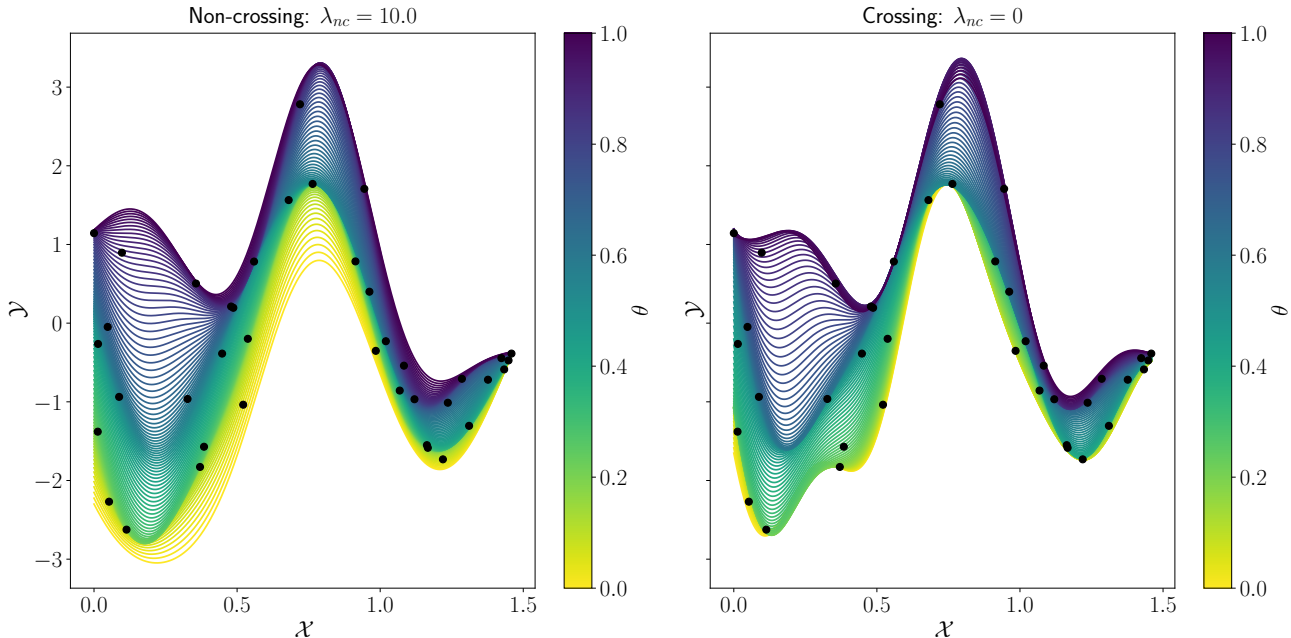


Figure 1: Impact of crossing penalty on toy data. Left plot: strong non-crossing penalty ( $\lambda_{nc} = 10$ ). Right plot: no non-crossing penalty ( $\lambda_{nc} = 0$ ). The plots show 100 quantiles of the continuum learned, linearly spaced between 0 (blue) and 1 (red). Notice that the non-crossing penalty does not provide crossings to occur in the regions where there is no points to enforce the penalty (e. g.  $x \in [0.13, 0.35]$ ). This phenomenon is alleviated by the regularity of the model.

(Schölkopf et al., 2000): the proportion of inliers has to be approximately  $1 - \theta$  ( $\forall \theta \in (0, 1)$ ). For the studied datasets (Wilt, Spambase) we used the raw inputs without applying any preprocessing. Our input kernel was the exponentiated  $\chi^2$  kernel  $k_X(x, z) := \exp\left(-\gamma_X \sum_{k=1}^d (x_k - z_k)^2 / (x_k + z_k)\right)$  with bandwidth  $\gamma_X = 0.25$ . A Gauss-Legendre quadrature rule provided the integral approximation in Eq. (8), with  $m = 100$  samples. We chose the Gaussian kernel for  $k_\Theta$ ; its bandwidth parameter  $\gamma_\Theta$  was the 0.2-quantile of the pairwise Euclidean distances between the  $\theta_j$ 's obtained via the quadrature rule. The margin (bias) kernel was  $k_b = k_\Theta$ . As it can be seen in Fig. 2, the  $\theta$ -property holds for the estimate which illustrates the efficiency of the proposed continuum approach for density level-set estimation.

## 6 Conclusion

In this work we proposed Infinite Task Learning, a novel nonparametric framework aiming at jointly solving parametrized tasks for a continuum of hyperparameters. We provided excess risk guarantees for the studied ITL scheme, and demonstrated its practical efficiency and flexibility in various tasks including cost-sensitive classification, quantile regression and density level set estimation.

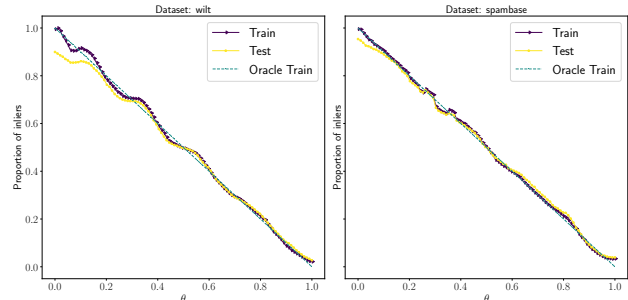


Figure 2: Density Level-Set Estimation: the  $\theta$ -property is approximately satisfied.

## Acknowledgments

The authors thank Arthur Tenenhaus for some insightful discussions. This work was supported by the Labex *DigiCosme* (project ANR-11-LABEX-0045-DIGICOSME) and the industrial chair *Machine Learning for Big Data* at Télécom ParisTech.

## References

Abadi, M. et al. (2016). “Tensorflow: Large-scale machine learning on heterogeneous distributed systems.” In: *USENIX Symposium on Operating Sys-*



- tems Design and Implementation (OSDI), pp. 265–283 (cit. on p. 7).
- Álvarez, M. A., L. Rosasco, and N. D. Lawrence (2012). “Kernels for vector-valued functions: a review.” In: *Foundations and Trends in Machine Learning* 4.3, pp. 195–266 (cit. on p. 3).
- Audiffren, J. and H. Kadri (2013). “Stability of Multi-Task Kernel Regression Algorithms.” In: *Asian Conference on Machine Learning (ACML)*. Vol. 29. PMLR, pp. 1–16 (cit. on p. 6).
- Bach, F. et al. (2012). “Optimization with sparsity-inducing penalties.” In: *Foundations and Trends in Machine Learning* 4.1, pp. 1–106 (cit. on p. 1).
- Bousquet, O. and A. Elisseeff (2002). “Stability and generalization.” In: *Journal of Machine Learning Research* 2, pp. 499–526 (cit. on p. 6).
- Carmeli, C. et al. (2010). “Vector valued reproducing kernel Hilbert spaces and universality.” In: *Analysis and Applications* 8 (1), pp. 19–61 (cit. on p. 2).
- Ciliberto, C. et al. (2017). “Consistent multitask learning with nonlinear output relations.” In: *Advances in Neural Information Processing Systems (NIPS)*, pp. 1986–1996 (cit. on p. 3).
- Evgeniou, T., C. A. Micchelli, and M. Pontil (2005). “Learning Multiple Tasks with kernel methods.” In: *JMLR* 6, pp. 615–637 (cit. on p. 3).
- Evgeniou, T. and M. Pontil (2004). “Regularized multi-task learning.” In: *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, pp. 109–117 (cit. on p. 1).
- Fei, Y. et al. (2014). “Parallel L-BFGS-B algorithm on GPU.” In: *Computers & Graphics* 40, pp. 1–9 (cit. on p. 7).
- Glazer, A., M. Lindenbaum, and S. Markovitch (2013). “q-OCSVM: A q-quantile estimator for high-dimensional distributions.” In: *Advances in Neural Information Processing Systems (NIPS)*, pp. 503–511 (cit. on pp. 1–3, 6).
- Hastie, T. et al. (2004). “The entire regularization path for the support vector machine.” In: *Journal of Machine Learning Research* 5, pp. 1391–1415 (cit. on p. 1).
- Jones, E., T. Oliphant, P. Peterson, et al. (2001). *SciPy: Open source scientific tools for Python* (cit. on p. 7).
- Kadri, H. et al. (2016). “Operator-valued Kernels for Learning from Functional Response Data.” In: *Journal of Machine Learning Research* 17, pp. 1–54 (cit. on pp. 6, 7).
- Keskar, N. and A. Wächter (2017). “A limited-memory quasi-Newton algorithm for bound-constrained non-smooth optimization.” In: *Optimization Methods and Software*, pp. 1–22 (cit. on p. 7).
- Koenker, R. and G. Bassett Jr (1978). “Regression quantiles.” In: *Econometrica: journal of the Econometric Society*, pp. 33–50 (cit. on pp. 1, 2).
- Pedrick, G. (1957). “Theory of reproducing kernels for Hilbert spaces of vector-valued functions.” PhD thesis. University of Kansas (cit. on p. 1).
- Sangnier, M., O. Fercoq, and F. d’Alché-Buc (2016). “Joint quantile regression in vector-valued RKHSs.” In: *Advances in Neural Information Processing Systems (NIPS)*, pp. 3693–3701 (cit. on pp. 2, 3, 5–7).
- Schölkopf, B. et al. (2000). “New support vector algorithms.” In: *Neural computation* 12.5, pp. 1207–1245 (cit. on pp. 1, 3, 8).
- Skajaa, A. (2010). “Limited memory BFGS for non-smooth optimization.” In: *Master’s thesis* (cit. on p. 7).
- Stein, O. (2012). “How to solve a semi-infinite optimization problem.” In: *European Journal of Operational Research* 223.2, pp. 312–320 (cit. on p. 1).
- Takeuchi, I., T. Hongo, et al. (2013). “Parametric task learning.” In: *Advances in Neural Information Processing Systems (NIPS)*, pp. 1358–1366 (cit. on pp. 1, 3).
- Takeuchi, I., Q. V. Le, et al. (2006). “Nonparametric quantile estimation.” In: *Journal of Machine Learning Research* 7, pp. 1231–1264 (cit. on p. 1).
- Tibshirani, R. (1996). “Regression shrinkage and selection via the Lasso.” In: *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 267–288 (cit. on p. 1).
- Tikhonov, A. N. and V. Y. Arsenin (1977). *Solution of Ill-posed Problems*. Winston & Sons (cit. on p. 1).
- Vapnik, V., S. E. Golowich, and A. J. Smola (1997). “Support vector method for function approximation, regression estimation and signal processing.” In: *Advances in Neural Information Processing Systems (NIPS)*, pp. 281–287 (cit. on p. 1).
- Vert, R. and J.-P. Vert (2006). “Consistency and convergence rates of one-class SVMs and related algorithms.” In: *Journal of Machine Learning Research* 7, pp. 817–854 (cit. on p. 3).
- Wen, U.-P. and S.-T. Hsu (1991). “Linear bi-level programming problems a review.” In: *Journal of the Operational Research Society* 42.2, pp. 125–133 (cit. on p. 1).
- Zadrozny, B. and C. Elkan (2001). “Learning and making decisions when costs and probabilities are both unknown.” In: *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 204–213 (cit. on p. 1).
- Zhu, C. et al. (1997). “Algorithm 778: L-BFGS-B: Fortran subroutines for large-scale bound-constrained optimization.” In: *ACM Transactions on Mathematical Software (TOMS)* 23.4, pp. 550–560 (cit. on p. 7).