

---

# Kernel Bayesian Inference with Posterior Regularization

---

Yang Song<sup>†</sup>, Jun Zhu<sup>‡\*</sup>, Yong Ren<sup>‡</sup>

<sup>†</sup> Dept. of Physics, Tsinghua University, Beijing, China

<sup>‡</sup> Dept. of Comp. Sci. & Tech., TNList Lab; Center for Bio-Inspired Computing Research  
State Key Lab for Intell. Tech. & Systems, Tsinghua University, Beijing, China  
yangsong@cs.stanford.edu; {dcszj@, renyong15@mails}.tsinghua.edu.cn

## Abstract

We propose a vector-valued regression problem whose solution is equivalent to the reproducing kernel Hilbert space (RKHS) embedding of the Bayesian posterior distribution. This equivalence provides a new understanding of kernel Bayesian inference. Moreover, the optimization problem induces a new regularization for the posterior embedding estimator, which is faster and has comparable performance to the squared regularization in kernel Bayes' rule. This regularization coincides with a former thresholding approach used in kernel POMDPs whose consistency remains to be established. Our theoretical work solves this open problem and provides consistency analysis in regression settings. Based on our optimizational formulation, we propose a flexible Bayesian posterior regularization framework which for the first time enables us to put regularization at the distribution level. We apply this method to nonparametric state-space filtering tasks with extremely nonlinear dynamics and show performance gains over all other baselines.

## 1 Introduction

Kernel methods have long been effective in generalizing linear statistical approaches to nonlinear cases by embedding a sample to the reproducing kernel Hilbert space (RKHS) [1]. In recent years, the idea has been generalized to embedding probability distributions [2, 3]. Such embeddings of probability measures are usually called *kernel embeddings* (a.k.a. *kernel means*). Moreover, [4, 5, 6] show that statistical operations of distributions can be realized in RKHS by manipulating kernel embeddings via linear operators. This approach has been applied to various statistical inference and learning problems, including training hidden Markov models (HMM) [7], belief propagation (BP) in tree graphical models [8], planning Markov decision processes (MDP) [9] and partially observed Markov decision processes (POMDP) [10].

One of the key workhorses in the above applications is the *kernel Bayes' rule* [5], which establishes the relation among the RKHS representations of the priors, likelihood functions and posterior distributions. Despite empirical success, the characterization of kernel Bayes' rule remains largely incomplete. For example, it is unclear how the estimators of the posterior distribution embeddings relate to optimizers of some loss functions, though the vanilla Bayes' rule has a nice connection [11]. This makes generalizing the results especially difficult and hinders the intuitive understanding of kernel Bayes' rule.

To alleviate this weakness, we propose a vector-valued regression [12] problem whose optimizer is the posterior distribution embedding. This new formulation is inspired by the progress in two fields: 1) the alternative characterization of conditional embeddings as regressors [13], and 2) the

---

\*Corresponding author.

introduction of posterior regularized Bayesian inference (RegBayes) [14] based on an optimizational reformulation of the Bayes' rule.

We demonstrate the novelty of our formulation by providing a new understanding of kernel Bayesian inference, with theoretical, algorithmic and practical implications. On the theoretical side, we are able to prove the (weak) consistency of the estimator obtained by solving the vector-valued regression problem under reasonable assumptions. As a side product, our proof can be applied to a thresholding technique used in [10], whose consistency is left as an open problem. On the algorithmic side, we propose a new regularization technique, which is shown to run faster and has comparable accuracy to squared regularization used in the original kernel Bayes' rule [5]. Similar in spirit to RegBayes, we are also able to derive an extended version of the embeddings by directly imposing regularization on the posterior distributions. We call this new framework kRegBayes. Thanks to RKHS embeddings of distributions, this is the first time, to the best of our knowledge, people can do posterior regularization without invoking linear functionals (such as moments) of the random variables. On the practical side, we demonstrate the efficacy of our methods on both simple and complicated synthetic state-space filtering datasets.

Same to other algorithms based on kernel embeddings, our kernel regularized Bayesian inference framework is nonparametric and general. The algorithm is nonparametric, because the priors, posterior distributions and likelihood functions are all characterized by weighted sums of data samples. Hence it does not need the explicit mechanism such as differential equations of a robot arm in filtering tasks. It is general in terms of being applicable to a broad variety of domains as long as the kernels can be defined, such as strings, orthonormal matrices, permutations and graphs.

## 2 Preliminaries

### 2.1 Kernel embeddings

Let  $(\mathcal{X}, \mathcal{B}_\mathcal{X})$  be a measurable space of random variables,  $p_X$  be the associated probability measure and  $\mathcal{H}_\mathcal{X}$  be a RKHS with kernel  $k(\cdot, \cdot)$ . We define the *kernel embedding* of  $p_X$  to be  $\mu_X = \mathbb{E}_{p_X}[\phi(X)] \in \mathcal{H}_\mathcal{X}$ , where  $\phi(X) = k(X, \cdot)$  is the feature map. Such a vector-valued expectation always exists if the kernel is bounded, namely  $\sup_x k_\mathcal{X}(x, x) < \infty$ .

The concept of kernel embeddings has several important statistical merits. Inasmuch as the reproducing property, the expectation of  $f \in \mathcal{H}$  w.r.t.  $p_X$  can be easily computed as  $\mathbb{E}_{p_X}[f(X)] = \mathbb{E}_{p_X}[\langle f, \phi(X) \rangle] = \langle f, \mu_X \rangle$ . There exists *universal kernels* [15] whose corresponding RKHS  $\mathcal{H}$  is dense in  $\mathcal{C}_\mathcal{X}$  in terms of sup norm. This means  $\mathcal{H}$  contains a rich range of functions  $f$  and their expectations can be computed by inner products without invoking usually intractable integrals. In addition, the inner product structure of the embedding space  $\mathcal{H}$  provides a natural way to measure the differences of distributions through norms.

In much the same way we can define kernel embeddings of linear operators. Let  $(\mathcal{X}, \mathcal{B}_\mathcal{X})$  and  $(\mathcal{Y}, \mathcal{B}_\mathcal{Y})$  be two measurable spaces,  $\phi(x)$  and  $\psi(y)$  be the measurable feature maps of corresponding RKHS  $\mathcal{H}_\mathcal{X}$  and  $\mathcal{H}_\mathcal{Y}$  with bounded kernels, and  $p$  denote the joint distribution of a random variable  $(X, Y)$  on  $\mathcal{X} \times \mathcal{Y}$  with product measures. The *covariance operator*  $\mathcal{C}_{XY}$  is defined as  $\mathcal{C}_{XY} = \mathbb{E}_p[\phi(X) \otimes \psi(Y)]$ , where  $\otimes$  denotes the tensor product. Note that it is possible to identify  $\mathcal{C}_{XY}$  with  $\mu_{(XY)}$  in  $\mathcal{H}_\mathcal{X} \otimes \mathcal{H}_\mathcal{Y}$  with the kernel function  $k((x_1, y_1), (x_2, y_2)) = k_\mathcal{X}(x_1, x_2)k_\mathcal{Y}(y_1, y_2)$  [16]. There is an important relation between kernel embeddings of distributions and covariance operators, which is fundamental for the sequel:

**Theorem 1** ([4, 5]). *Let  $\mu_X, \mu_Y$  be the kernel embeddings of  $p_X$  and  $p_Y$  respectively. If  $\mathcal{C}_{XX}$  is injective,  $\mu_X \in \mathcal{R}(\mathcal{C}_{XX})$  and  $\mathbb{E}[g(Y) | X = \cdot] \in \mathcal{H}_\mathcal{X}$  for all  $g \in \mathcal{H}_\mathcal{Y}$ , then*

$$\mu_Y = \mathcal{C}_{YX} \mathcal{C}_{XX}^{-1} \mu_X. \quad (1)$$

*In addition,  $\mu_{Y|X=x} = \mathbb{E}[\psi(Y)|X=x] = \mathcal{C}_{YX} \mathcal{C}_{XX}^{-1} \phi(x)$ .*

On the implementation side, we need to estimate these kernel embeddings via samples. An intuitive estimator for the embedding  $\mu_X$  is  $\hat{\mu}_X = \frac{1}{N} \sum_{i=1}^N \phi(x_i)$ , where  $\{x_i\}_{i=1}^N$  is a sample from  $p_X$ . Similarly, the covariance operators can also be estimated by  $\hat{\mathcal{C}}_{XY} = \frac{1}{N} \sum_{i=1}^N \phi(x_i) \otimes \psi(y_i)$ . Both operators are shown to converge in the RKHS norm at a rate of  $O_p(N^{-\frac{1}{2}})$  [4].

## 2.2 Kernel Bayes' rule

Let  $\pi(Y)$  be the prior distribution of a random variable  $Y$ ,  $p(X = x | Y)$  be the likelihood,  $p^\pi(Y | X = x)$  be the posterior distribution given  $\pi(Y)$  and observation  $x$ , and  $p^\pi(X, Y)$  be the joint distribution incorporating  $\pi(Y)$  and  $p(X | Y)$ . Kernel Bayesian inference aims to obtain the posterior embedding  $\mu_Y^\pi(X = x)$  given a prior embedding  $\pi_Y$  and a covariance operator  $\mathcal{C}_{XY}$ . By Bayes' rule,  $p^\pi(Y | X = x) \propto \pi(Y)p(X = x | Y)$ . We assume that there exists a joint distribution  $p$  on  $\mathcal{X} \times \mathcal{Y}$  whose conditional distribution matches  $p(X | Y)$  and let  $\mathcal{C}_{XY}$  be its covariance operator. Note that we do not require  $p = p^\pi$  hence  $p$  can be any convenient distribution.

According to Thm. 1,  $\mu_Y^\pi(X = x) = \mathcal{C}_{YX}^\pi \mathcal{C}_{XX}^\pi{}^{-1} \phi(x)$ , where  $\mathcal{C}_{YX}^\pi$  corresponds to the joint distribution  $p^\pi$  and  $\mathcal{C}_{XX}^\pi$  to the marginal probability of  $p^\pi$  on  $X$ . Recall that  $\mathcal{C}_{YX}^\pi$  can be identified with  $\mu_{(YX)}$  in  $\mathcal{H}_Y \otimes \mathcal{H}_X$ , we can apply Thm. 1 to obtain  $\mu_{(YX)} = \mathcal{C}_{(YX)Y} \mathcal{C}_{YY}^{-1} \pi_Y$ , where  $\mathcal{C}_{(YX)Y} := \mathbb{E}[\psi(Y) \otimes \phi(X) \otimes \psi(Y)]$ . Similarly,  $\mathcal{C}_{XX}^\pi$  can be represented as  $\mu_{(XX)} = \mathcal{C}_{(XX)Y} \mathcal{C}_{YY}^{-1} \pi_Y$ . This way of computing posterior embeddings is called the *kernel Bayes' rule* [5].

Given estimators of the prior embedding  $\hat{\pi}_Y = \sum_{i=1}^m \tilde{\alpha}_i \psi(y_i)$  and the covariance operator  $\hat{\mathcal{C}}_{YX}$ , The posterior embedding can be obtained via  $\hat{\mu}_Y^\pi(X = x) = \hat{\mathcal{C}}_{YX}^\pi ([\hat{\mathcal{C}}_{XX}^\pi]^2 + \lambda I)^{-1} \hat{\mathcal{C}}_{XX}^\pi \phi(x)$ , where squared regularization is added to the inversion. Note that the regularization for  $\hat{\mu}_Y^\pi(X = x)$  is not unique. A thresholding alternative is proposed in [10] without establishing the consistency. We will discuss this thresholding regularization in a different perspective and give consistency results in the sequel.

## 2.3 Regularized Bayesian inference

Regularized Bayesian inference (RegBayes [14]) is based on a variational formulation of the Bayes' rule [11]. The posterior distribution can be viewed as the solution of  $\min_{p(Y|X=x)} \text{KL}(p(Y|X = x) \| \pi(Y)) - \int \log p(X = x|Y) dp(Y|X = x)$ , subjected to  $p(Y|X = x) \in \mathcal{P}_{\text{prob}}$ , where  $\mathcal{P}_{\text{prob}}$  is the set of valid probability measures. RegBayes combines this formulation and posterior regularization [17] in the following way

$$\begin{aligned} \min_{p(Y|X=x), \xi} \quad & \text{KL}(p(Y|X = x) \| \pi(Y)) - \int \log p(X = x|Y) dp(Y|X = x) + U(\xi) \\ \text{s.t.} \quad & p(Y|X = x) \in \mathcal{P}_{\text{prob}}(\xi), \end{aligned}$$

where  $\mathcal{P}_{\text{prob}}(\xi)$  is a subset depending on  $\xi$  and  $U(\xi)$  is a loss function. Such a formulation makes it possible to regularize Bayesian posterior distributions, smoothing the gap between Bayesian generative models and discriminative models. Related applications include max-margin topic models [18] and infinite latent SVMs [14].

Despite the flexibility of RegBayes, regularization on the posterior distributions is practically imposed indirectly via expectations of a function. We shall see soon in the sequel that our new framework of kernel Regularized Bayesian inference can control the posterior distribution in a direct way.

## 2.4 Vector-valued regression

The main task for vector-valued regression [12] is to minimize the following objective

$$E(f) := \sum_{i=1}^n \|y_j - f(x_j)\|_{\mathcal{H}_Y}^2 + \lambda \|f\|_{\mathcal{H}_K}^2,$$

where  $y_j \in \mathcal{H}_Y$ ,  $f : \mathcal{X} \rightarrow \mathcal{H}_Y$ . Note that  $f$  is a function with RKHS values and we assume that  $f$  belongs to a *vector-valued* RKHS  $\mathcal{H}_K$ . In vector-valued RKHS, the kernel function  $k$  is generalized to linear operators  $\mathcal{L}(\mathcal{H}_Y) \ni K(x_1, x_2) : \mathcal{H}_Y \rightarrow \mathcal{H}_Y$ , such that  $K(x_1, x_2)y := (K_{x_2}y)(x_1)$  for every  $x_1, x_2 \in \mathcal{X}$  and  $y \in \mathcal{H}_Y$ , where  $K_{x_2}y \in \mathcal{H}_K$ . The reproducing property is generalized to  $\langle y, f(x) \rangle_{\mathcal{H}_Y} = \langle K_x y, f \rangle_{\mathcal{H}_K}$  for every  $y \in \mathcal{H}_Y$ ,  $f \in \mathcal{H}_K$  and  $x \in \mathcal{X}$ . In addition, [12] shows that the representer theorem still holds for vector-valued RKHS.

## 3 Kernel Bayesian inference as a regression problem

One of the unique merits of the posterior embedding  $\mu_Y^\pi(X = x)$  is that expectations w.r.t. posterior distributions can be computed via inner products, *i.e.*,  $\langle h, \mu_Y^\pi(X = x) \rangle = \mathbb{E}_{p^\pi(Y|X=x)}[h(Y)]$  for all

$h \in \mathcal{H}_Y$ . Since  $\mu_Y^\pi(X = x) \in \mathcal{H}_Y$ ,  $\mu_Y^\pi$  can be viewed as an element of a vector-valued RKHS  $\mathcal{H}_K$  containing functions  $f : \mathcal{X} \rightarrow \mathcal{H}_Y$ .

A natural optimization objective [13] thus follows from the above observations

$$\mathcal{E}[\mu] := \sup_{\|h\|_Y \leq 1} \mathbb{E}_X [\langle \mathbb{E}_Y[h(Y)|X] - \langle h, \mu(X) \rangle_{\mathcal{H}_Y} \rangle^2], \quad (2)$$

where  $\mathbb{E}_X[\cdot]$  denotes the expectation w.r.t.  $p^\pi(X)$  and  $\mathbb{E}_Y[\cdot|X]$  denotes the expectation w.r.t. the Bayesian posterior distribution, i.e.,  $p^\pi(Y | X) \propto \pi(Y)p(X | Y)$ . Clearly,  $\mu_Y^\pi = \arg \inf_{\mu} \mathcal{E}[\mu]$ . Following [13], we introduce an upper bound  $\mathcal{E}_s$  for  $\mathcal{E}$  by applying Jensen's and Cauchy-Schwarz's inequalities consecutively

$$\mathcal{E}_s[\mu] := \mathbb{E}_{(X,Y)} [\|\psi(Y) - \mu(X)\|_{\mathcal{H}_Y}^2], \quad (3)$$

where  $(X, Y)$  is the random variable on  $\mathcal{X} \times \mathcal{Y}$  with the joint distribution  $p^\pi(X, Y) = \pi(Y)p(X | Y)$ .

The first step to make this optimizational framework practical is to find finite sample estimators of  $\mathcal{E}_s[\mu]$ . We will show how to do this in the following section.

### 3.1 A consistent estimator of $\mathcal{E}_s[\mu]$

Unlike the conditional embeddings in [13], we do not have i.i.d. samples from the joint distribution  $p^\pi(X, Y)$ , as the priors and likelihood functions are represented with samples from different distributions. We will eliminate this problem using a kernel trick, which is one of our main innovations in this paper.

The idea is to use the inner product property of a kernel embedding  $\mu_{(X,Y)}$  to represent the expectation  $\mathbb{E}_{(X,Y)} [\|\psi(Y) - \mu(X)\|_{\mathcal{H}_Y}^2]$  and then use finite sample estimators of  $\mu_{(X,Y)}$  to estimate  $\mathcal{E}_s[\mu]$ . Recall that we can identify  $\mathcal{C}_{XY} := \mathbb{E}_{XY}[\phi(X) \otimes \psi(Y)]$  with  $\mu_{(X,Y)}$  in a product space  $\mathcal{H}_X \otimes \mathcal{H}_Y$  with a product kernel  $k_X k_Y$  on  $\mathcal{X} \times \mathcal{Y}$  [16]. Let  $f(x, y) = \|\psi(y) - \mu(x)\|_{\mathcal{H}_Y}^2$  and assume that  $f \in \mathcal{H}_X \otimes \mathcal{H}_Y$ . The optimization objective  $\mathcal{E}_s[\mu]$  can be written as

$$\mathcal{E}_s[\mu] = \mathbb{E}_{(X,Y)} [\|\psi(Y) - \mu(X)\|_{\mathcal{H}_Y}^2] = \langle f, \mu_{(X,Y)} \rangle_{\mathcal{H}_X \otimes \mathcal{H}_Y}. \quad (4)$$

From Thm. 1, we assert that  $\mu_{(X,Y)} = \mathcal{C}_{(X,Y)Y} \mathcal{C}_{YY}^{-1} \pi_Y$  and a natural estimator follows to be  $\hat{\mu}_{(X,Y)} = \hat{\mathcal{C}}_{(X,Y)Y} (\hat{\mathcal{C}}_{YY} + \lambda I)^{-1} \hat{\pi}_Y$ . As a result,  $\hat{\mathcal{E}}_s[\mu] := \langle \hat{\mu}_{(X,Y)}, f \rangle_{\mathcal{H}_X \otimes \mathcal{H}_Y}$  and we introduce the following proposition to write  $\hat{\mathcal{E}}_s$  in terms of Gram matrices.

**Proposition 1** (Proof in Appendix). *Suppose  $(X, Y)$  is a random variable in  $\mathcal{X} \times \mathcal{Y}$ , where the prior for  $Y$  is  $\pi(Y)$  and the likelihood is  $p(X | Y)$ . Let  $\mathcal{H}_X$  be a RKHS with kernel  $k_X$  and feature map  $\phi(x)$ ,  $\mathcal{H}_Y$  be a RKHS with kernel  $k_Y$  and feature map  $\psi(y)$ ,  $\phi(x, y)$  be the feature map of  $\mathcal{H}_X \otimes \mathcal{H}_Y$ ,  $\hat{\pi}_Y = \sum_{i=1}^l \tilde{\alpha}_i \psi(\tilde{y}_i)$  be a consistent estimator of  $\pi_Y$  and  $\{(x_i, y_i)\}_{i=1}^n$  be a sample representing  $p(X | Y)$ . Under the assumption that  $f(x, y) = \|\psi(y) - \mu(x)\|_{\mathcal{H}_Y}^2 \in \mathcal{H}_X \otimes \mathcal{H}_Y$ , we have*

$$\hat{\mathcal{E}}_s[\mu] = \sum_{i=1}^n \beta_i \|\psi(y_i) - \mu(x_i)\|_{\mathcal{H}_Y}^2, \quad (5)$$

where  $\beta = (\beta_1, \dots, \beta_n)^\top$  is given by  $\beta = (G_Y + n\lambda I)^{-1} \tilde{G}_Y \tilde{\alpha}$ , where  $(G_Y)_{ij} = k_Y(y_i, y_j)$ ,  $(\tilde{G}_Y)_{ij} = k_Y(y_i, \tilde{y}_j)$ , and  $\tilde{\alpha} = (\tilde{\alpha}_1, \dots, \tilde{\alpha}_l)^\top$ .

The consistency of  $\hat{\mathcal{E}}_s[\mu]$  is a direct consequence of the following theorem adapted from [5], since the Cauchy-Schwarz inequality ensures  $|\langle \mu_{(X,Y)}, f \rangle - \langle \hat{\mu}_{(X,Y)}, f \rangle| \leq \|\mu_{(X,Y)} - \hat{\mu}_{(X,Y)}\| \|f\|$ .

**Theorem 2** (Adapted from [5], Theorem 8). *Assume that  $\mathcal{C}_{YY}$  is injective,  $\hat{\pi}_Y$  is a consistent estimator of  $\pi_Y$  in  $\mathcal{H}_Y$  norm, and that  $\mathbb{E}[k(\tilde{X}, Y), (\tilde{X}, \tilde{Y}) | Y = y, \tilde{Y} = \tilde{y}]$  is included in  $\mathcal{H}_X \otimes \mathcal{H}_Y$  as a function of  $(y, \tilde{y})$ , where  $(\tilde{X}, \tilde{Y})$  is an independent copy of  $(X, Y)$ . Then, if the regularization coefficient  $\lambda_n$  decays to 0 sufficiently slowly,*

$$\left\| \hat{\mathcal{C}}_{(X,Y)Y} (\hat{\mathcal{C}}_{YY} + \lambda_n I)^{-1} \hat{\pi}_Y - \mu_{(X,Y)} \right\|_{\mathcal{H}_X \otimes \mathcal{H}_Y} \rightarrow 0 \quad (6)$$

in probability as  $n \rightarrow \infty$ .

Although  $\widehat{\mathcal{E}}_s[\mu]$  is a consistent estimator of  $\mathcal{E}_s[\mu]$ , it does not necessarily have minima, since the coefficients  $\beta_i$  can be negative. One of our main contributions in this paper is the discovery that we can ignore data points  $(x_i, y_i)$  with a negative  $\beta_i$ , *i.e.*, replacing  $\beta_i$  with  $\beta_i^+ := \max(0, \beta_i)$  in  $\widehat{\mathcal{E}}_s[\mu]$ . We will give explanations and theoretical justifications in the next section.

### 3.2 The thresholding regularization

We show in the following theorem that  $\widehat{\mathcal{E}}_s^+[\mu] := \sum_{i=1}^n \beta_i^+ \|\psi(y_i) - \mu(x_i)\|^2$  converges to  $\mathcal{E}_s[\mu]$  in probability in discrete situations. The trick of replacing  $\beta_i$  with  $\beta_i^+$  is named *thresholding regularization*.

**Theorem 3** (Proof in Appendix). *Assume that  $\mathcal{X}$  is compact and  $|\mathcal{Y}| < \infty$ ,  $k$  is a strictly positive definite continuous kernel with  $\sup_{(x,y)} k((x,y), (x,y)) < \kappa$  and  $f(x,y) = \|\psi(y) - \mu(x)\|_{\mathcal{H}_Y}^2 \in \mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_Y$ . With the conditions in Thm. 2, we assert that  $\widehat{\mu}_{(\mathcal{X},\mathcal{Y})}^+$  is a consistent estimator of  $\mu_{(\mathcal{X},\mathcal{Y})}$  and  $|\widehat{\mathcal{E}}_s^+[\mu] - \mathcal{E}_s[\mu]| \rightarrow 0$  in probability as  $n \rightarrow \infty$ .*

In the context of partially observed Markov decision processes (POMDPs) [10], a similar thresholding approach, combined with normalization, was proposed to make the Bellman operator isotonic and contractive. However, the authors left the consistency of that approach as an open problem. The justification of normalization has been provided in [13], Lemma 2.2 under the finite space assumption. A slight modification of our proof of Thm. 3 (change the probability space from  $\mathcal{X} \times \mathcal{Y}$  to  $\mathcal{X}$ ) can complete the other half as a side product, under the same assumptions.

Compared to the original squared regularization used in [5], thresholding regularization is more computational efficient because 1) it does not need to multiply the Gram matrix twice, and 2) it does not need to take into consideration those data points with negative  $\beta_i$ 's. In many cases a large portion of  $\{\beta_i\}_{i=1}^n$  is negative but the sum of their absolute values is small. The finite space assumption in Thm. 3 may also be weakened, but it requires deeper theoretical analyses.

### 3.3 Minimizing $\widehat{\mathcal{E}}_s^+[\mu]$

Following the standard steps of solving a RKHS regression problem, we add a Tikhonov regularization term to  $\widehat{\mathcal{E}}_s^+[\mu]$  to provide a well-proposed problem,

$$\widehat{\mathcal{E}}_{\lambda,n}[\mu] = \sum_{i=1}^n \beta_i^+ \|\psi(y_i) - \mu(x_i)\|_{\mathcal{H}_Y}^2 + \lambda \|\mu\|_{\mathcal{H}_K}^2. \quad (7)$$

Let  $\widehat{\mu}_{\lambda,n} = \arg \min_{\mu} \widehat{\mathcal{E}}_{\lambda,n}[\mu]$ . Note that  $\widehat{\mathcal{E}}_{\lambda,n}[\mu]$  is a vector-valued regression problem, and the representer theorems in vector-valued RKHS apply here. We summarize the matrix expression of  $\widehat{\mu}_{\lambda,n}$  in the following proposition.

**Proposition 2** (Proof in Appendix). *Without loss of generality, we assume that  $\beta_i^+ \neq 0$  for all  $1 \leq i \leq n$ . Let  $\mu \in \mathcal{H}_K$  and choose the kernel of  $\mathcal{H}_K$  to be  $K(x_i, x_j) = k_{\mathcal{X}}(x_i, x_j)\mathcal{I}$ , where  $\mathcal{I} : \mathcal{H}_K \rightarrow \mathcal{H}_K$  is an identity map. Then*

$$\widehat{\mu}_{\lambda,n}(x) = \Psi(K_X + \lambda_n \Lambda^+)^{-1} K_{:x}, \quad (8)$$

where  $\Psi = (\psi(y_1), \dots, \psi(y_n))$ ,  $(K_X)_{ij} = k_{\mathcal{X}}(x_i, x_j)$ ,  $\Lambda^+ = \text{diag}(1/\beta_1^+, \dots, 1/\beta_n^+)$ ,  $K_{:x} = (k_{\mathcal{X}}(x, x_1), \dots, k_{\mathcal{X}}(x, x_n))^{\top}$  and  $\lambda_n$  is a positive regularization constant.

### 3.4 Theoretical justifications for $\widehat{\mu}_{\lambda,n}$

In this section, we provide theoretical explanations for using  $\widehat{\mu}_{\lambda,n}$  as an estimator of the posterior embedding under specific assumptions. Let  $\mu^* = \arg \min_{\mu} \mathcal{E}[\mu]$ ,  $\mu' = \arg \min_{\mu} \mathcal{E}_s[\mu]$ , and recall that  $\widehat{\mu}_{\lambda,n} = \arg \min_{\mu} \widehat{\mathcal{E}}_{\lambda,n}[\mu]$ . We first show the relations between  $\mu^*$  and  $\mu'$  and then discuss the relations between  $\widehat{\mu}_{\lambda,n}$  and  $\mu'$ .

The forms of  $\mathcal{E}$  and  $\mathcal{E}_s$  are exactly the same for posterior kernel embeddings and conditional kernel embeddings. As a consequence, the following theorem in [13] still hold.

**Theorem 4** ([13]). *If there exists a  $\mu^* \in \mathcal{H}_K$  such that for any  $h \in \mathcal{H}_Y$ ,  $\mathbb{E}[h|X] = \langle h, \mu^*(X) \rangle_{\mathcal{H}_Y}$   $p_X$ -a.s., then  $\mu^*$  is the  $p_X$ -a.s. unique minimiser of both objectives:*

$$\mu^* = \arg \min_{\mu \in \mathcal{H}_K} \mathcal{E}[\mu] = \arg \min_{\mu \in \mathcal{H}_K} \mathcal{E}_s[\mu].$$

This theorem shows that if the vector-valued RKHS  $\mathcal{H}_K$  is rich enough to contain  $\mu_{Y|X=x}^\pi$ , both  $\mathcal{E}$  and  $\mathcal{E}_s$  can lead us to the correct embedding. In this case, it is reasonable to use  $\mu'$  instead of  $\mu^*$ . For the situation where  $\mu_{Y|X=x}^\pi \notin \mathcal{H}_K$ , we refer the readers to [13].

Unfortunately, we cannot obtain the relation between  $\widehat{\mu}_{\lambda,n}$  and  $\mu'$  by referring to [19], as in [13]. The main difficulty here is that  $\{(x_i, y_i)\}_{i=1}^n$  is not an i.i.d. sample from  $p^\pi(X, Y) = \pi(Y)p(X | Y)$  and the estimator  $\widehat{\mathcal{E}}_s^+[\mu]$  does not use i.i.d. samples to estimate expectations. Therefore the concentration inequality ([19], Prop. 2) used in the proofs of [19] cannot be applied.

To solve the problem, we propose Thm. 9 (in Appendix) which can lead to a consistency proof for  $\widehat{\mu}_{\lambda,n}$ . The relation between  $\widehat{\mu}_{\lambda,n}$  and  $\mu'$  can now be summarized in the following theorem.

**Theorem 5** (Proof in Appendix). *Assume Hypothesis 1 and Hypothesis 2 in [20] and our Assumption 1 (in the Appendix) hold. With the conditions in Thm. 3, we assert that if  $\lambda_n$  decreases to 0 sufficiently slowly,*

$$\mathcal{E}_s[\widehat{\mu}_{\lambda_n, n}] - \mathcal{E}_s[\mu'] \rightarrow 0 \quad (9)$$

in probability as  $n \rightarrow \infty$ .

## 4 Kernel Bayesian inference with posterior regularization

Based on our optimizational formulation of kernel Bayesian inference, we can add additional regularization terms to control the posterior embeddings. This technique gives us the possibility to incorporate rich side information from domain knowledge and to enforce supervisions on Bayesian inference. We call our framework of imposing posterior regularization *kRegBayes*.

As an example of the framework, we study the following optimization problem

$$\mathcal{L} := \underbrace{\sum_{i=1}^m \beta_i^+ \|\mu(x_i) - \psi(y_i)\|_{\mathcal{H}_Y}^2 + \lambda \|\mu\|_{\mathcal{H}_K}^2}_{\widehat{\mathcal{E}}_{\lambda, n}[\mu]} + \delta \underbrace{\sum_{i=m+1}^n \|\mu(x_i) - \psi(t_i)\|_{\mathcal{H}_Y}^2}_{\text{The regularization term}}, \quad (10)$$

where  $\{(x_i, y_i)\}_{i=1}^m$  is the sample used for representing the likelihood,  $\{(x_i, t_i)\}_{i=m+1}^n$  is the sample used for posterior regularization and  $\lambda, \delta$  are the regularization constants. Note that in RKHS embeddings,  $\psi(t)$  is identified as a point distribution at  $t$  [2]. Hence the regularization term in (10) encourages the posterior distributions  $p(Y | X = x_i)$  to be concentrated at  $t_i$ . More complicated regularization terms are also possible, such as  $\|\mu(x_i) - \sum_{i=1}^l \alpha_i \psi(t_i)\|_{\mathcal{H}_Y}$ .

Compared to vanilla RegBayes, our kernel counterpart has several obvious advantages. First, the difference between two distributions can be naturally measured by RKHS norms. This makes it possible to regularize the posterior distribution as a whole, rather than through expectations of discriminant functions. Second, the framework of kernel Bayesian inference is totally nonparametric, where the priors and likelihood functions are all represented by respective samples. We will further demonstrate the properties of kRegBayes through experiments in the next section.

Let  $\widehat{\mu}_{reg} = \arg \min_{\mu} \mathcal{L}$ . It is clear that solving  $\mathcal{L}$  is substantially the same as  $\widehat{\mathcal{E}}_{\lambda, n}[\mu]$  and we summarize it in the following proposition.

**Proposition 3.** *With the conditions in Prop. 2, we have*

$$\widehat{\mu}_{reg}(x) = \Psi(K_X + \lambda\Lambda^+)^{-1}K_{:x}, \quad (11)$$

where  $\Psi = (\psi(y_1), \dots, \psi(y_n))$ ,  $(K_X)_{ij} = k_{\mathcal{X}}(x_i, x_j)_{1 \leq i, j \leq n}$ ,  $\Lambda^+ = \text{diag}(1/\beta_1^+, \dots, 1/\beta_m^+, 1/\delta, \dots, 1/\delta)$ , and  $K_{:x} = (k_{\mathcal{X}}(x, x_1), \dots, k_{\mathcal{X}}(x, x_n))^T$ .

## 5 Experiments

In this section, we compare the results of kRegBayes and several other baselines for two state-space filtering tasks. The mechanism behind kernel filtering is stated in [5] and we provide a detailed introduction in Appendix, including all the formula used in implementation.

**Toy dynamics** This experiment is a twist of that used in [5]. We report the results of extended Kalman filter (EKF) [21] and unscented Kalman filter (UKF) [22], kernel Bayes’ rule (KBR) [5], kernel Bayesian learning with thresholding regularization (pKBR) and kRegBayes.

The data points  $\{(\theta_t, x_t, y_t)\}$  are generated from the dynamics

$$\theta_{t+1} = \theta_t + 0.4 + \xi_t \pmod{2\pi}, \quad \begin{pmatrix} x_{t+1} \\ y_{t+1} \end{pmatrix} = (1 + \sin(8\theta_{t+1})) \begin{pmatrix} \cos \theta_{t+1} \\ \sin \theta_{t+1} \end{pmatrix} + \zeta_t, \quad (12)$$

where  $\theta_t$  is the hidden state,  $(x_t, y_t)$  is the observation,  $\xi_t \sim \mathcal{N}(0, 0.04)$  and  $\zeta_t \sim \mathcal{N}(0, 0.04)$ . Note that this dynamics is nonlinear for both transition and observation functions. The observation model is an oscillation around the unit circle. There are 1000 training data and 200 validation/test data for each algorithm.

We suppose that EKF, UKF and kRegBayes know the true dynamics of the model and the first hidden state  $\theta_1$ . In this case, we use  $\tilde{\theta}_{t+1} = \theta_1 + 0.4t \pmod{2\pi}$  and  $(\tilde{x}_{t+1}, \tilde{y}_{t+1})^\top = (1 + \sin(8\tilde{\theta}_{t+1}))(\cos \tilde{\theta}_{t+1}, \sin \tilde{\theta}_{t+1})^\top$  as the supervision data point for the  $(t + 1)$ -th step. We follow [5] to set our parameters.

The results are summarized in Fig. 5. pKBR has lower errors compared to KBR, which means the thresholding regularization is practically no worse than the original squared regularization. The lower MSE of kRegBayes compared with pKBR shows that the posterior regularization successfully incorporates information from equations of the dynamics. Moreover, pKBR and kRegBayes run faster than KBR. The total running times for 50 random datasets of pKBR, kRegBayes and KBR are respectively 601.3s, 677.5s and 3667.4s.

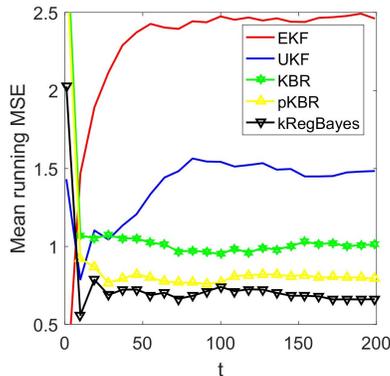


Figure 1: Mean running MSEs against time steps for each algorithm. (Best view in color)

**Camera position recovery** In this experiment, we build a scene containing a table and a chair, which is derived from `classchair.pov` (<http://www.oyonale.com>). With a fixed focal point, the position of the camera uniquely determines the view of the scene. The task of this experiment is to estimate the position of the camera given the image. This is a problem with practical applications in remote sensing and robotics.

We vary the position of the camera in a plane with a fixed height. The transition equations of the hidden states are

$$\theta_{t+1} = \theta_t + 0.2 + \xi_\theta, \quad r_{t+1} = \max(R_2, \min(R_1, r_t + \xi_r)), \quad x_{t+1} = \cos \theta_{t+1}, \quad y_{t+1} = \sin \theta_{t+1},$$

where  $\xi_\theta \sim \mathcal{N}(0, 4e - 4)$ ,  $\xi_r \sim \mathcal{N}(0, 1)$ ,  $0 \leq R_1 < R_2$  are two constants and  $\{(x_t, y_t)\}_{t=1}^m$  are treated as the hidden variables. As the observation at  $t$ -th step, we render a  $100 \times 100$  image with the camera located at  $(x_t, y_t)$ . For training data, we set  $R_1 = 0$  and  $R_2 = 10$  while for validation data and test data we set  $R_1 = 5$  and  $R_2 = 7$ . The motivation is to distinguish the efficacy of enforcing the posterior distribution to concentrate around distance 6 by kRegBayes. We show a sample set of training and test images in Fig. 2.

We compare KBR, pKBR and kRegBayes with the traditional linear Kalman filter (KF [23]). Following [4] we down-sample the images and train a linear regressor for observation model. In all experiments, we flatten the images to a column vector and apply Gaussian RBF kernels if needed. The kernel band widths are set to be the median distances in the training data. Based on experiments on the validation dataset, we set  $\lambda_T = 1e - 6 = 2\delta_T$  and  $\mu_T = 1e - 5$ .



Figure 2: First several frames of training data (upper row) and test data (lower row).

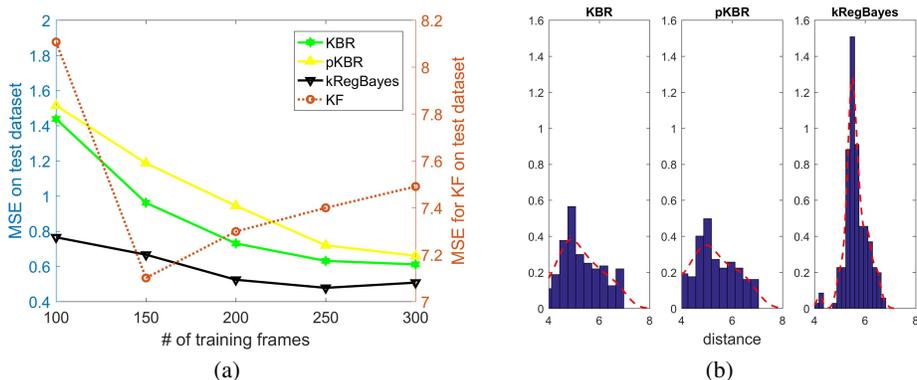


Figure 3: (a) MSEs for different algorithms (best view in color). Since KF performs much worse than kernel filters, we use a different scale and plot it on the right  $y$ -axis. (b) Probability histograms for the distance between each state and the scene center. All algorithms use 100 training data.

To provide supervision for kRegBayes, we uniformly generate 2000 data points  $\{(\hat{x}_i, \hat{y}_t)\}_{i=1}^{2000}$  on the circle  $r = 6$ . Given the previous estimate  $(\tilde{x}_t, \tilde{y}_t)$ , we first compute  $\hat{\theta}_t = \arctan(\hat{y}_t/\hat{x}_t)$  (where the value  $\hat{\theta}_t$  is adapted according to the quadrant of  $(\hat{x}_t, \hat{y}_t)$ ) and estimate  $(\tilde{x}_{t+1}, \tilde{y}_{t+1}) = (\cos(\hat{\theta}_t + 0.4), \sin(\hat{\theta}_t + 0.4))$ . Next, we find the nearest point to  $(\tilde{x}_{t+1}, \tilde{y}_{t+1})$  in the supervision set  $(\tilde{x}_k, \tilde{y}_k)$  and add the regularization  $\mu_T \|\mu(\mathcal{I}_{t+1}) - \phi(\tilde{x}_k, \tilde{y}_k)\|$  to the posterior embedding, where  $\mathcal{I}_{t+1}$  denotes the  $(t + 1)$ -th image.

We vary the size of training dataset from 100 to 300 and report the results of KBR, pKBR, kRegBayes and KF on 200 test images in Fig. 3. KF performs much worse than all three kernel filters due to the extreme non-linearity. The result of pKBR is a little worse than that of KBR, but the gap decreases as the training dataset becomes larger. kRegBayes always performs the best. Note that the advantage becomes less obvious as more data come. This is because kernel methods can learn the distance relation better with more data, and posterior regularization tends to be more useful when data are not abundant and domain knowledge matters. Furthermore, Fig. 3(b) shows that the posterior regularization helps the distances to concentrate.

## 6 Conclusions

We propose an optimizational framework for kernel Bayesian inference. With thresholding regularization, the minimizer of the framework is shown to be a reasonable estimator of the posterior kernel embedding. In addition, we propose a posterior regularized kernel Bayesian inference framework called kRegBayes. These frameworks are applied to non-linear state-space filtering tasks and the results of different algorithms are compared extensively.

## Acknowledgements

We thank all the anonymous reviewers for valuable suggestions. The work was supported by the National Basic Research Program (973 Program) of China (No. 2013CB329403), National NSF of China Projects (Nos. 61620106010, 61322308, 61332007), the Youth Top-notch Talent Support Program, and Tsinghua Initiative Scientific Research Program (No. 20141080934).

## References

- [1] Alex J Smola and Bernhard Schölkopf. *Learning with kernels*. Citeseer, 1998.
- [2] Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media, 2011.
- [3] Alex Smola, Arthur Gretton, Le Song, and Bernhard Schölkopf. A hilbert space embedding for distributions. In *Algorithmic learning theory*, pages 13–31. Springer, 2007.
- [4] 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.
- [5] Kenji Fukumizu, Le Song, and Arthur Gretton. Kernel bayes’ rule. In *Advances in neural information processing systems*, pages 1737–1745, 2011.
- [6] Le Song, Kenji Fukumizu, and Arthur Gretton. Kernel embeddings of conditional distributions: A unified kernel framework for nonparametric inference in graphical models. *Signal Processing Magazine, IEEE*, 30(4):98–111, 2013.
- [7] Le Song, Byron Boots, Sajid M Siddiqi, Geoffrey J Gordon, and Alex Smola. Hilbert space embeddings of hidden markov models. 2010.
- [8] Le Song, Arthur Gretton, and Carlos Guestrin. Nonparametric tree graphical models. In *International Conference on Artificial Intelligence and Statistics*, pages 765–772, 2010.
- [9] Steffen Grunewalder, Guy Lever, Luca Baldassarre, Massi Pontil, and Arthur Gretton. Modelling transition dynamics in mdps with rkhs embeddings. *arXiv preprint arXiv:1206.4655*, 2012.
- [10] Yu Nishiyama, Abdeslam Boularias, Arthur Gretton, and Kenji Fukumizu. Hilbert space embeddings of pomdps. *arXiv preprint arXiv:1210.4887*, 2012.
- [11] Peter M. Williams. Bayesian conditionalisation and the principle of minimum information. *The British Journal for the Philosophy of Science*, 31(2), 1980.
- [12] Charles A Micchelli and Massimiliano Pontil. On learning vector-valued functions. *Neural computation*, 17(1):177–204, 2005.
- [13] Steffen Grunewalder, Guy Lever, Luca Baldassarre, Sam Patterson, Arthur Gretton, and Massimiliano Pontil. Conditional mean embeddings as regressors. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, pages 1823–1830, 2012.
- [14] Jun Zhu, Ning Chen, and Eric P Xing. Bayesian inference with posterior regularization and applications to infinite latent svms. *The Journal of Machine Learning Research*, 15(1):1799–1847, 2014.
- [15] Charles A Micchelli, Yuesheng Xu, and Haizhang Zhang. Universal kernels. *The Journal of Machine Learning Research*, 7:2651–2667, 2006.
- [16] Nachman Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, 68(3):337–404, 1950.
- [17] Kuzman Ganchev, Joao Graça, Jennifer Gillenwater, and Ben Taskar. Posterior regularization for structured latent variable models. *The Journal of Machine Learning Research*, 11:2001–2049, 2010.
- [18] Jun Zhu, Amr Ahmed, and Eric Xing. MedLDA: Maximum margin supervised topic models. *JMLR*, 13:2237–2278, 2012.
- [19] Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- [20] Ernesto De Vito and Andrea Caponnetto. Risk bounds for regularized least-squares algorithm with operator-value kernels. Technical report, DTIC Document, 2005.
- [21] Simon J Julier and Jeffrey K Uhlmann. New extension of the kalman filter to nonlinear systems. In *AeroSense’97*, pages 182–193. International Society for Optics and Photonics, 1997.
- [22] Eric A Wan and Ronell Van Der Merwe. The unscented kalman filter for nonlinear estimation. In *Adaptive Systems for Signal Processing, Communications, and Control Symposium 2000. AS-SPCC. The IEEE 2000*, pages 153–158. Ieee, 2000.
- [23] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1):35–45, 1960.
- [24] Alexander J Smola and Risi Kondor. Kernels and regularization on graphs. In *Learning theory and kernel machines*, pages 144–158. Springer, 2003.
- [25] Heinz Werner Engl, Martin Hanke, and Andreas Neubauer. *Regularization of inverse problems*, volume 375. Springer Science & Business Media, 1996.

## A Appendix

### A.1 Kernel filtering

We first review how to use kernel techniques to do state-space filtering [5]. Assume that a sample  $(y_1, x_1, \dots, y_{T+1}, x_{T+1})$  is given, in which  $y_i \in \mathcal{Y}$  is the state and  $x_i \in \mathcal{X}$  is the corresponding observation. The transition and observation probabilities are estimated empirically in a nonparametric way:

$$\widehat{\mathcal{C}}_{Y Y_+} = \frac{1}{T} \sum_{i=1}^T \psi(y_i) \otimes \psi(y_{i+1}), \quad \widehat{\mathcal{C}}_{Y X} = \frac{1}{T} \sum_{i=1}^T \psi(y_i) \otimes \phi(x_i).$$

The filtering task is composed of two steps. The first step is to predict the next state based on current state, *i.e.*,  $p(Y_{t+1} | X_1, \dots, X_t) = \int p(Y_{t+1} | Y_t) p(Y_t | X_1, \dots, X_t) dY_t$ . The second step is to update the state based on a new observation  $x_{t+1}$  via Bayes' rule, *i.e.*,  $p(Y_{t+1} | X_1, \dots, X_{t+1}) \propto p(Y_{t+1} | X_1, \dots, X_t) p(X_{t+1} | Y_{t+1})$ . Following these two steps, we can obtain a recursive kernel update formula under different assumptions of the forms of kernel embedding  $\widehat{m}_{y_t | x_1, \dots, x_t}$ .

For kernel embeddings without posterior regularization, we suppose  $\widehat{m}_{y_t | x_1, \dots, x_t} = \sum_{i=1}^T \alpha_i^{(t)} \psi(y_i)$ . According to Thm. 1, the prediction step is realized by  $\widehat{m}_{y_{t+1} | x_1, \dots, x_t} = \widehat{\mathcal{C}}_{Y_+ Y} (\widehat{\mathcal{C}}_{Y Y} + \lambda_T I)^{-1} \widehat{m}_{y_t | x_1, \dots, x_t} = \Psi_+ (G_Y + T \lambda_T I)^{-1} G_Y \alpha^{(t)}$ , where  $\Psi_+ = (\psi(y_2), \dots, \psi(y_{T+1}))$ ,  $G_Y$  is the Gram matrix of  $\{y_1, \dots, y_T\}$  and  $\alpha^{(t)}$  is the vector of coefficients. The update step can be realized by invoking Prop. 2, *i.e.*,  $\widehat{m}_{y_{t+1} | x_1, \dots, x_{t+1}} = \Psi (K_X + \delta_T \Lambda^+)^{-1} K_{:x_{t+1}}$ , where  $K_X$  is the Gram matrix for  $(x_1, \dots, x_t)$ ,  $\Lambda^+ = \text{diag}(1/\beta^+)$  and  $\beta = (G_Y + T \lambda_T I)^{-1} G_{Y Y_+} (G_Y + T \lambda_T I)^{-1} G_Y \alpha^{(t)}$ , where  $(G_{Y Y_+})_{ij} = k_Y(y_i, y_{i+1})$ . The update formula of  $\alpha^{(t+1)}$  can then be summarized as follows

$$\alpha^{(t+1)} = (K_X + \delta_T \Lambda^+)^{-1} K_{:x_{t+1}}. \quad (13)$$

For kernel embeddings with posterior regularization, we suppose that for each step  $t$ , the regularization  $\mu_T \|\mu(\tilde{x}_t) - \psi(\tilde{y}_t)\|$  is used, meaning that  $p(Y_t | X_1, \dots, X_t = \tilde{x}_t)$  is encouraged to concentrate around  $\delta(Y_t = \tilde{y}_t)$ . To obtain a recursive formula, we assume that  $\widehat{m}_{y_t | x_1, \dots, x_t} = \sum_{i=1}^T \alpha_i^{(t)} \psi(y_i) + \sum_{i=1}^N \tilde{\alpha}_i^{(t)} \psi(\tilde{y}_i)$ , where  $N$  is the number of supervision data points  $(\tilde{x}_i, \tilde{y}_i)$ . Following a similar logic except replacing Prop. 2 with Prop. 3, we get the update rule for  $\alpha^{(t+1)}$  and  $\tilde{\alpha}^{(t+1)}$

$$\gamma = (K_X' + \delta_T \Lambda^\dagger)^{-1} K_{:x_{t+1}}' \quad (14)$$

$$\alpha^{(t+1)} = \gamma[1 : m] \quad (15)$$

$$\tilde{\alpha}^{(t+1)} = (0, \dots, \gamma[m+1], 0, \dots)^\top, \quad (16)$$

where  $\Lambda^\dagger = \text{diag}(1/\beta^+, 1/\mu_T)$ ,  $\beta = (G_Y + T \lambda_T I)^{-1} G_{Y Y_+} (G_Y + T \lambda_T I)^{-1} (G_{Y Y} \alpha^{(t)} + G_{Y \tilde{Y}} \tilde{\alpha}^{(t)})$ .  $K_X'$  and  $K_{:x_{t+1}}'$  are augmented Gram matrices, which incorporate  $(\tilde{x}_i, \tilde{y}_i)$ . The position of  $\gamma[m+1]$  in  $\tilde{\alpha}^{(t+1)}$  corresponds to the index of supervision  $(\tilde{x}_k, \tilde{y}_k)$  at  $t+1$  step in  $\{(\tilde{x}_i, \tilde{y}_i)\}_{i=m+1}^n$ .

To obtain  $\alpha^{(1)}$ , we use conditional operators [4] to estimate  $m_{y_1}$  without priors. We set  $\alpha^{(1)} = (K_X + T \lambda_T I)^{-1} K_{:x_1}$  for both types of kernel filtering and  $\tilde{\alpha}^{(1)} = \mathbf{0}$ . To decode the state from kernel embeddings, we solve an optimization problem  $\hat{y}_t = \arg \min_y \|m(x) - \psi(y)\|$ , which can be computed using an iteration scheme as depicted in [4].

### A.2 Proofs

**Proposition 1.** *Suppose  $(X, Y)$  is a random variable in  $\mathcal{X} \times \mathcal{Y}$ , where the prior for  $Y$  is  $\pi(Y)$  and the likelihood is  $p(X | Y)$ . Let  $\mathcal{H}_X$  be a RKHS with kernel  $k_X$  and feature map  $\phi(x)$ ,  $\mathcal{H}_Y$  be a RKHS with kernel  $k_Y$  and feature map  $\psi(y)$ ,  $\phi(x, y)$  be the feature map of  $\mathcal{H}_X \otimes \mathcal{H}_Y$ ,  $\widehat{\pi}_Y = \sum_{i=1}^l \tilde{\alpha}_i \psi(\tilde{y}_i)$  be an estimator for  $\pi_Y$  and  $\{(x_i, y_i)\}_{i=1}^n$  be a sample representing  $p(X | Y)$ . Under the assumption that  $f(x, y) = \|\psi(y) - \mu(x)\|_{\mathcal{H}_Y}^2 \in \mathcal{H}_X \otimes \mathcal{H}_Y$ , we have*

$$\widehat{\mathcal{E}}_s[\mu] = \sum_{i=1}^n \beta_i \|\psi(y_i) - \mu(x_i)\|_{\mathcal{H}_Y}^2, \quad (17)$$

where  $\beta = (\beta_1, \dots, \beta_n)^\top$  is given by  $\beta = (G_Y + n\lambda I)^{-1} \tilde{G}_Y \tilde{\alpha}$ , where  $(G_Y)_{ij} = k_{\mathcal{Y}}(y_i, y_j)$ ,  $(\tilde{G}_Y)_{ij} = k_{\mathcal{Y}}(y_i, \tilde{y}_j)$ , and  $\tilde{\alpha} = (\tilde{\alpha}_1, \dots, \tilde{\alpha}_l)^\top$ .

*Proof.* The reasoning is similar to [5], Prop. 5. We only need to show that  $\hat{\mu}_{(X,Y)} = \Phi_{X,Y} \beta = \Phi_{X,Y} (G_Y + n\lambda I)^{-1} \tilde{G}_Y \tilde{\alpha}$ , where  $\Phi_{X,Y} = (\phi(x_1, y_1), \dots, \phi(x_n, y_n))$ . Recall that  $\hat{\mu}_{(X,Y)} = \hat{C}_{(X,Y)Y} (\hat{C}_{YY} + \lambda I)^{-1} \hat{\pi}_Y$ . Let  $h = (\hat{C}_{YY} + \lambda I)^{-1} \hat{\pi}_Y$  and decompose it as  $h = \sum_{i=1}^n a_i \psi(y_i) + h_\perp$ , where  $h_\perp$  is perpendicular to  $\text{span}\{\psi(y_1), \dots, \psi(y_n)\}$ . Expanding  $(\hat{C}_{YY} + \lambda I)h = \hat{\pi}_Y$ , we obtain

$$\frac{1}{n} \sum_{i,j \leq n} a_i k_{\mathcal{Y}}(y_i, y_j) \psi(y_j) + \lambda \left( \sum_{i \leq n} a_i \psi(y_i) + h_\perp \right) = \sum_{i \leq l} \tilde{\alpha}_i \psi(\tilde{y}_i). \quad (18)$$

Multiplying both sides with  $\psi(y_k)|_{k=1}^n$ , we get  $\frac{1}{n} G_Y^2 \mathbf{a} + \lambda G_Y \mathbf{a} = \tilde{G}_Y \tilde{\alpha}$ . Therefore  $\hat{\mu}_{(X,Y)}$  can be written as  $\hat{\mu}_{(X,Y)} = \frac{1}{n} [\sum_{i \leq n} \phi(x_i, y_i) \otimes \psi(y_i)] h = \frac{1}{n} \Phi_{X,Y} G_Y \mathbf{a} = \Phi_{X,Y} (G_Y + n\lambda I)^{-1} \tilde{G}_Y \tilde{\alpha}$ .  $\square$

**Proposition 2.** *Without loss of generality, we assume that  $\beta_i^+ \neq 0$  for all  $1 \leq i \leq n$ . Let  $\mu \in \mathcal{H}_K$  and choose the kernel of  $\mathcal{H}_K$  to be  $K(x_i, x_j) = k_{\mathcal{X}}(x_i, x_j) \mathcal{I}$ , where  $\mathcal{I} : \mathcal{H}_K \rightarrow \mathcal{H}_K$  is an identity map. Then*

$$\hat{\mu}_{\lambda,n}(x) = \Psi (K_X + \lambda_n \Lambda^+)^{-1} K_{:x}, \quad (19)$$

where  $\Psi = (\psi(y_1), \dots, \psi(y_n))$ ,  $(K_X)_{ij} = k_{\mathcal{X}}(x_i, x_j)$ ,  $\Lambda^+ = \text{diag}(1/\beta_1^+, \dots, 1/\beta_n^+)$ ,  $K_{:x} = (k_{\mathcal{X}}(x, x_1), \dots, k_{\mathcal{X}}(x, x_n))^\top$  and  $\lambda_n$  is a positive regularization constant.

*Proof.* If  $\beta_i^+ = 0$  for any  $i$ , we can discard the data point  $(x_i, y_i)$  without affecting results. Let  $\mu = \mu_0 + g$ , where  $\mu_0 = \sum_{i=1}^n K_{x_i} c_i$ . Plugging  $\mu = \mu_0 + g$  into  $\hat{\mathcal{E}}_{\lambda,n}[\mu]$  and expand, we obtain  $\hat{\mathcal{E}}_{\lambda,n}[\mu] = \sum_{i=1}^n \beta_i^+ \|\psi(y_i) - \mu_0(x_i)\|^2 + \lambda_n \|\mu_0\|^2 + \sum_{i=1}^n \beta_i^+ \|g(x_i)\|^2 + \lambda_n \|g\|^2 + 2\lambda_n \langle \mu_0, g \rangle - 2 \sum_{i=1}^n \beta_i^+ \langle g(x_i), \psi(y_i) - \mu_0(x_i) \rangle$ .

We conjecture that  $\psi(y_i) - \sum_{j=1}^n k_{\mathcal{X}}(x_i, x_j) c_j = \frac{\lambda_n}{\beta_i^+} c_i$ , for all  $1 \leq i \leq n$ . Actually, substituting these equations into  $\hat{\mathcal{E}}_{\lambda,n}[\mu]$  gives the relation  $\lambda_n \langle \mu_0, g \rangle - \sum_{i=1}^n \beta_i^+ \langle g(x_i), \psi(y_i) - \mu_0(x_i) \rangle = 0$ . As a result,  $\hat{\mathcal{E}}_{\lambda,n}[\mu] = \hat{\mathcal{E}}_{\lambda,n}[\mu_0] + \sum_{i=1}^n \beta_i^+ \|g(x_i)\|^2 + \lambda_n \|g\|^2 \geq \hat{\mathcal{E}}_{\lambda,n}[\mu_0]$ , which means that  $\mu_0 = \sum_{i=1}^n K_{x_i} c_i$  with  $c_i$  satisfying the conjectured equations is the solution. The equation  $\psi(y_i) - \sum_{j=1}^n k_{\mathcal{X}}(x_i, x_j) c_j = \frac{\lambda_n}{\beta_i^+} c_i$  implies that  $(K_X + \lambda_n \Lambda^+) c = \Psi$  and  $\mu_0(x) = \sum_{i=1}^n k_{\mathcal{X}}(x, x_i) c_i = \Psi (K_X + \lambda_n \Lambda^+)^{-1} K_{:x}$ .  $\square$

**Theorem 6.** *Assume that  $|\mathcal{X} \times \mathcal{Y}| < \infty$ ,  $k$  is strictly positive definite with  $\sup_{(x,y)} k((x,y), (x,y)) < \kappa$  and  $f(x,y) = \|\psi(y) - \mu(x)\|_{\mathcal{H}_Y}^2 \in \mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}$ . With the conditions in Thm. 2, we assert that  $\hat{\mu}_{(X,Y)}^+$  is a consistent estimator of  $\mu_{(X,Y)}$  and  $|\hat{\mathcal{E}}_s^+[\mu] - \mathcal{E}_s[\mu]| \rightarrow 0$  in probability as  $n \rightarrow \infty$ .*

*Proof.* We only need to show that  $\hat{\mu}_{(X,Y)}^+ := \sum_{i=1}^n \beta_i^+ \phi(x_i) \otimes \psi(y_i)$  converges to  $\mu_{(X,Y)}$  in probability as  $n \rightarrow \infty$ , since  $|\hat{\mathcal{E}}_s^+[\mu] - \mathcal{E}_s[\mu]| = |\langle f, \hat{\mu}_{(X,Y)}^+ - \mu_{(X,Y)} \rangle| \leq \|f\| \left\| \hat{\mu}_{(X,Y)}^+ - \mu_{(X,Y)} \right\|$ . From Thm. 2 we know that  $\hat{\mu}_{(X,Y)}$  converges to  $\mu_{(X,Y)}$  in probability, hence it is sufficient to show that  $\hat{\mu}_{(X,Y)}^+$  converges to  $\hat{\mu}_{(X,Y)}$  in RKHS norm as  $n \rightarrow \infty$ .

Let  $|\mathcal{X} \times \mathcal{Y}| = M$ . Without losing generality, we assume  $\mathcal{X} \times \mathcal{Y} = \{(x_1, y_1), \dots, (x_M, y_M)\}$  and  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  is a sample representing  $p(X | Y)$ . According to Theorem 4 in [24],  $k$  is strictly positive definite on a finite set implies that  $\mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}$  consists of all bounded functions on  $\mathcal{X} \times \mathcal{Y}$ . In particular,  $\mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}$  contains the function

$$g(x_i, y_i) = \begin{cases} 1, & \beta_i < 0 \\ 0, & \text{otherwise.} \end{cases} \quad (20)$$

We denote  $b := \max_g \|g\|_{\mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}} = \max_{\mathbf{g}} \mathbf{g}^\top K^{-1} \mathbf{g}$  for all possibilities of  $\beta$ . Here  $\mathbf{g}$  represents the point evaluations of  $g$  on  $\{(x_i, y_i)\}_{i=1}^M$  and  $K_{ij} = k((x_i, y_i), (x_j, y_j))$ . Note that  $g(x,y)$  is non-negative, thus  $\mathbb{E}[g(X,Y)] = \langle g, \mu_{(X,Y)} \rangle \geq 0$ . For sufficiently large  $n$ ,  $|\langle g, \hat{\mu}_{(X,Y)} -$

$\mu_{(X,Y)}\rangle| \leq \|g\| \|\widehat{\mu}_{(X,Y)} - \mu_{(X,Y)}\| \leq \epsilon b$  in arbitrarily high probability. In this case  $\langle g, \widehat{\mu}_{(X,Y)} \rangle = -\sum_{i=1}^n \beta_i^- \geq -\epsilon b$ , where  $\beta_i^- = -\min(0, \beta_i)$ , and  $\|\widehat{\mu}_{(X,Y)}^+ - \widehat{\mu}_{(X,Y)}\| = \|\sum_{i=1}^n \beta_i^- \phi(x_i, y_i)\| = \sqrt{\sum_{i,j} \beta_i^- \beta_j^- k((x_i, y_i), (x_j, y_j))} \leq \sqrt{\kappa} \sum_{i=1}^n \beta_i^- \leq \epsilon b \sqrt{\kappa}$ . The inequalities can now be linked and the theorem proved.  $\square$

**Theorem 7.** Assume that  $|\mathcal{X} \times \mathcal{Y}| < \infty$ ,  $k$  is strictly positive definite with  $\sup_{(x,y)} k((x,y), (x,y)) < \kappa$ , we assert  $\sum_{i=1}^n \beta_i^+ \rightarrow 1$  in probability as  $n \rightarrow \infty$ .

*Proof.* The proof follows a similar reasoning to that in Thm. 6. Let  $|\mathcal{X} \times \mathcal{Y}| = M$  and  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  be a sample representing  $p(X | Y)$ . According to Theorem 4 in [24],  $k$  is strictly positive definite on a finite set implies that  $\mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}$  consists of all bounded functions on  $\mathcal{X} \times \mathcal{Y}$ . In particular,  $\mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}$  contains the function  $f(x, y) \equiv 1$ . From Thm. 6 we know that  $\widehat{\mu}^+ = \sum_{i=1}^n \beta_i^+ \phi(x_i) \otimes \psi(y_i) \rightarrow \mu$  in probability. Therefore,  $|\sum_{i=1}^n \beta_i^+ - 1| = |\langle f, \widehat{\mu}_{(X,Y)}^+ - \mu_{(X,Y)} \rangle| \leq \|f\| \|\widehat{\mu}_{(X,Y)}^+ - \mu_{(X,Y)}\| \rightarrow 0$  in probability.  $\square$

Since  $\beta_i$ 's do not depend on  $X_1, \dots, X_n$ , we have the following corollary:

**Corollary 1.** Assume that  $|\mathcal{Y}| < \infty$ ,  $k$  is strictly positive definite with  $\sup_{(x,y)} k((x,y), (x,y)) < \kappa$ , we assert  $\sum_{i=1}^n \beta_i^+ \rightarrow 1$  in probability as  $n \rightarrow \infty$ .

Next, we will relax the finite space condition on  $\mathcal{X} \times \mathcal{Y}$  in Thm. 6. To this end, we introduce the following convenient concept of  $\epsilon$ -partition.

**Definition 1** ( $\epsilon$ -partition). An  $\epsilon$ -partition of a metric space  $\mathcal{X}$  is a partition whose elements are all within  $\epsilon$ -balls of  $\mathcal{X}$ .

Since a compact space is totally bounded, we have the more general result.

**Theorem 3.** Assume that  $\mathcal{X}$  is compact and  $|\mathcal{Y}| < \infty$ ,  $k$  is a strictly positive definite continuous kernel with  $\sup_{(x,y)} k((x,y), (x,y)) < \kappa$  and  $f(x, y) = \|\psi(y) - \mu(x)\|_{\mathcal{H}_{\mathcal{Y}}}^2 \in \mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}$ . With the conditions in Thm. 2, we assert that  $\widehat{\mu}_{(X,Y)}^+$  is a consistent estimator of  $\mu_{(X,Y)}$  and  $|\widehat{\mathcal{E}}_s^+[\mu] - \mathcal{E}_s[\mu]| \rightarrow 0$  in probability as  $n \rightarrow \infty$ .

*Proof.* From the condition that  $\phi(x, y)$  is continuous on the compact space  $\mathcal{X} \times \mathcal{Y}$ , we know  $\phi(x, y)$  and  $\phi(x)$  are uniformly continuous.

For any probability measure  $p$  and  $\epsilon$ -partition of  $\mathcal{X}$ , we can construct a new discretized probability measure in the following way. Suppose the  $\epsilon$ -partition is  $\{B_1^\epsilon, B_2^\epsilon, \dots\}$ , we identify each set  $B_i^\epsilon$  with a representative element  $x_i^\epsilon \in B_i^\epsilon$ . The resulting probability measure is denoted as  $p^\epsilon$  and satisfies  $p^\epsilon(A) = \sum_{x_i^\epsilon \in A} p(B_i^\epsilon)$ . We also define the discretization  $x_i^\epsilon$  of  $x_i$  to be  $x_i^\epsilon = x_j^\epsilon$  if  $x_i \in B_j$ . Let the kernel embedding of  $p$  be  $\mu$  and  $p^\epsilon$  be  $\mu^\epsilon$ . Suppose  $\forall \delta > 0, \exists \epsilon > 0$  such that  $\|x_1 - x_2\| \leq \epsilon$  implies  $\|\phi(x_1) - \phi(x_2)\|_{\mathcal{H}_{\mathcal{X}}} \leq \delta$ . We assert that  $\|\mu - \mu^\epsilon\| \leq \delta$ . To prove this, we observe that an i.i.d. sample  $\{x_1, \dots, x_n\}$  from  $p$  is also an i.i.d. sample of  $p^\epsilon$  if we replace  $x_i$  with  $x_i^\epsilon$ . Since the estimator  $\widehat{\mu} = \frac{1}{n} \sum_{i=1}^n \phi(x_i)$  is a consistent estimator of  $\mu$ , we know that  $\widehat{\mu}^\epsilon = \frac{1}{n} \sum_{i=1}^n \phi(x_i^\epsilon)$  is also consistent. Via consistency, we have that with no less than any high probability  $1 - \Delta$ , for any  $n > N(\Delta, \delta', \epsilon)$ ,  $\|\widehat{\mu} - \mu\| \leq \delta'$  and  $\|\widehat{\mu}^\epsilon - \mu^\epsilon\| \leq \delta'$  holds. Since  $\|\widehat{\mu} - \widehat{\mu}^\epsilon\| \leq \frac{1}{n} \sum_{i=1}^n \|\phi(x_i) - \phi(x_i^\epsilon)\|$  and  $\|x_i - x_i^\epsilon\| \leq \epsilon$ , we have  $\|\widehat{\mu} - \widehat{\mu}^\epsilon\| \leq \delta$  from uniform continuity. Combining this with  $\|\widehat{\mu} - \mu\| \leq \delta'$  and  $\|\widehat{\mu}^\epsilon - \mu^\epsilon\| \leq \delta'$  we know  $\|\mu - \mu^\epsilon\| \leq \|\mu - \widehat{\mu}\| + \|\widehat{\mu} - \widehat{\mu}^\epsilon\| + \|\widehat{\mu}^\epsilon - \mu^\epsilon\| \leq \delta + 2\delta'$  with high probability  $1 - \Delta$ . Note that  $\|\mu - \mu^\epsilon\| \leq \delta + 2\delta'$  is a deterministic event and holds for any  $\delta' > 0$ , we have  $\|\mu - \mu^\epsilon\| \leq \delta$ .

Now we would like to discretize  $X$  for  $\mu_{(X,Y)}$ . For any  $\epsilon > 0$ , we have  $\widehat{\mu}_{(X,Y)} = \sum_{i=1}^n \beta_i \phi(x_i, y_i) \rightarrow \mu_{(X,Y)}$  in probability and  $\widehat{\mu}_{(X,Y)}^\epsilon = \sum_{i=1}^n \beta_i^\epsilon \phi(x_i^\epsilon, y_i) \rightarrow \mu_{(X,Y)}^\epsilon$  in probability. Since  $\beta_i$  depends only on  $y_1, \dots, y_n$ , we have  $\beta_i = \beta_i^\epsilon$ . From the last paragraph we suppose that

$\epsilon$  is chosen such that  $\forall \|(x_i^\epsilon, y_i) - (x_i, y_i)\| \leq \epsilon, \|\phi(x_i^\epsilon, y_i) - \phi(x_i, y_i)\| \leq \delta$ . Note that

$$\begin{aligned} \left\| \sum_{i=1}^n \beta_i^+ \phi(x_i, y_i) - \mu_{(X,Y)} \right\| &\leq \left\| \sum_{i=1}^n \beta_i^+ \phi(x_i, y_i) - \sum_{i=1}^n \beta_i^+ \phi(x_i^\epsilon, y_i) \right\| \\ &\quad + \left\| \sum_{i=1}^n \beta_i^+ \phi(x_i^\epsilon, y_i) - \sum_{i=1}^n \beta_i \phi(x_i^\epsilon, y_i) \right\| \\ &\quad + \left\| \sum_{i=1}^n \beta_i \phi(x_i^\epsilon, y_i) - \mu_{(X,Y)}^\epsilon \right\| + \left\| \mu_{(X,Y)}^\epsilon - \mu_{(X,Y)} \right\| \\ &\leq \delta \sum_{i=1}^n \beta_i^+ + \delta + \left\| \sum_{i=1}^n \beta_i^+ \phi(x_i^\epsilon, y_i) - \sum_{i=1}^n \beta_i \phi(x_i^\epsilon, y_i) \right\| \\ &\quad + \left\| \sum_{i=1}^n \beta_i \phi(x_i^\epsilon, y_i) - \mu_{(X,Y)}^\epsilon \right\|. \end{aligned}$$

From Corollary 1, Thm. 6 and the consistency of  $\sum_{i=1}^n \beta_i \phi(x_i^\epsilon, y_i)$  we see  $\left\| \sum_{i=1}^n \beta_i^+ \phi(x_i, y_i) - \mu_{(X,Y)} \right\|$  can be arbitrarily small with arbitrarily high probability. This proves the consistency of  $\sum_{i=1}^n \beta_i^+ \phi(x_i, y_i)$ .  $\square$

**Corollary 2.** Assume that  $\mathcal{X}$  is compact and  $|\mathcal{Y}| < \infty$ ,  $k$  is a bounded strictly positive definite continuous kernel,  $k_{\mathcal{X}}$  is a bounded kernel with  $\sup_x k_{\mathcal{X}}(x, x) \leq \kappa_{\mathcal{X}}$ , we assert that  $\widehat{\mu}_X^+ = \sum_{i=1}^n \beta_i^+ \phi(x_i)$  is a consistent estimator of  $\mu_X$ , i.e., the kernel embedding of the marginal distribution on  $X$ .

**Theorem 8.** Let  $\mathcal{B}_1, \mathcal{B}_2$  be Banach spaces. For any linear operator  $\mathcal{A} : \mathcal{B}_1 \rightarrow \mathcal{B}_2$ , we assert that there exists a subset  $\mathcal{F} \subseteq \mathcal{B}_1$  such that  $\mathcal{F}$  is dense in  $\mathcal{B}_1$  and  $\|\mathcal{A}f\|_{\mathcal{B}_2} \leq N \|f\|_{\mathcal{B}_1}$  for some constant  $N$  and any  $f \in \mathcal{F}$ .

*Proof.* Let  $M_k$  be the set of  $f \in \mathcal{B}_1$  satisfying  $\|\mathcal{A}f\|_{\mathcal{B}_2} \leq k \|f\|_{\mathcal{B}_1}$ . Clearly we have  $\mathcal{B}_1 = \bigcup_{k=1}^{\infty} M_k$ . Since  $\mathcal{B}_1$  is complete, we can invoke Baire category theorem to conclude that there exists an integer  $n$  such that  $M_n$  is dense in some sphere  $S_0 \subseteq \mathcal{B}_1$ . Consider the spherical shell  $P$  in  $S_0$  consisting of the points  $z$  for which

$$\beta < \|z - y_0\| < \alpha,$$

where  $0 < \beta < \alpha, y_0 \in M_n$ . Next, translate the spherical shell  $P$  so that its center coincides with the origin of coordinates to obtain spherical shell  $P_0$ . We now show that there is some set  $M_N$  dense in  $P_0$ . For every  $z \in M_n \cap P$ , we have

$$\begin{aligned} \|\mathcal{A}(z - y_0)\|_{\mathcal{B}_2} &\leq \|\mathcal{A}z\|_{\mathcal{B}_2} + \|\mathcal{A}y_0\|_{\mathcal{B}_2} \leq n(\|z\|_{\mathcal{B}_1} + \|y_0\|_{\mathcal{B}_1}) \leq n(\|z - y_0\|_{\mathcal{B}_1} + 2\|y_0\|_{\mathcal{B}_1}) \\ &= n \|z - y_0\|_{\mathcal{B}_1} [1 + 2\|y_0\|_{\mathcal{B}_1} / \|z - y_0\|_{\mathcal{B}_1}] \leq n \|z - y_0\|_{\mathcal{B}_1} [1 + 2\|y_0\|_{\mathcal{B}_1} / \beta]. \end{aligned}$$

Let  $N = n(1 + 2\|y_0\|_{\mathcal{B}_1} / \beta)$ , we have  $z - y_0 \in M_N$ . Since  $z - y_0 \in M_N$  is obtained from  $z \in M_n$  and  $M_n$  is dense in  $P$ , it is easy to see that  $M_N$  is dense in  $P_0$ . For any  $y \in \mathcal{B}_1$  except  $\|y\|_{\mathcal{B}_1} = 0$ , it is always possible to choose  $\lambda$  so that  $\beta < \|\lambda y\| < \alpha$  and we can construct a sequence  $y_k \in M_N$  that converges to  $\lambda y$ . This means there exists a sequence  $(1/\lambda)y_k$  converging to  $y$ . By virtue of  $(1/\lambda)y_k \in M_N$  and  $0 \in M_N$ , we conclude  $M_N$  is dense in  $\mathcal{B}_1$ .  $\square$

**Theorem 9.** Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $\xi$  be a random variable on  $\Omega$  taking values in a Hilbert space  $\mathcal{K}$ . Define  $\mathcal{A} : f \in \mathcal{K} \mapsto \langle f, \xi(\cdot) \rangle \in \mathcal{H}$ , where  $\mathcal{H}$  is a RKHS with feature maps  $\phi(\omega)$ . Let  $\mu$  be a kernel embedding for  $P^\pi$  and  $\widehat{\mu} = \sum_{i=1}^n \beta_i^+ \phi(\omega_i)$  be a consistent estimator of  $\mu$ . Assume  $\sum_{i=1}^n \beta_i^+ \rightarrow 1$  in probability and there are two positive constants  $H$  and  $\sigma$  such that  $\|\xi(\omega)\|_{\mathcal{K}} \leq \frac{H}{2}$  a.s. and  $\mathbb{E}_{P^\pi}[\|\xi\|_{\mathcal{K}}^2] \leq \sigma^2$ . Then for any  $\epsilon > 0$ ,

$$\lim_{n \rightarrow 0} P^l \left[ (\omega_1, \dots, \omega_l) \in \Omega^l \mid \left\| \sum_{i=1}^n \beta_i^+ \xi(\omega_i) - \mathbb{E}_{P^\pi}[\xi] \right\|_{\mathcal{K}} > \epsilon \right] = 0 \quad (21)$$

*Proof.* From the consistency of  $\widehat{\mu}$ , we know for every  $\epsilon_1$ , there exists  $N_{\epsilon_1}(\delta_1)$  such that  $\forall n > N_{\epsilon_1}(\delta_1)$ ,  $\|\widehat{\mu} - \mu\|_{\mathcal{H}} < \epsilon_1$  with probability no less than  $1 - \delta_1$ . Similarly, for every  $\epsilon_2$ , there exists  $N_{\epsilon_2}(\delta_2)$  such that  $\forall n > N_{\epsilon_2}(\delta_2)$ ,  $|\sum_{i=1}^n \beta_i^+ - 1| < \epsilon_2$  with probability no less than  $1 - \delta_2$ . Furthermore, with probability no less than  $1 - \delta_2$ ,  $\|\sum_{i=1}^n \beta_i^+ \xi(w_i) - \mathbb{E}_{P^\pi}[\xi]\|_{\mathcal{K}} \leq \sum_{i=1}^n \beta_i^+ \|\xi(w_i)\|_{\mathcal{K}} + \|\mathbb{E}_{P^\pi}[\xi]\|_{\mathcal{K}} \leq (1 + \epsilon_2) \|\xi(\omega)\|_{\mathcal{K}} + \mathbb{E}_{P^\pi}[\|\xi\|] \leq \frac{H(1+\epsilon_2)}{2} + \sqrt{\mathbb{E}_{P^\pi}[\|\xi\|_{\mathcal{K}}^2]} = \frac{H(1+\epsilon_2)}{2} + \sigma$ , where the last two inequalities follow from Jensen's inequality.

Let  $f = \sum_{i=1}^n \beta_i^+ \xi(w_i) - \mathbb{E}_{P^\pi}[\xi]$  and clearly  $\|f\|_{\mathcal{K}} \leq \frac{H(1+\epsilon_2)}{2} + \sigma$ . Consider  $\Delta_f := \sum_{i=1}^n \beta_i^+ \langle f, \xi(w_i) \rangle - \langle f, \mathbb{E}_{P^\pi}[\xi] \rangle = \sum_{i=1}^n \beta_i^+ \langle \mathcal{A}f \rangle(w_i) - \mathbb{E}_{P^\pi}[\mathcal{A}f] = \langle \widehat{\mu} - \mu, \mathcal{A}f \rangle$ . In virtue of Thm. 8, for any  $\epsilon_3$ , there exists an element  $g \in \mathcal{K}$  and constant  $N$  (only depends on  $\mathcal{A}$ ) such that  $\|g - f\|_{\mathcal{K}} < \epsilon_3$  and  $\|\mathcal{A}g\|_{\mathcal{H}} \leq N \|g\|_{\mathcal{K}}$ . Similarly define  $\Delta_g := \sum_{i=1}^n \beta_i^+ \langle g, \xi(w_i) \rangle - \langle g, \mathbb{E}_{P^\pi}[\xi] \rangle = \langle \widehat{\mu} - \mu, \mathcal{A}g \rangle$ . It is easy to see that  $|\Delta_g - \Delta_f| \leq (1 + \epsilon_2)\epsilon_3 \|\xi(\omega)\|_{\mathcal{K}} + \epsilon_3 \|\mathbb{E}_{P^\pi}[\xi]\|_{\mathcal{K}} \leq \frac{H\epsilon_3(1+\epsilon_2)}{2} + \epsilon_3\sigma$  and  $\Delta_g = \langle \widehat{\mu} - \mu, \mathcal{A}g \rangle \leq \epsilon_1 N \|g\|_{\mathcal{K}} \leq \epsilon_1 N(\epsilon_3 + \|f\|_{\mathcal{K}}) \leq \epsilon_1 N(\epsilon_3 + \sigma + \frac{H(1+\epsilon_2)}{2})$  with probability no less than  $1 - \delta_1 - \delta_2$ . Hence  $\|\sum_{i=1}^n \beta_i^+ \xi(w_i) - \mathbb{E}_{P^\pi}[\xi]\|_{\mathcal{K}} = \sqrt{|\Delta_f|} \leq \sqrt{\epsilon_1 N(\epsilon_3 + \sigma + \frac{H(1+\epsilon_2)}{2}) + \frac{H\epsilon_3(1+\epsilon_2)}{2} + \epsilon_3\sigma}$  with probability no less than  $1 - \delta_1 - \delta_2$  for all  $n > \max(N_{\epsilon_1}(\delta_1), N_{\epsilon_2}(\delta_2))$ . The theorem now gets proved.  $\square$

The proof of Thm. 5 is based on the proof of Thm. 5 in [20], with more assumptions and different concentration results. For convenience, we borrow some notations in their paper and refer the readers to [20] for definitions. We suggest the readers to be familiar with [20] because we modify and skip some details of the proofs to make the reasoning clearer.

Let  $\mathcal{X}, \mathcal{Y}$  be Polish spaces,  $\mathcal{H}_Y$  be a separable Hilbert space,  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ ,  $\mathcal{H}_K$  be a real Hilbert space of functions  $\mu : \mathcal{X} \rightarrow \mathcal{H}_Y$  satisfying  $\mu(x) = K_x^* \mu$  where  $K_x : \mathcal{H}_Y \rightarrow \mathcal{H}_K$  is the bounded operator  $K_x v = K(\cdot, x)v$ ,  $v \in \mathcal{H}_Y$ . Moreover, let  $T_x = K_x K_x^* \in \mathcal{L}_2(\mathcal{H}_K)$  be a positive Hilbert-Schmidt operator.

Let  $\rho$  be a probability measure on  $\mathcal{Z}$  and  $\rho_X$  denotes the marginal distribution of  $\rho$  on  $\mathcal{X}$ . We suppose that  $\rho = p(X | Y)\pi(Y)$  and thus it incorporates the information of the prior. In contrast, we are given a sample  $\mathbf{z} = ((x_1, y_1), \dots, (x_n, y_n))$  from another distribution on  $\mathcal{Z}$  with the same  $p(X | Y)$ .

The optimization objective now becomes  $\mathcal{E}_s[\mu] = \int_{\mathcal{Z}} \|\mu(x) - \phi(y)\|_{\mathcal{H}_Y}^2 d\rho(x, y)$ . Denote  $T = \int_{\mathcal{X}} T_x d\rho_X(x)$ ,  $T_{\mathbf{x}} = \sum_{i=1}^n \beta_i^+ T_{x_i}$ ,  $\mu_{\mathcal{H}_K} = \arg \min_{\mu} \mathcal{E}_s[\mu]$ ,  $\mu^\lambda = \mathcal{E}_s[\mu] + \lambda \|\mu\|_{\mathcal{H}_K}^2$  and  $\mu_\lambda^\lambda = \widehat{\mathcal{E}}_{\lambda, n}[\mu]$ . Additionally, let  $A : \mathcal{H}_K \rightarrow L^2(\mathcal{Z}, \rho, \mathcal{H}_Y)$  be the linear operator  $(Af)(x, y) = K_x^* f \quad \forall (x, y) \in \mathcal{Z}$  and  $A_{\mathbf{z}} := A_{\rho = \sum_{i=1}^n \beta_i^+ \delta_{x_i}}$ . Finally, let  $\mathcal{A}(\lambda) = \|\mu^\lambda - \mu_{\mathcal{H}_K}\|_{\rho}^2 = \|\sqrt{T}(\mu^\lambda - \mu_{\mathcal{H}_K})\|$ ,  $\mathcal{B}(\lambda) = \|\mu^\lambda - \mu_{\mathcal{H}_K}\|_{\mathcal{H}_K}^2$  and  $\mathcal{N}(\lambda) = \text{Tr}((T + \lambda)^{-1}T)$ .

**Assumption 1.** Let  $\mathcal{A}_1 : f \in \mathcal{L}_2(\mathcal{H}_K) \mapsto \langle f, (T + \lambda)^{-1}T \rangle \in \mathcal{H}_1$ ,  $\mathcal{A}_2 : f \in \mathcal{L}(\mathcal{H}_K) \mapsto \langle f, T(\mu^\lambda - \mu_{\mathcal{H}_K}) \rangle \in \mathcal{H}_2$ ,  $\mathcal{A}_3 : f \in \mathcal{H}_K \mapsto \langle f, (T + \lambda)^{-\frac{1}{2}} K_{\#1}(\psi(\#2) - \mu_{\mathcal{H}_K}(\#1)) \rangle \in \mathcal{H}_3$ , where  $\#1$  and  $\#2$  denote two arguments of the function. We assume that  $\mathcal{H}_1 = \mathcal{H}_2 = \mathcal{H}_X$ ,  $\mathcal{H}_3 = \mathcal{H}_X \otimes \mathcal{H}_Y$ .

**Assumption 2.** We assume that  $\widehat{\mu}_{(X, Y)}^+ = \sum_{i=1}^n \beta_i^+ \phi(x_i) \otimes \psi(y_i)$  is a consistent estimator of  $\mu_{(X, Y)}$  and  $\widehat{\mu}_X^+ = \sum_{i=1}^n \beta_i^+ \phi(x_i)$  is also consistent for the kernel embedding of the marginal distribution on  $X$ . Furthermore, we assume  $\sum_{i=1}^n \beta_i^+ \xrightarrow{p} 1$ . Note that as shown in Thm. 3, Thm. 7 and Corollary 2, this hypothesis holds when  $\mathcal{X}$  is compact and  $\mathcal{Y}$  is finite.

**Theorem 10.** With the above Assumption 1, Assumption 2 and Hypothesis 1, Hypothesis 2 in [20], we assert that if  $\lambda_n$  decreases to 0,

$$\mathcal{E}_s[\mu_{\mathbf{z}}^{\lambda_n}] - \mathcal{E}_s[\mu_{\mathcal{H}_K}] \rightarrow 0 \quad (22)$$

in probability as  $n \rightarrow \infty$ .

*Proof.* This proof is adapted from that of Thm. 5 in [20]. We split the proof to 3 steps.

**Step 1:** Given a training set  $\mathbf{z} = (\mathbf{x}, \mathbf{y}) \in \mathcal{Z}^n$ , Prop. 2 in [20] gives

$$\mathcal{E}_s[\mu_{\mathbf{z}}^\lambda] - \mathcal{E}_s[\mu_{\mathcal{H}_K}] = \left\| \sqrt{T}(\mu_{\mathbf{z}}^\lambda - \mu_{\mathcal{H}_K}) \right\|_{\mathcal{H}_K}^2.$$

As usual,

$$\mu_{\mathbf{z}}^\lambda - \mu_{\mathcal{H}_K} = (\mu_{\mathbf{z}}^\lambda - \mu^\lambda) + (\mu^\lambda - \mu_{\mathcal{H}_K})$$

Another application of Prop. 2 in [20] gives

$$\begin{aligned} \mu_{\mathbf{z}}^\lambda - \mu^\lambda &= (T_{\mathbf{x}} + \lambda)^{-1} A_{\mathbf{z}}^* \psi(\mathbf{y}) - (T + \lambda)^{-1} A^* \psi(y) \\ &= (T_{\mathbf{x}} + \lambda)^{-1} (A_{\mathbf{z}}^* \psi(\mathbf{y}) - T_{\mathbf{x}} \mu_{\mathcal{H}_K}) + (T_{\mathbf{x}} + \lambda)^{-1} (T - T_{\mathbf{x}}) (\mu^\lambda - \mu_{\mathcal{H}_K}). \end{aligned}$$

From  $\|\mu_1 + \mu_2 + \mu_3\|_{\mathcal{H}_K}^2 \leq 3(\|\mu_1\|_{\mathcal{H}_K}^2 + \|\mu_2\|_{\mathcal{H}_K}^2 + \|\mu_3\|_{\mathcal{H}_K}^2)$ ,

$$\mathcal{E}_s[\mu_{\mathbf{z}}^\lambda] - \mathcal{E}_s[\mu_{\mathcal{H}_K}] \leq 3(\mathcal{A}(\lambda) + \mathcal{S}_1(\lambda, \mathbf{z}) + \mathcal{S}_2(\lambda, \mathbf{z})), \quad (23)$$

where

$$\begin{aligned} \mathcal{S}_1(\lambda, \mathbf{z}) &= \left\| \sqrt{T} (T_{\mathbf{x}} + \lambda)^{-1} (A_{\mathbf{z}}^* \psi(\mathbf{y}) - T_{\mathbf{x}} \mu_{\mathcal{H}_K}) \right\|_{\mathcal{H}_K}^2 \\ \mathcal{S}_2(\lambda, \mathbf{z}) &= \left\| \sqrt{T} (T_{\mathbf{x}} + \lambda)^{-1} (T - T_{\mathbf{x}}) (\mu^\lambda - \mu_{\mathcal{H}_K}) \right\|_{\mathcal{H}_K}^2. \end{aligned}$$

**Step 2:** probabilistic bound on  $\mathcal{S}_2(\lambda, \mathbf{z})$ . First

$$\mathcal{S}_2(\lambda, \mathbf{z}) \leq \left\| \sqrt{T} (T_{\mathbf{x}} + \lambda)^{-1} \right\|_{\mathcal{L}(\mathcal{H}_K)}^2 \left\| (T - T_{\mathbf{x}}) (\mu^\lambda - \mu_{\mathcal{H}_K}) \right\|_{\mathcal{H}_K}^2. \quad (24)$$

**Step 2.1:** probabilistic bound on  $\left\| \sqrt{T} (T_{\mathbf{x}} + \lambda)^{-1} \right\|_{\mathcal{L}(\mathcal{H}_K)}$ . We introduce an auxiliary quantity

$$\Theta(\lambda, \mathbf{z}) = \left\| (T + \lambda)^{-1} (T - T_{\mathbf{x}}) \right\|_{\mathcal{L}(\mathcal{H}_K)}$$

and assume

$$\Theta(\lambda, \mathbf{z}) \leq \frac{1}{2}.$$

Invoking the Neumann series,

$$\begin{aligned} \left\| \sqrt{T} (T_{\mathbf{x}} + \lambda)^{-1} \right\|_{\mathcal{L}(\mathcal{H}_K)} &= \sqrt{T} (T + \lambda)^{-1} \sum_{n=0}^{\infty} ((T + \lambda)^{-1} (T - T_{\mathbf{x}}))^n \\ &\leq \left\| \sqrt{T} (T + \lambda)^{-1} \right\|_{\mathcal{L}(\mathcal{H}_K)} \sum_{n=0}^{\infty} \Theta(\lambda, \mathbf{z})^n \\ \text{(By spectral theorem)} &\leq \frac{1}{2\sqrt{\lambda}} \frac{1}{1 - \Theta(\lambda, \mathbf{z})} \leq \frac{1}{\sqrt{\lambda}} \end{aligned} \quad (25)$$

We now claim that  $\Theta(\lambda, \mathbf{z}) \leq \frac{1}{2}$  with high probability as  $n \rightarrow \infty$ . Let  $\xi_1 : \mathcal{X} \rightarrow \mathcal{L}_2(\mathcal{H}_K)$  be the random variable

$$\xi_1(x) = (T + \lambda)^{-1} T_x.$$

By the same reasoning in the proof of Thm. 5 in [20], we have  $\|\xi_1\|_{\mathcal{L}_2(\mathcal{H}_K)} \leq \frac{\kappa}{\lambda} = \frac{H_1}{2}$  and  $\mathbb{E}[\|\xi_1\|_{\mathcal{L}_2(\mathcal{H}_K)}^2] \leq \frac{\kappa}{\lambda} \mathcal{N}(\lambda) = \sigma_1^2$ . Our assumptions and Thm. 9 ensure that for any  $\delta_1$  there exists  $N_1(\delta_1)$  such that

$$\Theta(\lambda, \mathbf{z}) = \left\| (T + \lambda)^{-1} T_{\mathbf{x}} - (T + \lambda)^{-1} T \right\|_{\mathcal{L}_2(\mathcal{H}_K)} \leq \frac{1}{2}$$

with probability greater than  $1 - \delta_1$  as long as  $n > N_1(\delta_1)$ .

**Step 2.2:** probabilistic bound on  $\left\| (T - T_{\mathbf{x}}) (\mu^\lambda - \mu_{\mathcal{H}_K}) \right\|_{\mathcal{L}(\mathcal{H}_K)}$ . Let  $\xi_2 : \mathcal{X} \rightarrow \mathcal{H}_K$  be the random variable

$$\xi_2(x) = T_x (\mu^\lambda - \mu_{\mathcal{H}_K}).$$

By the same reasoning, we have  $\|\xi_2(x)\|_{\mathcal{H}_K} \leq \kappa\sqrt{\mathcal{B}(\lambda)} = \frac{H_2}{2}$  and  $\mathbb{E}[\|\xi_2\|_{\mathcal{H}_K}^2] \leq \kappa\mathcal{A}(\lambda) = \sigma_2^2$ . Applying our assumptions and Thm. 9 we conclude that for any  $\delta_2, \epsilon_2$  there exists  $N_2(\delta_2, \epsilon_2)$  such that

$$\|(T - T_{\mathbf{x}})(\mu^\lambda - \mu_{\mathcal{H}_K})\|_{\mathcal{H}_K} \leq \epsilon_2 \quad (26)$$

with probability greater than  $1 - \delta_2$  as long as  $n > N_2(\delta_2, \epsilon_2)$ .

**Step 3:** probabilistic bound on  $\mathcal{S}_1(\lambda, \mathbf{z})$ . As usual,

$$\mathcal{S}_1(\lambda, \mathbf{z}) \leq \left\| \sqrt{T}(T_{\mathbf{x}} + \lambda)^{-1}(T + \lambda)^{1/2} \right\|_{\mathcal{L}(\mathcal{H}_K)}^2 \left\| (T + \lambda)^{-1/2}(A_{\mathbf{z}}^*\psi(\mathbf{y}) - T_{\mathbf{x}}\mu_{\mathcal{H}_K}) \right\|_{\mathcal{H}_K}^2.$$

**Step 3.1:** bound  $\left\| \sqrt{T}(T_{\mathbf{x}} + \lambda)^{-1}(T + \lambda)^{1/2} \right\|_{\mathcal{L}(\mathcal{H}_K)}$ . Let

$$\Omega(\lambda, \mathbf{z}) = \left\| (T + \lambda)^{1/2}(T - T_{\mathbf{x}})(T + \lambda)^{-1/2} \right\|_{\mathcal{L}(\mathcal{H}_K)}$$

and assume  $\Omega(\lambda, \mathbf{z}) \leq \frac{1}{2}$ . Clearly,

$$\begin{aligned} & \left\| \sqrt{T}(T_{\mathbf{x}} + \lambda)^{-1}(T + \lambda)^{1/2} \right\|_{\mathcal{L}(\mathcal{H}_K)} \quad (27) \\ &= \left\| \sqrt{T}(T + \lambda)^{-1/2} \{I - (T + \lambda)^{1/2}(T - T_{\mathbf{x}})(T + \lambda)^{-1/2}\}^{-1} \right\|_{\mathcal{L}(\mathcal{H}_K)} \\ &\leq \left\| \sqrt{T}(T + \lambda)^{-1/2} \right\|_{\mathcal{L}(\mathcal{H}_K)} \sum_{i=1}^{\infty} \Omega(\lambda, \mathbf{z})^i \end{aligned}$$

$$\text{(By spectral theorem)} \leq \frac{1}{1 - \Omega(\lambda, \mathbf{z})} = 2. \quad (28)$$

On the other hand,

$$\begin{aligned} \Omega(\lambda, \mathbf{z})^2 &= \langle (T + \lambda)^{-1}(T - T_{\mathbf{x}}), ((T + \lambda)^{-1}(T - T_{\mathbf{x}}))^* \rangle_{\mathcal{L}_2(\mathcal{H}_K)} \\ &\leq \left\| (T + \lambda)^{-1}(T - T_{\mathbf{x}}) \right\|_{\mathcal{L}_2(\mathcal{H}_K)}^2 = \Theta(\lambda, \mathbf{z})^2. \end{aligned}$$

As a result, we have  $\Omega(\lambda, \mathbf{z}) \leq \frac{1}{2}$  with probability greater than  $1 - \delta_1$  as long as  $n > N_1(\delta_1)$ .

**Step 3.2:** probabilistic bound on  $\left\| (T + \lambda)^{-1/2}(A_{\mathbf{z}}^*\psi(\mathbf{y}) - T_{\mathbf{x}}\mu_{\mathcal{H}_K}) \right\|_{\mathcal{H}_K}$ . Let  $\xi_3 : \mathcal{Z} \rightarrow \mathcal{H}_K$  be the random variable

$$\xi_3(x, y) = (T + \lambda)^{-1/2}K_x(\psi(y) - \mu_{\mathcal{H}_K}(x)).$$

Via the same reasoning in the proof of Thm. 5 in [20], we have  $\|\xi_3\|_{\mathcal{H}_K} \leq \sqrt{\frac{\kappa M}{\lambda}} = \frac{H_3}{2}$  and  $\mathbb{E}[\|\xi_3\|_{\mathcal{H}_K}^2] \leq MN(\lambda) = \sigma_3^2$ . From our assumptions and Thm. 9 we know for each  $\epsilon_3$  and  $\delta_3$  there exists  $N_3(\delta_3, \epsilon_3)$  such that

$$\left\| (T + \lambda)^{-1/2}(A_{\mathbf{z}}^*\psi(\mathbf{y}) - T_{\mathbf{x}}\mu_{\mathcal{H}_K}) \right\|_{\mathcal{H}_K} \leq \epsilon_3 \quad (29)$$

with probability greater than  $1 - \delta_3$  as long as  $n > N_3(\delta_3, \epsilon_3)$ .

Linking bounds (23), (25), (26), (28), and (29) we obtain that for every  $\epsilon_1, \epsilon_2, \epsilon_3 > 0$  and  $\delta_1, \delta_2, \delta_3 > 0$  there exists  $N = \max\{N_1(\delta_1), N_2(\delta_2, \epsilon_2), N_3(\delta_3, \epsilon_3)\}$  such that for each  $n > N$ ,

$$\mathcal{E}_s[\mu_{\mathbf{z}}^\lambda] - \mathcal{E}_s[\mu_{\mathcal{H}_K}] \leq 3[\mathcal{A}(\lambda) + \frac{\epsilon_2^2}{\lambda} + 4\epsilon_3^2]$$

with probability greater than  $1 - \delta_1 - \delta_2 - \delta_3$ . This means that for any  $\epsilon > 0$  and fixed  $\lambda$

$$\lim_{n \rightarrow 0} p(\mathcal{E}_s[\mu_{\mathbf{z}}^\lambda] - \mathcal{E}_s[\mu_{\mathcal{H}_K}] > 3\mathcal{A}(\lambda) + \epsilon) = 0 \quad (30)$$

From [25] we know

$$\lim_{\lambda \rightarrow 0} \mathcal{A}(\lambda) = 0. \quad (31)$$

Combining (30) and (31) we can conclude that as long as  $\lambda$  decreases to 0,  $\mathcal{E}_s[\mu_{\mathbf{z}}^\lambda]$  converges to  $\mathcal{E}_s[\mu_{\mathcal{H}_K}]$  in probability.  $\square$

**Theorem 5.** Assume Hypothesis 1 and Hypothesis 2 in [20] and our Assumption 1 hold. With the conditions in Thm. 3, we assert that if  $\lambda_n$  decreases to 0 sufficiently slowly,

$$\mathcal{E}_s[\widehat{\mu}_{\lambda_n, n}] - \mathcal{E}_s[\mu'] \rightarrow 0 \tag{32}$$

in probability as  $n \rightarrow \infty$ .

*Proof.* The proof follows directly from Thm. 2, Thm. 3, Thm. 7, Corollary 2 and Thm. 10.  $\square$