# Kernel Bayesian Inference with Posterior Regularization

Yang Song[1]　　Jun Zhu[2]　　Yong Ren[2]

[1]Dept. of Physics, Tsinghua University　　[2]Dept. of CS and Tech., Tsinghua University

## 1. INTRODUCTION

**Observations:**

- Kernel methods can be used to embed a distribution to a Hilbert space and probability rules can be replaced by corresponding linear operators
- The kernel embedding of a conditional distribution has an optimizational formulation
- The posterior distribution in Bayes' rule has an optimizational formulation

**Does the kernel embedding of a posterior distribution have an optimizational formulation?**

**Contributions:**

- A theoretically justified affirmative answer to the question
- A simpler but faster regularization technique called thresholding regularization
- Posterior regularization for kernel Bayesian inference called kRegBayes, analogous to RegBayes

## 2. PRELIMINARIES

**Kernel embedding:**

$$p_X \mapsto \mathbb{E}_{p_X}[\phi(X)] =: \mu_X \in \mathcal{H}_\mathcal{X},$$

where $\phi(X) := k(X, \cdot)$. (i) When $p_X$ is a conditional distribution, $\mu_X$ is called *conditional embedding*. (ii) When $p_X$ is a posterior distribution in a Bayesian setting, $\mu_X$ is called *posterior embedding*.

**Optimizationl formulation of conditional embedding**

$$\mu_{Y|X} = \arg\inf_\mu \mathcal{E}_s[\mu] = \arg\inf_\mu \mathbb{E}_{(X,Y)}[\|\psi(Y) - \mu(X)\|^2_{\mathcal{H}_\mathcal{Y}}]$$

Given i.i.d. samples $\{(x_1, y_1), \cdots, (x_n, y_n)\}$ from $p(X, Y)$, the *estimator* is

$$\widehat{\mathcal{E}}_s[\mu] = \frac{1}{n}\|\psi(y_i) - \mu(x_i)\|^2_{\mathcal{H}_\mathcal{Y}}$$

**Optimizational formulation of posterior distribution**

$$p(Y \mid X = x) = \arg\min_{q(Y)} \left\{ \mathrm{KL}(q(Y)\|\pi(Y)) - \int \log p(X = x \mid Y)dq(Y) \right\}$$
$$s.t. \quad q(Y) \in \mathcal{P}_{\mathrm{prob}}$$

**Posterior regularization for Bayesian inferece (RegBayes)**

$$\min_{q(Y), \xi} \left\{ \mathrm{KL}(q(Y)\|\pi(Y)) - \int \log p(X = x \mid Y)dq(Y) + U(\xi) \right\}$$
$$s.t. \quad q(Y) \in \mathcal{P}_{\mathrm{prob}}(\xi)$$

## 3. POSTERIOR EMBEDDING AS A REGRESSOR

Let $\pi(Y)$ be the prior, $p(X \mid Y)$ be the likelihood, $p^\pi(X, Y)$ be the joint distribution and suppose we have samples to directly estimate $\pi(Y)$ and $p(X \mid Y)$. The posterior embedding $\mu^\pi_{Y|X}$ is the same as conditional embedding

$$\mu^\pi_{Y|X} = \arg\inf_\mu \mathcal{E}_s[\mu] = \arg\inf_\mu \mathbb{E}_{(X,Y)}[\|\psi(Y) - \mu(X)\|^2_{\mathcal{H}_\mathcal{Y}}]$$

**How to get a reasonable estimator of $\mathcal{E}_s$ when we do not have i.i.d. samples from $p^\pi(X, Y)$?**

Assuming $f(x, y) = \|\psi(y) - \mu(x)\|^2_{\mathcal{H}_\mathcal{Y}} \in \mathcal{H}_\mathcal{X} \otimes \mathcal{H}_\mathcal{Y}$, we have

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

We show in the paper that $\mu_{(X,Y)}$ can be estimated by $\sum_{i=1}^n \beta_i \psi(Y_i) \otimes \phi(X_i)$.

**Theorem 1** (Proof in Appendix). *Under some conditions (details in paper), we have the following consistent estimator of $\mathcal{E}_s[\mu]$:*

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

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

**What if some $\beta_i$'s are negative and $\widehat{\mathcal{E}}_s[\mu]$ has no minima?**

Under some conditions, $\widehat{\mathcal{E}}_s^+[\mu] = \sum_{i=1}^n \beta_i^+ \|\psi(y_i) - \mu(x_i)\|^2_{\mathcal{H}_\mathcal{Y}}$, where $\beta_i^+ = \max(0, \beta_i)$ is also consistent. This is called *thresholding regularization*.

Finally, we can establish the consistency of $\widehat{\mu}_{\lambda, n} = \arg\inf_\mu \widehat{\mathcal{E}}_{\lambda, n}[\mu]$, where

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

## 4. kREGBAYES

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

where $\{(x_i, y_i)\}_{i=1}^m$ is the sample used for *representing likelihood*, $\{(x_i, t_i)\}_{i=m+1}^n$ is the sample used for nonparametric *posterior regularization*. $\psi(t_i)$ is the kernel embedding of $\delta(Y = t_i)$ and encourages $p(Y \mid X = x_i)$ to be close to $\delta(Y = t_i)$.

## 5. EXPERIMENTS

We apply the framework to state-space filtering tasks, since Bayesian inference is a key element of filtering.

**Toy dynamics**

- We compare results of EKF, UKF, KBR (kernel Bayes' rule), pKBR (KBR with thresholding regularization) and kRegBayes
- The data points $\{(\theta_t, x_t, y_t)\}$ are generated from the dynamics

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

- Use samples from the true dynamics as regularization for kRegBayes.



Results for toy dynamics

**Camera position recovery**

- We compare results of KF, KBR, pKBR and kRegBayes
- The camera has a fixed height and is in a circular region with bounded radii.
- The dynamics is

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

- During training we choose $R_1 = 0$ and $R_2 = 10$ while during testing we use $R_1 = 5$ and $R_2 = 7$.
- We generate positions with radii 6 and use them as the regularization for kRegBayes.



First several training and testing frames for camera position recovery



MSEs for camera position recovery　　Probability histograms of distances