

"Hilbert Space Embedding of Conditional Distribution with Application to Dynamical Systems"

report by Mohsen Nabian

The main contributions of this paper is as follows:

- ① ^{the} paper introduced the concept of embedding conditional distributions in an RKHS (Reproducing kernel Hilbert Space) and presented a novel method for estimating such embeddings from training data.
- ② The paper consider several useful probabilistic inference operations such as marginalization and conditioning and show, using the author's theory that these operations can be performed only in RKHS.

(3) The article, applied their own inference algorithm to learn non-parametric models and perform inference for dynamical systems. These algorithms are general because it handles wide variety of possible nonlinear non gaussian models and also they apply in any setting in which an appropriate kernel function can be defined.

reproducing kernel Hilbert space (RKHS) F on X with kernel K is a Hilbert space of function $f: X \rightarrow \mathbb{R}$. Its dot product $\langle \cdot, \cdot \rangle_F$ satisfy the reproducing property:

$$\begin{aligned} \textcircled{1} \quad \langle f(\cdot), K(x, \cdot) \rangle_F &= f(x) \\ \Downarrow \\ \textcircled{2} \quad \langle K(x, \cdot), K(x', \cdot) \rangle_F &= K(x, x') \end{aligned}$$

also true map and its empirical estimates are

$$(a) \mu_x := E_x[\phi(x)] \quad E_x : \text{expectation}$$

$$(b) \hat{\mu}_x := \frac{1}{m} \sum_{i=1}^m \phi(x_i) \quad \text{feature map}$$

Def 1) When the mean map $\mu_x : P \rightarrow F$ is injective, the kernel function k is called characteristic

Th2) The empirical mean $\hat{\mu}_x$ converges to μ_x in the RKHS norm at rate of $O_p(\underbrace{m^{-1/2}}_{\substack{\text{size of training} \\ \text{set}}})$
 $D_x = \{x_1, \dots, x_m\}$

Cross-Covariance Operator:

$$C_{XY} : \mathcal{G} \rightarrow F$$

$$C_{XY} = E_{XY}[\phi(x) \otimes \phi(y)] - \mu_x \otimes \mu_y$$

\otimes \Rightarrow Tensor product

given two functions $f \in F$ and $g \in \mathcal{G}$ their cross-covariance

$$\text{Cov}_{xy}[f(X), g(Y)] := E_{XY}[f(X)g(Y)] - E_X[f(X)]E_Y[g(Y)] \text{ will be}$$

Computed as: $\langle f, C_{XY}g \rangle_F$ or $\langle f \otimes g, C_{XY} \rangle_{F \otimes g}$

For example:

Given ~~n~~ m pairs of training examples

$$D_{XY} = \{(x_1, y_1), \dots, (x_m, y_m)\}$$

iid from $P(X, Y)$,

$$Y = (\varphi(x_1), \dots, \varphi(x_m))$$

$$\Phi = (\phi(y_1), \dots, \phi(y_m))$$

Theorem 5: let $k_x := Y^T \varphi(x)$

then $\hat{\mu}_{Y|X}$ can be estimated as ~~$\hat{\mu}_{Y|X}$~~

$$\hat{\mu}_{Y|X} = \Phi (HK + \lambda m I)^{-1} H k_x = \sum_{i=1}^m \beta_i(x) \phi(y_i)$$

$\beta_i(x)$
 \downarrow
 real values
 $\in \mathbb{R}$

Theorem 6:

Assume $k(x, \cdot)$ is in range of C_{xx} .

The empirical conditional embedding $\hat{\mu}_{Y|X}$ converges to $\mu_{Y|X}$ in the RKHS norm.

Operations on RKHS Embeddings:

we have the following:

$$\begin{aligned}\mu_x &= E_Y [U_{x|Y} \phi(Y)] = U_{x|Y} E_Y [\phi(Y)] = \\ &= U_{x|Y} \mu_Y\end{aligned}$$

Chain Rule:

Similar to $P(X, Y) = P(X|Y) P(Y) = P(Y|X) P(X)$

we have:

$$\mu_{XY} = U_{X|Y} \mu_Y^{\otimes} = U_{Y|X} \mu_X^{\otimes}$$

here: $\mu_X^{\otimes} = E_X [\phi(X) \otimes \phi(X)]$ & $\mu_Y^{\otimes} = E_Y [\phi(Y) \otimes \phi(Y)]$

(App)

The paper also extract kernel formulation for Conditional cross validation.

Application: Dynamical Systems:

obj: Learning and inference in a dynamical system.

The article

~~models~~ models uncertainty in a dynamic system using partially observable Markov model which is:

$$P(s^1, \dots, s^T, o^1, \dots, o^T)$$

s^t is hidden state at time t .

\downarrow
 o^t corresponding observation

~~Learn~~
Theorem: Hilbert space prediction is given by:

$$\mu_{s^{t+1}|h} = U_{s^{t+1}} \mu_{s^t|s^t} \mu_{s^t|h}$$

Theorem: Hilbert space conditioning step is given by:

$$\mu_{s^{t+1}|h}$$

Learning Algorithm:

$$\text{input } \gamma = (\phi(s^t)) \quad \Phi = (\phi(s^{t+1}))$$

$$\Psi = (\psi(o^{t+1})) \quad t=1, 2, \dots, m$$

- ① Compute $K = \gamma^T \gamma$, $U = \Psi^T \Psi$, $G_\gamma = \gamma^T \Phi$
- ② Compute $T_2 = (HK + HU + \lambda m I)^{-1} H$
- ③ Compute $T_1 = T_2 G_\gamma$

Experiments: Experiment on Two dynamical systems

- ① a synthetic dataset generated from a linear dynamical ~~dataset~~ dynamical sys.
- ② Camera Tracking Problem

⑦ Synthetic data exp:
a particle rotating around the origin:

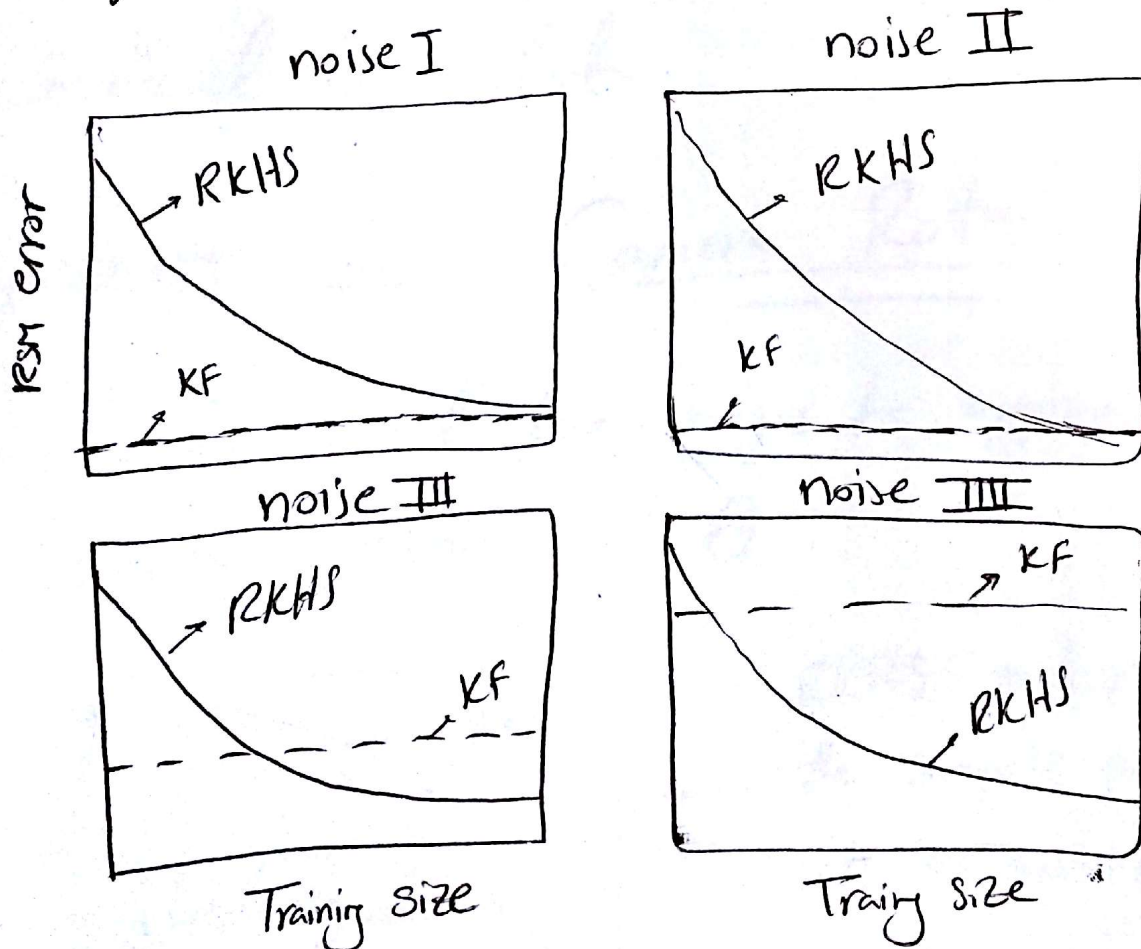
$$S^{t+1} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} S^t + \xi \rightarrow \text{Process noise}$$

$$O^{t+1} = S^{t+1} + \eta \rightarrow \text{observation noise}$$

$$\theta = 0.02$$

Goal: Compare the performance of ~~an~~ the article method with the Kalman Filter in estimating the position of particle.

Four Types of Process noise and Observation noise were tested. The result of KF & RKHS were compared:



←
nonlinear

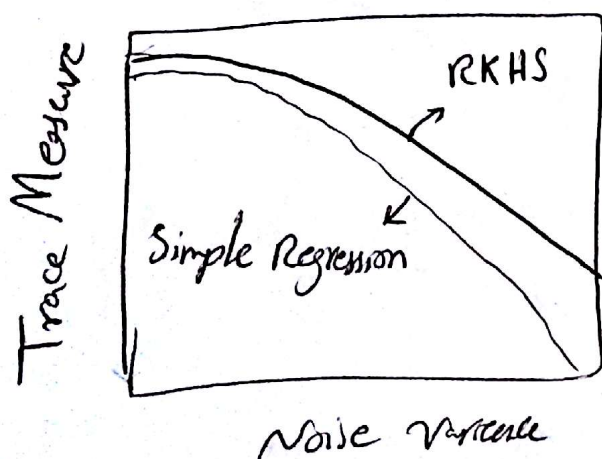
for

Result:

in cases when kalman filter is optimal (Case I, II)
the performance of the RKHS method approaches
the performance of KF with higher training data
However for nonlinear systems (Case III & IV)
~~RKHS~~ RKHS outperform KF significantly with
sufficient training data.

Experiment 2 Camera Rotation

used Cornell box images to approximate
camera rotation angle θ .



RKHS outperforms
the simple regression
in estimating camera
rotation