

Multi-Kernel Probability Distribution Regressions

Pingping Zhu[†], Hongchuan Wei[‡], Wenjie Lu[‡], and Silvia Ferrari[†]
 Laboratory for Intelligent Systems and Controls

[†] Department of Mechanical and Aerospace Engineering, Cornell University, Ithaca, NY, USA
[‡] Department of Mechanical Engineering and Materials Science, Duke University, Durham, NC, USA



Introduction and Motivation

This paper presents a multi-kernel probability distribution regression methodology that uses multi-layer reproducing kernel Hilbert space (RKHS) mappings to perform probability distribution to real and probability distribution to function regressions. The approach maps the distributions into RKHS by distribution embeddings. Then based on this RKHS, a multi-layer RKHS is constructed, within which the multi-kernel distribution regression can be implemented using an existing kernel regression algorithm, such as kernel recursive least squares (KRLS). The proposed algorithms are demonstrated through numerical simulations on synthetic data sets and compared with an existing algorithm. The results show that the proposed algorithm can outperform the existing algorithm.

Formulation of Probability Distribution Regression Problems

Distribution to Real Regression (DRR):

$$\mathbf{z}_k = \mathcal{F}(P_k) + \boldsymbol{\epsilon}_k, \quad k = 1, \dots, T \quad (1)$$

Distribution to Function Regression (DFR):

$$f_k = \mathcal{G}(P_k) + \boldsymbol{\epsilon}_k, \quad k = 1, \dots, T \quad (2)$$

where $P_k \in \mathcal{I}$ is a probability distribution defined in the probability space \mathcal{I} , $\mathbf{z}_k \in \mathbb{R}^{n_z}$ is the corresponding output response, and $\boldsymbol{\epsilon}_k$ is a zero mean Gaussian noise variable; $f_k \in \mathbb{F}$ is the function defined in a function space \mathbb{F} , $\boldsymbol{\epsilon}_k$ is a zero mean Gaussian process.

DRR operator $\mathcal{F}: \mathcal{I} \mapsto \mathbb{R}^{n_z}$ **DFR operator** $\mathcal{G}: \mathcal{I} \mapsto \mathbb{F}$

Training Data sets:

- Distribution P_k is approximated by using data set $\mathcal{X}_k = \{\mathbf{x}_{k,1}, \dots, \mathbf{x}_{k,s}, \dots, \mathbf{x}_{k,N_k}\}$
- DRR operator \mathcal{F} is approximated by using data set $\mathcal{D}_{\mathcal{F}} = \{(\mathcal{X}_1, \mathbf{z}_1), \dots, (\mathcal{X}_T, \mathbf{z}_T)\}$
- Output function f_k is approximated by using data set $\mathcal{Y}_k = \{(\mathbf{y}_{k,1}, \mathbf{z}_{k,1}), \dots, (\mathbf{y}_{k,M}, \mathbf{z}_{k,M})\}$ where $\mathbf{z}_{k,m} = f_k(\mathbf{y}_{k,m})$
- DFR operator \mathcal{G} is approximated by using data set $\mathcal{D}_{\mathcal{G}} = \{(\mathcal{X}_1, \mathcal{Y}_1), \dots, (\mathcal{X}_T, \mathcal{Y}_T)\}$.

Modified DFR operator:

According to (2)

$$\mathbf{z}_{k,m} = \mathcal{G}(P_k)(\mathbf{y}_{k,m}) + \boldsymbol{\epsilon}_k(\mathbf{y}_{k,m}) \quad (3)$$

Define new DFR operator $\mathcal{G}'(P_k, \mathbf{y}_{k,m}) = \mathcal{G}(P_k)(\mathbf{y}_{k,m})$, then

$$\mathbf{z}_{k,m} = \mathcal{G}'(P_k, \mathbf{y}_{k,m}) + \boldsymbol{\epsilon}'_k \quad (4)$$

where the evaluation of $\boldsymbol{\epsilon}(\mathbf{y}_{k,m})$ is denoted by $\boldsymbol{\epsilon}'_k$, which is a zero mean Gaussian noise.

Goal:

- Learn the operator \mathcal{F} by using data set $\mathcal{D}_{\mathcal{F}}$
- Learn the operator \mathcal{G}' by using data set $\mathcal{D}_{\mathcal{G}}$ instead of the operator \mathcal{G}

Multi-Kernel Distribution Regressions Methodology

Distribution Embeddings: Given a random variable (R.V.) $X \in \mathbb{R}^{n_x}$ associated with a distribution P_X and a corresponding Probability Density Function (PDF) p_X , an embedding $\boldsymbol{\mu}_X$ in RKHS can be defined as,

$$\boldsymbol{\mu}_X := \mathbf{E}_X[k_X(X, \cdot)] = \int p_X(\mathbf{x})k_X(\mathbf{x}, \cdot)d(\mathbf{x}), \quad (5)$$

where $\mathbf{E}_X[\cdot]$ indicates the expectation operator, $k_X(\cdot, \cdot)$ is a kernel defined on $\mathbb{R}^{n_x} \times \mathbb{R}^{n_x}$ associated with RKHS \mathcal{H}_X . The distribution embedding $\boldsymbol{\mu}_X$ is also in the RKHS \mathcal{H}_X , provided $\mathbf{E}_X[k_X(X, X)] < \infty$. Its empirical estimate is

$$\hat{\boldsymbol{\mu}}_X = \frac{1}{N} \sum_{n=1}^N k_X(\mathbf{x}_n, \cdot) \quad (6)$$

where $\mathcal{D}_X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is a training set that is assumed to have been drawn *i.i.d* from P_X . With a characteristic kernel, the mapping from distribution P_X to the distribution embedding $\boldsymbol{\mu}_X \in \mathcal{H}_X$ is *injective*. A famous characteristic kernel is the *Gaussian kernel*, which is used in this paper to specify the kernel function $k_X(\cdot, \cdot)$.

Kernel Design and Multi-Layer RKHS

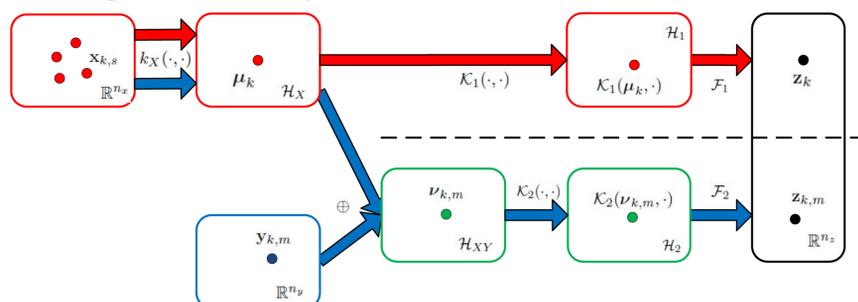


Figure 1: The frameworks of the multi-layer RKHS to implement MKDRR (red arrows) and MKDFR (blue arrows)

For DRR:

- Map the distribution data sets \mathcal{X}_i and \mathcal{X}_j , $i, j = 1, \dots, T$, into \mathcal{H}_X as $\boldsymbol{\mu}_i, \boldsymbol{\mu}_j$, respectively.
- Define new kernel \mathcal{K}_1 and the corresponding RKHS \mathcal{H}_1 ,

$$\mathcal{K}_1(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j) = \exp\left[-\frac{D(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j)}{2\sigma_{\boldsymbol{\mu}}^2}\right] \quad (7)$$

where $D(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j) = \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|_{\mathcal{H}_X}^2$

- Approximate the mapping \mathcal{F} using kernel methods in RKHS \mathcal{H}_1 with input signal $\boldsymbol{\mu}_k$ and desired signal $\mathbf{z}_k \in \mathbb{R}^{n_z}$, such as kernel adaptive filters.

For DFR:

- Map the distribution data sets \mathcal{X}_i and into \mathcal{H}_X as $\boldsymbol{\mu}_i$.
- Combine the distribution embedding $\boldsymbol{\mu}_i$ and $\mathbf{y}_{i,m}$ as a new term $\boldsymbol{\nu}_{i,m} = [\boldsymbol{\mu}_i^T, \mathbf{y}_{i,m}^T]^T \in \mathcal{H}_X \oplus \mathbb{R}^{n_y}$.
- Define new kernel \mathcal{K}_2 and the corresponding RKHS \mathcal{H}_2 ,

$$\mathcal{K}_2(\boldsymbol{\nu}_{i,m}, \boldsymbol{\nu}_{j,n}) = \exp\left[-\frac{(\boldsymbol{\nu}_{i,m} - \boldsymbol{\nu}_{j,n})^T \boldsymbol{\Sigma}_{XY}^{-1} (\boldsymbol{\nu}_{i,m} - \boldsymbol{\nu}_{j,n})}{2}\right] \quad (8)$$

- Approximate the mapping \mathcal{G} using kernel methods in RKHS \mathcal{H}_2 with input signal $\boldsymbol{\nu}_{k,m}$ and desired signal $\mathbf{z}_{k,m} \in \mathbb{R}^{n_z}$, such as kernel adaptive filters.

Multi-Kernel Distribution to Real Regression based on KRLS

Multi-Kernel Distribution to Real Regression based on KRLS

Like the standard KRLS algorithm, the following cost function is minimized to learn the DRR defined in (1) from data sets $\mathcal{D}_{\mathcal{F}}$,

$$J_{DRR} = \min_{\boldsymbol{\omega}_{\mathcal{F}}} \left[\sum_{k=1}^T \|\mathbf{z}_k - \langle \boldsymbol{\omega}_{\mathcal{F}}, \mathcal{K}_{\mathcal{F}}(\boldsymbol{\mu}_k, \cdot) \rangle_{\mathcal{H}_X}\|^2 + \lambda \|\boldsymbol{\omega}_{\mathcal{F}}\|_{\mathcal{H}_X}^2 \right] \quad (9)$$

where λ is a regularization factor. Then, the feature weight $\boldsymbol{\omega}_{\mathcal{F}}$ can be approximated at the k th iteration by

$$\boldsymbol{\omega}_{\mathcal{F}} = \boldsymbol{\Phi}_k [\boldsymbol{\Phi}_k^T \boldsymbol{\Phi}_k + \lambda \mathbf{I}_k]^{-1} \mathbf{Z}_k = \boldsymbol{\Phi}_k \mathbf{Q}_{\boldsymbol{\mu}}(k) \mathbf{Z}_k \quad (10)$$

where \mathbf{I}_k is a $k \times k$ identity matrix, the feature matrices $\boldsymbol{\Phi}_k = [\mathcal{K}_{\mathcal{F}}(\boldsymbol{\mu}_1, \cdot), \dots, \mathcal{K}_{\mathcal{F}}(\boldsymbol{\mu}_k, \cdot)]$ and desired matrices $\mathbf{Z}_k = [\mathbf{z}_1, \dots, \mathbf{z}_k]^T$, and

Multi-Kernel Distribution to Function Regression based on KRLS Similarly, the following cost function is minimized to learn the DFR defined in (4) from data sets $\mathcal{D}_{\mathcal{G}}$,

$$J_{DFR} = \min_{\boldsymbol{\omega}_{\mathcal{G}}} \left[\sum_{k=1}^T \sum_{m=1}^M \|\mathbf{z}_{k,m} - \mathcal{F}_{\mathcal{G}}(P_k, \mathbf{y}_{k,m})\|^2 + \lambda \|\boldsymbol{\omega}_{\mathcal{G}}\|_{\mathcal{H}_{\mathcal{G}}}^2 \right]. \quad (11)$$

we can approximate the feature weight $\boldsymbol{\omega}_{\mathcal{G}}$ at the k th by

$$\boldsymbol{\omega}_{\mathcal{G}} = \boldsymbol{\Upsilon}_k [\mathbf{K}_k + \lambda \mathbf{I}_{(kM)}]^{-1} \mathbf{V}_k = \boldsymbol{\Upsilon}_k \mathbf{Q}(\mathbf{V}_k) \mathbf{V}_k, \quad (12)$$

where input feature matrices $\boldsymbol{\Psi}_k = [\mathcal{K}_{\mathcal{G}}(\boldsymbol{\nu}_{k,1}, \cdot), \dots, \mathcal{K}_{\mathcal{G}}(\boldsymbol{\nu}_{k,M}, \cdot)]$ and $\boldsymbol{\Upsilon}_k = [\boldsymbol{\Psi}_1, \dots, \boldsymbol{\Psi}_k]$, and the desired matrices $\mathbf{U}_k = [\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,M}]^T$, $\mathbf{V}_k = [\mathbf{U}_1^T, \dots, \mathbf{U}_k^T]^T$, and $\mathbf{K}_k = \boldsymbol{\Upsilon}_k^T \boldsymbol{\Upsilon}_k$.

Simulations and Results

Simulation settings

- the initial agent location distributions P_k , $k = 1, 2, \dots$ on $\xi\eta$ -coordinate plane., generated by 2D Mixture Gaussian distributions with two equivalently weighted components
- The means $[\mu_{\xi,i}, \mu_{\eta,i}]$, $i = 1, 2$, and covariance matrices $\boldsymbol{\Sigma}_i = \text{diag}([\sigma_{\xi,i}, \sigma_{\eta,i}])$, $i = 1, 2$ are selected randomly,
- generate N_{train} training data sets, $N_{valid} = 25$ validation data sets, and $N_{test} = 50$ testing data sets.

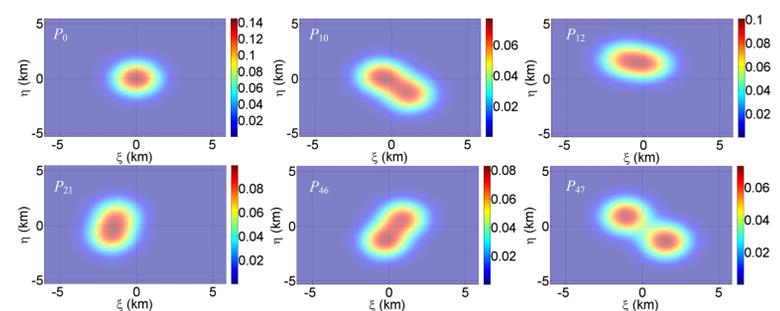


Figure 2: Examples of goal distribution and testing distributions generated by Gaussian mixture in DR problem.

Experiment of MKDRR-KRLS

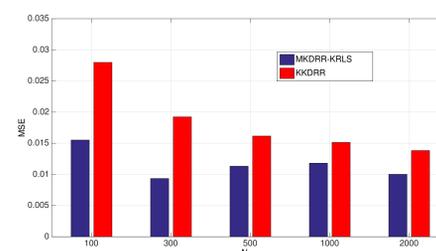


Figure 3: Performance comparison between MKDRR-KRLS and MKDFR-KRLS algorithms

$$D_{CS}(P_k || P_0) = -\log \frac{\int p_k(\xi, \eta) p_0(\xi, \eta) d\xi d\eta}{\sqrt{\int p_k^2(\xi, \eta) d\xi d\eta \int p_0^2(\xi, \eta) d\xi d\eta}} \quad (13)$$

Experiment of MKDFR-KRLS

Table 1: Regression performance results

Output functions	NMSE
CDF $F(\xi, \eta)$	0.0030 ± 0.0024
Gradient $g_{\xi}(\xi, \eta)$	0.0990 ± 0.0616
Gradient $g_{\eta}(\xi, \eta)$	0.0974 ± 0.0623

$$\text{NMSE} = \frac{\sum_{m=1}^M \|\mathbf{z}_{k,m} - \hat{\mathbf{z}}_{k,m}\|^2}{\sum_{m=1}^M \|\mathbf{z}_{k,m}\|^2} \quad (14)$$

- Each sample set \mathcal{X}_k has $N_{sample} = 500$ samples.
- $N_{train} = 500$ training distribution sets are used.

References and Acknowledgement

- Y. Engel, S. Mannor, and R. Meir, The kernel recursive least-squares algorithm, *IEEE Transactions on Signal Processing*, vol. 52, no. 8, pp. 2275–2285, Aug. 2004.
 - W. Liu, P. Pokharel, and J. C. Principe, The Kernel Least Mean Square Algorithm, *IEEE Transactions on Signal Processing*, vol. 56, no. 2, pp. 543–554, Feb. 2009.
 - B. Chen, S. Zhao, P. Zhu, and J. C. Principe, Quantized Kernel Recursive Least Square Algorithm, *IEEE Transactions on Neural Networks and Learning Systems*, vol. 23, no. 1, pp. 22–32, Jan. 2012.
 - B. Chen, S. Zhao, P. Zhu, and J. C. Principe, Quantized Kernel Recursive Least Squares Algorithm, *IEEE Transactions on Neural Networks and Learning Systems*, vol. 24, no. 9, pp. 1484–1491, Sep. 2013.
 - L. Song, J. Huang, A. Smola, and K. Fukumizu, Hilbert Space Embeddings of Conditional Distributions with Applications to Dynamical Systems, In *International Conference on Machine Learning (ICML)*, pp. 961A–968, 2009.
 - P. Zhu, B. Chen, and J. C. Principe, Learning Nonlinear Generative Models of Time Series with a Kalman Filter in RKHS, *IEEE Transactions on Signal Processing*, vol. 62, no. 1, pp. 141–155, Jan. 2014.
 - P. Zhu, Kalman Filtering in Reproducing Kernel Hilbert Spaces. Gainesville, FL, USA: PhD Thesis, University of Florida, 2013.
 - K. Rudd, G. Foderaro, and S. Ferrari, A Generalized Reduced Gradient Method for the Optimal Control of Multiscale Dynamical Systems, *Proc. of the IEEE Conference on Decision and Control (CDC)*, Dec. 2013.
 - K. Kampha, E. Hasanbelliu, and J. C. Principe, Closed-form cauchy-schwarz PDF divergence for mixture of Gaussians, *International Joint Conference on Neural Networks (IJCNN)*, pp. 2578–2585, 2011.
- This work was supported by NSF grant ECCS 1408022 and NFS grant DGE 1068871.