

Automated Spectral Kernel Learning

Jian Li, Yong Liu* and Weiping Wang
 {lijian9026, liuyong, wangweiping}@iie.ac.cn



中国科学院 信息工程研究所
 INSTITUTE OF INFORMATION ENGINEERING, CAS

Introduction

The generalization performance of kernel methods is largely determined by the kernel, but spectral representations of stationary kernels are both input-independent and output-independent, which limits their applications on complicated tasks.

In this paper, to achieve better performance ability for kernel methods, we propose an efficient algorithm, namely Automated Spectral Kernel Learning (ASKL), **learning suitable kernels and model weights together**.

1. Core Idea: **Spectral kernel representations + Rademacher complexity**.

2. Algorithmic contributions:

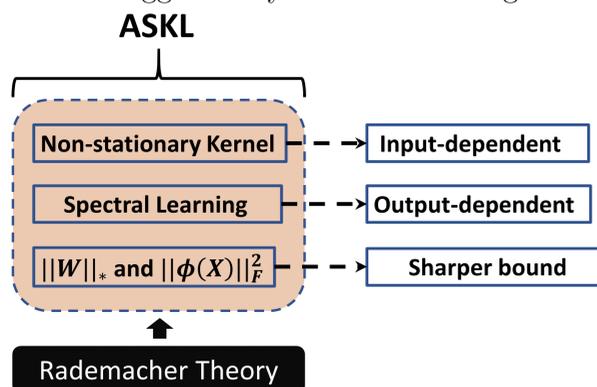
(1) non-stationary kernels to obtain **input-dependent features**

(2) backpropagation w.r.t the objective to make features **output-dependent**

(3) regularization terms to achieve **sharper generalization error bounds**

3. Theoretical contributions:

Based on **Rademacher complexity theory**, we explore how the feature mappings affect the performance and suggests ways to refine the algorithm.



Problem Definition

In ordinary supervised learning settings, training samples $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ are drawn i.i.d. from a fixed but unknown distribution ρ over $\mathcal{X} \times \mathcal{Y}$, where $\mathcal{X} = \mathbb{R}^d$ is the input space and $\mathcal{Y} \subseteq \mathbb{R}^K$ is the output space in single-valued ($K = 1$) or vector-valued ($K > 1$) forms. The goal is to learn an estimator $f: \mathcal{X} \rightarrow \mathcal{Y}$, which outputs K predictive labels. We define a standard hypothesis space for kernel methods

$$\mathcal{H} = \left\{ \mathbf{x} \rightarrow f(\mathbf{x}) = \mathbf{W}^T \phi(\mathbf{x}) \right\},$$

where $\mathbf{W} \in \mathbb{R}^{D \times K}$ is the model weight, $\phi(\mathbf{x}): \mathbb{R}^d \rightarrow \mathbb{R}^D$ is a nonlinear feature mapping. For kernel methods, $\phi(\mathbf{x})$ is an implicit feature mapping associated with a Mercer kernel $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$. To improve the computational efficiency but also retain favorable statistical properties, random Fourier features were proposed to approximate kernel with explicit feature mappings $\phi(\mathbf{x})$ via $k(\mathbf{x}, \mathbf{x}') \approx \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$.

In statistical learning theory, the supervised learning problem is to minimize the expected loss on $\mathcal{X} \times \mathcal{Y}$

$$\inf_{f \in \mathcal{H}} \mathcal{E}(f), \quad \mathcal{E}(f) = \int_{\mathcal{X} \times \mathcal{Y}} \ell(f(\mathbf{x}), \mathbf{y}) d\rho(\mathbf{x}, \mathbf{y}),$$

where ℓ is a loss function related to specific tasks.

Learning Framework

Based on Yaglom's theorem, the nonlinear feature mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$ adopts Monte Carlo approximation via inverse Fourier transform

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= \int_{\mathbb{R}^d \times \mathbb{R}^d} \mathcal{E}_{\omega, \omega'}(\mathbf{x}, \mathbf{x}') \mu(d\omega, d\omega') \\ &= \mathbb{E}_{\omega, \omega' \sim s} [\mathcal{E}_{\omega, \omega'}(\mathbf{x}, \mathbf{x}')] \\ &= \mathbb{E}_{\omega, \omega' \sim s} \left[\frac{1}{4} [\cos(\omega^T \mathbf{x} - \omega'^T \mathbf{x}') + \cos(\omega'^T \mathbf{x} - \omega^T \mathbf{x}') \right. \\ &\quad \left. + \cos(\omega^T \mathbf{x} - \omega^T \mathbf{x}') + \cos(\omega'^T \mathbf{x} - \omega'^T \mathbf{x}')] \right] \\ &\approx \frac{1}{4D} \sum_{i=1}^D [\cos(\omega_i^T \mathbf{x} - \omega_i'^T \mathbf{x}') + \cos(\omega_i'^T \mathbf{x} - \omega_i^T \mathbf{x}') \\ &\quad + \cos(\omega_i^T \mathbf{x} - \omega_i^T \mathbf{x}') + \cos(\omega_i'^T \mathbf{x} - \omega_i'^T \mathbf{x}')] \\ &= \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle \end{aligned}$$

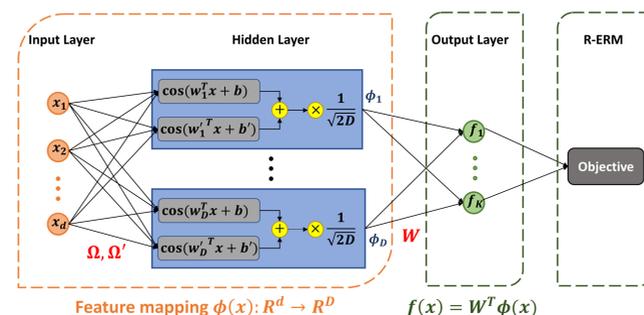
where $(\omega_i, \omega_i')_{i=1}^D \stackrel{\text{i.i.d.}}{\sim} s(\omega, \omega')$, the phase vectors \mathbf{b}, \mathbf{b}' are drawn uniformly from $[0, 2\pi]^D$ and

$$\phi(\mathbf{x}) = \frac{1}{\sqrt{2D}} [\cos(\Omega^T \mathbf{x} + \mathbf{b}) + \cos(\Omega'^T \mathbf{x} + \mathbf{b}')].$$

The trace norm $\|\mathbf{W}\|_*$ and the squared Frobenius norm $\|\phi(\mathbf{X})\|_F^2$ exerts constraints on updating model weights \mathbf{W} and frequency matrices Ω, Ω' . In this paper, we put two additional regularization terms into the ERM

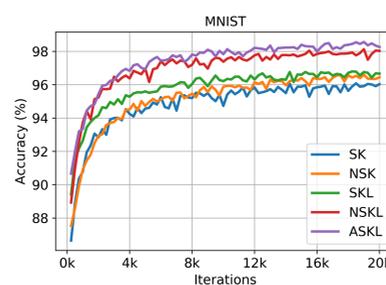
$$\arg \min_{\mathbf{W}, \Omega, \Omega'} \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), \mathbf{y}_i) + \lambda_1 \|\mathbf{W}\|_* + \lambda_2 \|\phi(\mathbf{X})\|_F^2,$$

where both feature mappings $\phi(\mathbf{X}) \in \mathbb{R}^{D \times n}$ on all data and $f(\mathbf{x}_i) = \mathbf{W}^T \phi(\mathbf{x}_i) \in \mathbb{R}^K$ use the non-stationary spectral representation.



Experimental Results

		SK	NSK	SKL	NSKL	ASKL
Accuracy (\uparrow)	segment	89.93 \pm 2.12	90.15 \pm 2.08	94.58 \pm 1.86	94.37 \pm 0.81	95.02\pm1.54
	satimage	74.54 \pm 1.35	75.15 \pm 1.38	83.61 \pm 1.08	83.74 \pm 1.34	85.32\pm1.45
	USPS	93.19 \pm 2.84	93.81 \pm 2.13	95.13 \pm 0.91	95.27 \pm 1.65	97.76\pm1.14
	pendigits	96.93 \pm 1.53	97.39 \pm 1.41	98.19 \pm 2.30	98.28 \pm 1.68	99.06\pm1.26
	letter	76.50 \pm 1.21	78.21 \pm 1.56	93.60 \pm 1.14	94.66 \pm 2.21	95.70\pm1.74
	porker	49.80 \pm 2.11	51.85 \pm 0.97	54.27 \pm 2.72	54.69 \pm 1.68	54.85\pm1.28
	shuttle	98.17 \pm 2.81	98.21 \pm 1.46	98.87 \pm 1.42	98.74 \pm 1.07	98.98\pm0.94
	MNIST	96.03 \pm 2.21	96.45 \pm 2.16	96.67 \pm 1.61	98.03 \pm 1.16	98.26\pm1.78
	abalone	10.09 \pm 0.42	9.71 \pm 0.28	8.35 \pm 0.28	7.85\pm0.42	7.88 \pm 0.16
	space_ga	11.86 \pm 0.26	11.58 \pm 0.42	11.40 \pm 0.18	11.39 \pm 0.46	11.34\pm0.27
RMSE (\downarrow)	cpusmall	2.77 \pm 0.71	2.84 \pm 0.38	2.56 \pm 0.72	2.57 \pm 0.63	2.42\pm0.48
	cadata	50.31 \pm 0.92	51.47 \pm 0.32	47.67 \pm 0.33	47.71 \pm 0.30	46.34\pm0.23



Generalization Analysis

Definition 1. The empirical Rademacher complexity of hypothesis space \mathcal{H} is

$$\widehat{\mathcal{R}}(\mathcal{H}) = \frac{1}{n} \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{H}} \sum_{i=1}^n \sum_{k=1}^K \epsilon_{ik} f_k(\mathbf{x}_i) \right],$$

where $f_k(\mathbf{x}_i)$ is the k -th value of the estimator $f(\mathbf{x}_i)$ with K outputs and ϵ_{ik} s are $n \times K$ independent Rademacher variables.

Theorem 1 (Excess Risk Bound). Assume that $B = \sup_{f \in \mathcal{H}} \|\mathbf{W}\|_* < \infty$ and assume the loss function ℓ is L -Lipschitz for \mathbb{R}^K , with probability at least $1 - \delta$, the excess risk bound holds

$$\mathcal{E}(\widehat{f}_n) - \mathcal{E}(f^*) \leq 4\sqrt{2}L\widehat{\mathcal{R}}(\mathcal{H}) + \mathcal{O}\left(\sqrt{\frac{\log 1/\delta}{n}}\right),$$

where $f^* \in \mathcal{H}$ is the most accurate estimator in the hypothesis space, \widehat{f}_n is the empirical estimator and the empirical Rademacher complexity is

$$\begin{aligned} \widehat{\mathcal{R}}(\mathcal{H}) &\leq \frac{B}{n} \sqrt{K \sum_{i=1}^n \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_i) \rangle} \\ &= \frac{B}{n} \sqrt{\frac{K}{D} \sum_{i=1}^n \sum_{j=1}^D \frac{1}{2} [\cos((\omega_j - \omega_j')^T \mathbf{x}_i) + 1]}. \end{aligned}$$

1. **The Influence of Non-Stationary Kernels.**

$k(\mathbf{x}_i, \mathbf{x}_i) = \cos(\omega^T(\mathbf{x}_i - \mathbf{x}_i)) = 1$, thus the trace of kernel matrix $\sum_{i=1}^n k(\mathbf{x}_i, \mathbf{x}_i) = n$, which corresponds to the worst cases. While for non-stationary kernels, $k(\mathbf{x}_i, \mathbf{x}_i) = \cos((\omega - \omega')^T \mathbf{x}_i) \in [-1, 1]$.

2. **Minimize the Trace Norm $\|\mathbf{W}\|_*$.** The convergence rate $B = \sup_{f \in \mathcal{H}} \|\mathbf{W}\|_* < \infty$ is also dependent on a constant B , that is the supremum of trace norm $\|\mathbf{W}\|_*$ in terms of the specific hypothesis space.

3. **Minimize the Squared Frobenius Norm $\|\phi(\mathbf{X})\|_F^2$.** Rademacher complexity is bounded by the trace of the kernel

$$\sum_{i=1}^n \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_i) \rangle = \sum_{i=1}^n \|\phi(\mathbf{x}_i)\|_2^2 = \|\phi(\mathbf{X})\|_F^2.$$