

# High-Dimensional Kernel Methods under Covariate Shift: Data-Dependent Implicit Regularization

Yihang Chen<sup>1</sup> Fanghui Liu<sup>2</sup> Taiji Suzuki<sup>3</sup> Volkan Cevher<sup>1</sup>  
<sup>1</sup>EPFL <sup>2</sup>University of Warwick <sup>3</sup>The University of Tokyo



## Overview

In this paper, we provide an initial analysis to the following question:

*How does IW affect bias-variance trade-off in high-capacity models?*

To summarize our contributions:

- We present the asymptotic expansion of high-dimensional kernels  $K(\mathbf{x}, \mathbf{x}')$  under covariate shifts, where the nonlinearity in kernels can be eliminated by the kernel function curvature.
- For variance, we demonstrate that the IW strategy can be regarded as an implicit data-dependent regularization on the respective kernel.
- For bias, we demonstrate two cases: 1) near interpolation, and 2) some proper regularization parameter.

## Kernel: Asymptotic Expansion

**Lemma 1.** Assume the kernel  $K$  is the inner-product kernel,  $K(\mathbf{x}, \mathbf{x}') := h(\langle \mathbf{x}, \mathbf{x}' \rangle / d)$ , or the radial kernel,  $K(\mathbf{x}, \mathbf{x}') := h(-\|\mathbf{x} - \mathbf{x}'\|_2^2 / d)$ , and the training data  $\mathbf{X} \sim p$ .

(1) Under suitable assumptions, we have  $\|\mathbf{K}(\mathbf{X}, \mathbf{X}) - \mathbf{K}^{\text{lin}}(\mathbf{X}, \mathbf{X})\|_2 \rightarrow 0$ , as  $n, d \rightarrow \infty, n/d \rightarrow \zeta$ , where  $\mathbf{K}^{\text{lin}}(\mathbf{X}, \mathbf{X})$  is defined by  $\mathbf{K}^{\text{lin}}(\mathbf{X}, \mathbf{X}) := \alpha_p \mathbf{1}\mathbf{1}^\top + \beta_p \frac{\mathbf{X}\mathbf{X}^\top}{d} + \gamma_p \mathbf{I} + \mathbf{T}_p$ , with non-negative parameters  $\alpha_p, \beta_p, \gamma_p$ , and the additional matrix  $\mathbf{T}_p$  given in Table 1.

(2) Under suitable assumptions, with  $c_{pq} < 2\theta_q - 1/2$ , with the training data  $\mathbf{X} \sim p$  and a test data  $\mathbf{x} \sim q$ , we have  $\mathbb{E}_q \|\mathbf{K}(\mathbf{X}, \mathbf{x}) - \mathbf{K}^{\text{lin}}(\mathbf{X}, \mathbf{x})\|_2 \rightarrow 0$ , as  $n, d \rightarrow \infty, n/d \rightarrow \zeta$ , where  $\mathbf{K}^{\text{lin}}(\mathbf{X}, \mathbf{x})$  is defined by  $\mathbf{K}^{\text{lin}}(\mathbf{X}, \mathbf{x}) := \beta_{pq} \frac{\mathbf{X}\mathbf{x}}{d} + \mathbf{T}_{pq}(\mathbf{X}, \mathbf{x})$ , with non-negative parameters  $\beta_{pq}$ , and the additional vector  $\mathbf{T}_{pq}$  given in Table 1.

Table 1: Parameters of the linearized kernel  $\mathbf{K}^{\text{lin}}$  involved with the curvature of  $h$ , when  $\mathbf{X} \sim p$ .

Parameters	Inner-Product Kernels	Radial Kernels
$\alpha_p$	$h(0) + h''(0) \frac{\text{Tr}(\Sigma_p^2)}{2d^2}$	$h(-2\tau_p) + 2h''(-2\tau_p) \frac{\text{Tr}(\Sigma_p^2)}{d^2}$
$\beta_p$	$h'(0)$	$2h'(-2\tau_p)$
$\gamma_p$	$h(\tau_p) - h(0) - \tau_p h'(0)$	$h(0) - 2\tau_p h'(-2\tau_p) - h(-2\tau_p)$
$\mathbf{T}_p$	$\mathbf{0}_{n \times n}$	$-h'(-2\tau_p) \mathbf{A} + \frac{1}{2} h''(-2\tau_p) \mathbf{A} \odot \mathbf{A}^1$
$\beta_{pq}$	$h'(0)$	$2h'(-(\tau_p + \tau_q))$
$\mathbf{T}_{pq}$	$\mathbf{0}_{n \times 1}$	$-h(-(\tau_p + \tau_q)) \cdot \mathbf{1} - \frac{\beta_{pq}}{2} \mathbf{A}(\mathbf{X}, \mathbf{x})^2$

<sup>1</sup>  $\mathbf{A} := \mathbf{1}\mathbf{\psi}^\top + \mathbf{\psi}\mathbf{1}^\top$ , where  $\mathbf{\psi} \in \mathbb{R}^n$  with  $\psi_i := \|\mathbf{x}_i\|_2^2 / d - \tau_p$ .

<sup>2</sup>  $\mathbf{A}(\mathbf{X}, \mathbf{x}) := \mathbf{\psi}_x + \mathbf{\psi}$ , where  $\psi_x = \|\mathbf{x}\|_2^2 / d - \tau_q$ .

## Problem Setting

**Notations:**

- **Data.**
  - Training distribution:  $p$ . Test distribution:  $q$ .
  - Re-weighting distribution  $\bar{q}$ . Re-weighting function  $\bar{w}(\mathbf{x}) = d\bar{q}(\mathbf{x})/dp(\mathbf{x})$ .
  - The label  $y$  is generated by  $f_\rho$ ,  $y(\mathbf{x}) = f_\rho(\mathbf{x}) + \varepsilon$ , and  $\mathbb{E}[\varepsilon] = 0$ ,  $\mathbb{V}[\varepsilon] \leq \sigma_\varepsilon^2$ .
- **Kernel.** The reproducing kernel Hilbert space (RKHS)  $\mathcal{H}$  is a Hilbert space  $\mathcal{H}$  endowed with the inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  of functions  $f: \mathcal{X} \rightarrow \mathbb{R}$  with a reproducing kernel  $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  where  $K(\cdot) \in \mathcal{H}$  and  $f(\mathbf{x}) = \langle f, K(\mathbf{x}, \cdot) \rangle_{\mathcal{H}}$ . Specifically, we consider the inner-product kernels,  $K(\mathbf{x}, \mathbf{x}') := h(\langle \mathbf{x}, \mathbf{x}' \rangle / d)$ .
- **Task:** Given  $n$  training data  $\mathbf{Z} = \{(\mathbf{x}_i, y_i) \sim p\}_{i=1}^n$ , the estimator of KRR in high dimensions under a general IW function  $\bar{w}(\mathbf{x})$  is given by  $\bar{f}_{\lambda, \mathbf{Z}} := \arg \min_{f \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^n \bar{w}(\mathbf{x}_i) (f(\mathbf{x}_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\}$ , where  $\lambda > 0$  is the regularization parameter.

**Assumptions (abbreviated):**

- Data.** Let  $\Sigma_p, \Sigma_q$  be the covariance matrix of the distribution  $p, q$ .
- **(Sub-Gaussian)** Let  $\mathbf{x} \sim p$ , then  $\Sigma_p^{-1} \mathbf{x}$  is sub-Gaussian with identity variance and bounded items, and similarly for  $q$ .
  - **(Covariance)** We assume  $\max\{\|\Sigma_p\|, \|\Sigma_q\|\} = O(1)$ . Define  $\Sigma_{pq} := \Sigma_p^{-1} \Sigma_q$ , and  $\exists c_{pq} \geq 0$  so that  $\text{Tr}(\Sigma_{pq})/d \lesssim d^{c_{pq}}$ . To limit the distribution shifts, we additionally assume  $c_{pq} < 2\theta_q - \frac{1}{2} = \frac{1}{2} - \frac{4}{8+m_q}$ .
  - **(Ratio)** The ratio  $w := dq/dp, \bar{w} := d\bar{q}/dp$ 's norm in some space is upper bounded by constants dependent on the dimension  $d$ .
- Model:**
- **(Source condition):** We have  $f_\rho \in \mathcal{H}$ , and there exists  $\frac{1}{2} \leq \bar{r} < 1, \bar{g}_\rho \in \mathcal{L}_q^2$  such that  $f_\rho = (L_{\bar{q}})^{\bar{r}} \bar{g}_\rho$ . We additionally assume  $\max\{\|f_\rho\|_{\mathcal{H}}, \|\bar{g}_\rho\|_q, \|f_\rho\|_\infty\} \lesssim d^{c_{\mathcal{H}}}$ .
  - **(Capacity condition):** For any  $\lambda > 0$ , there exists  $E_\mu > 0$  and  $s_\mu \in [0, 1]$  such that for distribution  $\mu \in \{q, \bar{q}\}$ ,  $\mathcal{N}_\mu(\lambda) := \text{Tr}((L_\mu + \lambda)^{-1} L_\mu) \leq E_\mu^2 \lambda^{-s_\mu}, \forall \lambda \in (0, 1]$ .

## Main Results

**Bias-variance decomposition:** We have the following bias-variance decomposition:

$$\mathbb{E}_{\mathbf{y}|\mathbf{X}} \|\bar{f}_{\lambda, \mathbf{Z}} - f_\rho\|_q^2 = \mathbb{E}_{\mathbf{y}|\mathbf{X}} \|\bar{f}_{\lambda, \mathbf{Z}} - \bar{f}_{\lambda, \mathbf{X}}\|_q^2 + \|\bar{f}_{\lambda, \mathbf{X}} - f_\rho\|_q^2 := \mathbf{V} + \mathbf{B}^2.$$

where  $\bar{f}_{\lambda, \mathbf{X}} := \arg \min_{f \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^n \bar{w}(\mathbf{x}_i) (f(\mathbf{x}_i) - f_\rho(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\}$ .

**Variance:**

**Theorem 2** (Data-dependent regularization). Let  $\delta \in (0, 1)$ , then for large  $d$ , with probability at least  $1 - \delta - 2d^{-2}$  with respect to a draw of  $\mathbf{X} \sim p$  and  $\varepsilon > 0$ , the variance can be estimated by

$$\mathbf{V} \leq \frac{8\sigma_\varepsilon^2 \|\Sigma_q\|}{d} \underbrace{\mathcal{N} \left( \frac{\mathbf{X}\mathbf{X}^\top}{d} + \frac{\lambda n \bar{\mathbf{W}}(\mathbf{X})^{-1}}{\beta_p}; \frac{\gamma_p}{\beta_p} \right)}_{\text{dominated term } \mathbf{V}_x} + \frac{8\sigma_\varepsilon^2}{\gamma_p^2} d^{-(4\theta_q - 1 - 2c_{pq})} \log^{4(1+\varepsilon)} d.$$

The dominated term in Theorem 2 can be represented as

$$\mathbf{V}_x \asymp \frac{1}{d} \mathcal{N} \left( \frac{\mathbf{X}\mathbf{X}^\top}{d} + \frac{\lambda n \bar{\mathbf{W}}(\mathbf{X})^{-1}}{\beta_p}; \frac{\gamma_p}{\beta_p} \right),$$

which implies that the variance is well controlled by the capacity of  $\mathbf{K}^{\text{lin}} + \lambda n \bar{\mathbf{W}}^{-1}$ .

**Bias:**

**Theorem 3** (Bias under arbitrary  $\lambda$ ). Let  $\delta \in (0, 1)$ , we have the bias  $\mathbf{B}$  is upper bounded as  $\mathbf{B} \leq \mathbf{B}_{\text{in}} + \mathbf{B}_{\text{iw}}$ , where  $\mathbf{B}_{\text{in}}$  is the intrinsic bias that only depends on the problem of covariate shift from  $p$  to  $q$  via the ratio  $w(\mathbf{x})$ , and  $\mathbf{B}_{\text{in}} := \text{Tr}(\mathbf{K}^{\text{lin}} \mathbf{W}) / n$ . The second term is the re-weighting bias  $\mathbf{B}_{\text{iw}}$  that depends on the choice of  $\bar{w}(\mathbf{x})$ ,  $w(\mathbf{x})$ , and  $\lambda$ . When  $w = \bar{w}$ , we have  $\mathbf{B}_{\text{iw}} := \lambda^2 n \mathcal{N}(\mathbf{K}^{\text{lin}} \mathbf{W}, n\lambda) + o(1)$ , with probability at least  $1 - 4\delta$  for sufficiently large  $d$ .

**Theorem 4** (Bias under some  $\lambda$ ). Under some assumptions on the data and model, with proper selection of  $c_\lambda$  and  $C_\lambda$ , when choosing  $\lambda := C_\lambda n^{-c_\lambda}$ , then with probability at least  $1 - \delta$ , for sufficiently large  $d$ , when  $c_{\mathcal{H}} < \bar{r} c_\lambda$ , it holds that

$$\mathbf{B} \lesssim n^{-\bar{r}c_\lambda + c_{\mathcal{H}}} \|L_q(L_{\bar{q}} + \lambda)^{-1}\|^{1/2}.$$

For general  $\lambda$ , we have, with  $\lesssim$  here hiding the dependence on  $n$ ,

$$\mathbf{B} \lesssim (\lambda^{\bar{r}} + \lambda^{-\frac{1}{2}}) \|L_q(L_{\bar{q}} + \lambda)^{-1}\|^{1/2}.$$

**Related works:**

1. Liang, T., Rakhlin, A. (2020). Just interpolate: Kernel "ridgeless" regression can generalize.
2. Liu, F., Liao, Z., Suykens, J. (2021, March). Kernel regression in high dimensions: Refined analysis beyond double descent. In International Conference on Artificial Intelligence and Statistics (pp. 649-657). PMLR.