

This is the supplementary material for the paper: “Kernel regression in high dimensions: Refined analysis beyond double descent”, by Fanghui Liu, Zhenyu Liao, and Johan A.K. Suykens. The supplementary material (Appendix) is organized as follows.

- Section A provides high dimensional linearizations of some typical smooth kernels as concrete examples of Table 2.
- In Section B, we demonstrate that, a kernel matrix in high dimensions admits the same eigenvalue decay as $\widetilde{\mathbf{X}}$ and $\mathbf{X}\mathbf{X}^\top/d$.
- Our proof framework includes the error decomposition in Section C, the error bound for the bias in Section D and for the variance in Section E, respectively.
- Section F discusses the quantity function $\mathcal{N}^{n\lambda+\gamma}$ based on three eigenvalue decays: *harmonic decay*, *polynomial decay*, and *exponential decay* in the $n < d$ and $n > d$ regimes.
- Some additional experiments are presented in Section G to further validate our theoretical results.

A Examples of kernels and their linearizations

In this section, we present linearization of some typical kernels by Eq. (5). Here we assume that $\alpha, \beta, \gamma \geq 0$ to ensure the positive definiteness of the approximated kernel matrix $\widetilde{\mathbf{K}}^{\text{lin}}$. Table 4 reports the results of three inner-product kernels including polynomial kernel, linear kernel, exponential kernel; as well as a radial kernel: the common-used Gaussian kernel. We can find that $\alpha, \gamma \geq 0$. Specifically, $\beta > 0$ avoids a trivial solution.

Table 4: Linearizations of typical kernels in high dimensions.

kernel	formulation	α	β	γ
polynomial kernels	$k(\mathbf{x}, \mathbf{x}') := (1 + \frac{1}{d} \langle \mathbf{x}, \mathbf{x}' \rangle)^p$	$1 + p(p-1) \frac{\text{tr}(\boldsymbol{\Sigma}_d^2)}{2d^2}$	p	$(1 + \tau)^p - 1 - p\tau$
linear kernel	$k(\mathbf{x}, \mathbf{x}') = \frac{1}{d} \langle \mathbf{x}, \mathbf{x}' \rangle$	0	1	0
exponential kernel	$k(\mathbf{x}, \mathbf{x}') = \exp(\frac{2}{d} \langle \mathbf{x}, \mathbf{x}' \rangle)$	$1 + 2 \frac{\text{tr}(\boldsymbol{\Sigma}_d^2)}{d^2}$	2	$\exp(2\tau) - 1 - 2\tau$
Gaussian kernel	$k(\mathbf{x}, \mathbf{x}') = \exp(-\frac{1}{d} \ \mathbf{x} - \mathbf{x}'\ _2^2)$	$\exp(-2\tau) \left[1 + 2 \frac{\text{tr}(\boldsymbol{\Sigma}_d^2)}{d^2} \right]$	$2 \exp(-2\tau)$	$1 - 2\tau \exp(-2\tau) - \exp(-2\tau)$

B Eigenvalue decay equivalence

In this section, we demonstrate that, in high dimensions, a kernel matrix induced by inner-product kernels or radial kernels admits the same eigenvalue decay as $\widetilde{\mathbf{X}} = \beta \mathbf{X}\mathbf{X}^\top/d + \alpha \mathbf{1}\mathbf{1}^\top$ and $\mathbf{X}\mathbf{X}^\top/d$.

For notational simplicity, denote the inner-product kernel matrix $\mathbf{K}_{\text{inner}}$ and its linearization $\widetilde{\mathbf{K}}_{\text{inner}}^{\text{lin}}$; the radial kernel matrix $\mathbf{K}_{\text{radial}}$ and its linearization $\widetilde{\mathbf{K}}_{\text{radial}}^{\text{lin}}$.

Proposition 2. *The inner-product kernel matrix $\mathbf{K}_{\text{inner}}$ admits the same eigenvalue decay as $\widetilde{\mathbf{X}}$ and $\mathbf{X}\mathbf{X}^\top/d$.*

Proof. According to Theorem 2.1 in [26], the inner-product kernel matrix $\mathbf{K}_{\text{inner}}$ can be well approximated by $\widetilde{\mathbf{K}}_{\text{inner}}^{\text{lin}}$ with

$$\widetilde{\mathbf{K}}_{\text{inner}}^{\text{lin}} := \beta \frac{\mathbf{X}\mathbf{X}^\top}{d} + \gamma \mathbf{I} + \alpha \mathbf{1}\mathbf{1}^\top,$$

in a spectral norm sense, where α, β, γ are given in Table 2. As a result, with high probability, the inner-product kernel matrix $\mathbf{K}_{\text{inner}}$ and its linearization $\widetilde{\mathbf{K}}_{\text{inner}}^{\text{lin}}$ has the same eigenvalue. That means, $\mathbf{K}_{\text{inner}}$ admits the same eigenvalue decay as $\widetilde{\mathbf{X}} := \beta \mathbf{X}\mathbf{X}^\top/d + \alpha \mathbf{1}\mathbf{1}^\top$ via a constant shift γ .

Next, we shall demonstrate that $\mathbf{K}_{\text{inner}}$ admits the same eigenvalue decay as $\mathbf{X}\mathbf{X}^\top/d$. Since $\mathbf{1}\mathbf{1}^\top$ is a rank-one matrix with $\lambda_1(\mathbf{1}\mathbf{1}^\top) = n$, with Weyl's inequality and $\lambda_n \leq \lambda_{n-1} \leq \dots \leq \lambda_1$, we have

$$\beta\lambda_1 \left(\frac{\mathbf{X}\mathbf{X}^\top}{d} \right) + \gamma \leq \lambda_1(\widetilde{\mathbf{K}}_{\text{inner}}^{\text{lin}}) \leq \beta\lambda_1 \left(\frac{\mathbf{X}\mathbf{X}^\top}{d} \right) + \gamma + \alpha n,$$

and

$$\beta\lambda_i \left(\frac{\mathbf{X}\mathbf{X}^\top}{d} \right) + \gamma \leq \lambda_i(\widetilde{\mathbf{K}}_{\text{inner}}^{\text{lin}}) \leq \beta\lambda_{i-1} \left(\frac{\mathbf{X}\mathbf{X}^\top}{d} \right) + \gamma, \quad i = 2, 3, \dots, n,$$

so that the eigenvalue of $\widetilde{\mathbf{K}}_{\text{inner}}^{\text{lin}}$ interlaced with those of $\beta\mathbf{X}\mathbf{X}^\top/d + \gamma\mathbf{I}$. We can thus conclude that the eigenvalue decay of $\widetilde{\mathbf{K}}_{\text{inner}}^{\text{lin}}$ is the same as that of $\mathbf{X}\mathbf{X}^\top/d$ with a constant shift and scaling, which do not effect the trend of eigenvalue decay. Accordingly, the inner-product-type kernel matrix $\mathbf{K}_{\text{inner}}$ and its linearization $\widetilde{\mathbf{K}}_{\text{inner}}^{\text{lin}}$, $\widetilde{\mathbf{X}}$ admit the same eigenvalue decay as $\mathbf{X}\mathbf{X}^\top/d$, which concludes the proof. \square

Proposition 2 also provides a justification to study the eigenvalue decay of a radial kernel matrix. According to Theorem 2.2 in [26], the radial kernel matrix $\mathbf{K}_{\text{radial}}$ can be well approximated by $\widetilde{\mathbf{K}}_{\text{radial}}^{\text{lin}}$ with

$$\widetilde{\mathbf{K}}_{\text{radial}}^{\text{lin}} := \beta \frac{\mathbf{X}\mathbf{X}^\top}{d} + \gamma\mathbf{I} + \alpha\mathbf{1}\mathbf{1}^\top + h'(2\tau)\mathbf{A} + \frac{1}{2}h''(2\tau)\mathbf{A} \odot \mathbf{A},$$

in a spectral norm sense, where α, β, γ are given in Table 2. Recall $\mathbf{A} := \mathbf{1}\boldsymbol{\psi}^\top + \boldsymbol{\psi}\mathbf{1}^\top$, where $\boldsymbol{\psi} \in \mathbb{R}^n$ with $\psi_i := \|\mathbf{x}_i\|_2^2/d - \tau$, we find that \mathbf{A} is a rank 2 matrix with its eigenvalues $\lambda(\mathbf{A}) = \mathbf{1}^\top\boldsymbol{\psi} \pm \sqrt{n}\|\boldsymbol{\psi}\|_2$, and thus we have $\text{rank}(\mathbf{A} \odot \mathbf{A}) = 3$.³ Hence, by virtue of Proposition 2, apart from the top 5 eigenvalues of the radial kernel matrix $\mathbf{K}_{\text{radial}}$, its remaining eigenvalues follow with

$$\beta\lambda_i \left(\frac{\mathbf{X}\mathbf{X}^\top}{d} \right) + \gamma \leq \lambda_i(\widetilde{\mathbf{K}}_{\text{radial}}^{\text{lin}}) \leq \beta\lambda_{i-1} \left(\frac{\mathbf{X}\mathbf{X}^\top}{d} \right) + \gamma, \quad i = 6, 7, \dots, n.$$

Accordingly, $\mathbf{K}_{\text{radial}}$ admits the same eigenvalue decay as $\mathbf{X}\mathbf{X}^\top/d$.

C Proof of Lemma 1

Proof. By virtue of the closed form of the KRR estimator in Eq. (3) and $\boldsymbol{\epsilon} := \mathbf{y} - f_\rho(\mathbf{X})$, we have

$$f_{\mathbf{z},\lambda}(\mathbf{x}) - f_\rho(\mathbf{x}) = k(\mathbf{x}, \mathbf{X})^\top (\mathbf{K} + n\lambda\mathbf{I})^{-1} \boldsymbol{\epsilon} + k(\mathbf{x}, \mathbf{X})^\top (\mathbf{K} + n\lambda\mathbf{I})^{-1} f_\rho(\mathbf{X}) - f_\rho(\mathbf{x}),$$

where $f_\rho(\mathbf{X}) = [f_\rho(\mathbf{x}_1), f_\rho(\mathbf{x}_1), \dots, f_\rho(\mathbf{x}_n)]^\top \in \mathbb{R}^n$. According to $\mathbb{E}_{y|\mathbf{x}}[\boldsymbol{\epsilon}] = 0$, we then have

$$\mathbb{E}_{y|\mathbf{x}} \|f_{\mathbf{z},\lambda} - f_\rho\|_{\mathcal{L}_{\rho_X}^2}^2 = \mathbb{E}_{\mathbf{x}} \|k(\mathbf{x}, \cdot)^\top (\mathbf{K} + n\lambda\mathbf{I})^{-1} f_\rho(\mathbf{X}) - f_\rho\|_{\mathcal{L}_{\rho_X}^2}^2 + \mathbb{E}_{y,\mathbf{x}} \|k(\mathbf{x}, \cdot)^\top (\mathbf{K} + n\lambda\mathbf{I})^{-1} \boldsymbol{\epsilon}\|_{\mathcal{L}_{\rho_X}^2}^2.$$

Based on the definition of \mathbf{B} , we decompose \mathbf{B} as

$$\begin{aligned} \mathbf{B} &:= \mathbb{E}_{\mathbf{x}} \|k(\mathbf{x}, \cdot)^\top (\mathbf{K} + n\lambda\mathbf{I})^{-1} f_\rho(\mathbf{X}) - f_\rho\|_{\mathcal{L}_{\rho_X}^2}^2 = \|f_{\mathbf{X},\lambda} - f_\rho\|_{\mathcal{L}_{\rho_X}^2}^2 \\ &\leq 2\|f_{\mathbf{X},\lambda} - f_\lambda\|_{\mathcal{L}_{\rho_X}^2}^2 + 2\|f_\lambda - f_\rho\|_{\mathcal{L}_{\rho_X}^2}^2, \end{aligned}$$

which concludes our proof. \square

D Proof for the bias

The error bound for the bias is given by the following theorem.

³This can be proved using rank-one decomposition of \mathbf{A} .

Theorem 3. (Bias) Under Assumption 4 (source condition with $0 < r \leq 1$), Assumption 5 (capacity condition with $0 \leq \eta \leq 1$), let $0 < \delta < 1/2$, taking the regularization parameter $\lambda := \bar{c}n^{-\vartheta}$ with $0 \leq \vartheta \leq \frac{1}{1+\eta}$, there holds with probability at least $1 - 2\delta$, we have

$$\mathbb{B} \leq 2 \left(\|f_{\mathbf{X},\lambda} - f_\lambda\|_{\mathcal{L}_{\rho_{\mathbf{X}}}^2}^2 + \|f_\lambda - f_\rho\|_{\mathcal{L}_{\rho_{\mathbf{X}}}^2}^2 \right) \lesssim n^{-2\vartheta r} \log^4 \left(\frac{2}{\delta} \right).$$

In our error decomposition, $\|f_\lambda - f_\rho\|_{\mathcal{L}_{\rho_{\mathbf{X}}}^2}^2$ is independent of data \mathbf{X} that corresponds to the approximation error in learning theory [30]; while the first term $\|f_{\mathbf{X},\lambda} - f_\lambda\|_{\mathcal{L}_{\rho_{\mathbf{X}}}^2}^2$ depends on \mathbf{X} , termed as bias-sample error. To prove Theorem 3, we need to bound the approximation error and the bias-sample error as follows.

D.1 Bound approximation error

In learning theory, the approximation error $\|f_\lambda - f_\rho\|_{\mathcal{L}_{\rho_{\mathbf{X}}}^2}$ can be estimated by the source condition in Assumption 4.

Lemma 2. (Lemma 3 in [61]) Under the source condition in Assumption 4 with $0 < r \leq 1$, the approximation error can be given by

$$\|f_\lambda - f_\rho\|_{\mathcal{L}_{\rho_{\mathbf{X}}}^2} = \|(L_K + \lambda I)^{-1} L_K f_\rho - f_\rho\|_{\mathcal{L}_{\rho_{\mathbf{X}}}^2} \leq \lambda^r \|L_K^{-r} f_\rho\|_{\mathcal{L}_{\rho_{\mathbf{X}}}^2} \leq R\lambda^r.$$

D.2 Bound bias-sample error

To bound the bias-sample error $\|f_{\mathbf{X},\lambda} - f_\lambda\|_{\mathcal{L}_{\rho_{\mathbf{X}}}^2}$, we need the following lemma.

Lemma 3. (Lemma 17 in [62]) For any $0 < \delta < 1$, it holds with probability at least $1 - \delta$ that

$$\|(L_K + \lambda I)^{-1/2} (L_K - L_{K,\mathbf{X}})\| \leq \frac{2\kappa}{\sqrt{n}} \left\{ \frac{\kappa}{\sqrt{n\lambda}} + \sqrt{\mathcal{N}(\lambda)} \right\} \log \left(\frac{2}{\delta} \right),$$

where $\kappa := \max\{1, \sup_{\mathbf{x} \in X} \sqrt{k(\mathbf{x}, \mathbf{x})}\}$.

Then the bias-sample error can be decomposed into several parts.

Lemma 4. Under Assumption 4, we have

$$\begin{aligned} \|f_{\mathbf{X},\lambda} - f_\lambda\| &\leq R\lambda^{1/2} \|(L_{K,\mathbf{X}} + \lambda I)^{-1/2} (L_K + \lambda I)^{1/2}\| \|(L_K + \lambda I)^{-1/2} (L_K - L_{K,\mathbf{X}})\|^r \\ &\quad \|(L_K + \lambda I)^{-1/2} (L_K - L_{K,\mathbf{X}}) (L_K + \lambda I)^{-1}\|^{1-r}. \end{aligned}$$

Proof of Lemma 4. According to the definition of $f_{\mathbf{X},\lambda}$ and f_λ , we have

$$f_{\mathbf{X},\lambda} - f_\lambda = (L_{K,\mathbf{X}} + \lambda I)^{-1} L_{K,\mathbf{X}} f_\rho - (L_K + \lambda I)^{-1} L_K f_\rho.$$

Due to $(A + \lambda I)^{-1} A = I - \lambda(A + \lambda I)^{-1}$ for any bounded positive operator A , we have

$$(L_K + \lambda I)^{-1} L_K f_\rho - (L_{K,\mathbf{X}} + \lambda I)^{-1} L_{K,\mathbf{X}} f_\rho = \lambda [(L_{K,\mathbf{X}} + \lambda I)^{-1} - (L_K + \lambda I)^{-1}] f_\rho.$$

Further, by virtue of the first order decomposition of operator difference: $A^{-1} - B^{-1} = A^{-1}(B - A)B^{-1}$ for any invertible bounded operator and using the source condition in Assumption 4, the above equation can be further expressed as

$$\begin{aligned} (L_K + \lambda I)^{-1} L_K f_\rho - (L_{K,\mathbf{X}} + \lambda I)^{-1} L_{K,\mathbf{X}} f_\rho &= \lambda (L_{K,\mathbf{X}} + \lambda I)^{-1} (L_K - L_{K,\mathbf{X}}) (L_K + \lambda I)^{-1} L_K^r g_\rho \\ &= \lambda^{1/2} \left(\lambda^{1/2} (L_{K,\mathbf{X}} + \lambda I)^{-1/2} \right) \left((L_{K,\mathbf{X}} + \lambda I)^{-1/2} (L_K + \lambda I)^{1/2} \right) \\ &\quad \left((L_K + \lambda I)^{-1/2} (L_K - L_{K,\mathbf{X}}) (L_K + \lambda I)^{-(1-r)} \right) \left((L_K + \lambda I)^{-r} L_K^r \right) g_\rho. \end{aligned}$$

Besides, using $\|AB^t\| \leq \|A\|^{1-t} \|AB\|^t$ with $t \in [0, 1]$ for any bounded linear operator A and positive semi-definite operator B in Proposition 9 in [37], we have

$$\begin{aligned} \|(L_K + \lambda I)^{-1/2} (L_K - L_{K,\mathbf{X}}) (L_K + \lambda I)^{-(1-r)}\| &\leq \|(L_K + \lambda I)^{-1/2} (L_K - L_{K,\mathbf{X}})\|^r \\ &\quad \|(L_K + \lambda I)^{-1/2} (L_K - L_{K,\mathbf{X}}) (L_K + \lambda I)^{-1}\|^{1-r}, \end{aligned}$$

where we choose $A := (L_K + \lambda I)^{-1/2}(L_K - L_{K,\mathbf{X}})$, $B := (L_K + \lambda I)^{-1}$, and $t := 1 - r \in [0, 1)$. Accordingly, we can conclude our proof due to $\|(L_{K,\mathbf{X}} + \lambda I)^{-1/2}\| \leq 1/\sqrt{\lambda}$ and $\|(L_K + \lambda I)^{-r} L_K^r\| \leq 1$. \square

Remark: The proof framework of Lemma 4 is similar to Lemma 4 in [37] but we consider a more general case $0 < r \leq 1$ than $1/2 \leq r \leq 1$ in [37]. Although $0 < r < 1/2$ appears to be unattainable as claimed in [37], we follow with [62, 63] on a quite general case with $r > 0$.

To prove Theorem 3, we also need the following two lemmas.

Lemma 5. (Proposition 6 in [37]) Let $\delta \in (0, 1/2]$, it holds with probability at least $1 - 2\delta$ that

$$\begin{aligned} & \|(L_K + \lambda I)^{-1/2}(L_K - L_{K,\mathbf{X}})(L_K + \lambda I)^{-1}\| \\ & \leq \|(L_K + \lambda I)^{-1/2}(L_K - L_{K,\mathbf{X}})(L_K + \lambda I)^{-1/2}\| \|(L_K + \lambda I)^{-1/2}\| \leq \left(\frac{\kappa^2}{3n\lambda} + \sqrt{\frac{\kappa^2}{n\lambda}} \right) \frac{1}{\sqrt{\lambda}}. \end{aligned}$$

Lemma 6. For any $0 < \delta < 1$, with probability at least $1 - \delta$, we have

$$\|(L_{K,\mathbf{X}} + \lambda I)^{-1/2}(L_K + \lambda I)^{1/2}\| \leq 1 + \frac{2\kappa}{\sqrt{n\lambda}} \left\{ \frac{\kappa}{\sqrt{n\lambda}} + \sqrt{\mathcal{N}(\lambda)} \right\} \log \left(\frac{2}{\delta} \right).$$

Proof of Lemma 6. By virtue of a second order decomposition of operator difference in Lemma 16 [62], we have

$$A^{-1} - B^{-1} = B^{-1}(B - A)A^{-1}(B - A)B^{-1} + B^{-1}(B - A)B^{-1},$$

which leads to

$$A^{-1}B = I + B^{-1}(B - A) + B^{-1}(B - A)A^{-1}(B - A). \quad (11)$$

Accordingly, denote $A := L_{K,\mathbf{X}} + \lambda I$ and $B := L_K + \lambda I$, we can derive that

$$\begin{aligned} & \|(L_{K,\mathbf{X}} + \lambda I)^{-1/2}(L_K + \lambda I)^{1/2}\| \leq \|(L_{K,\mathbf{X}} + \lambda I)^{-1}(L_K + \lambda I)\|^{1/2} \\ & \leq \sqrt{1 + \lambda^{-1/2} \|(L_K + \lambda I)^{-1/2}(L_K - L_{K,\mathbf{X}})\| + \lambda^{-1} \|(L_K + \lambda I)^{-1/2}(L_K - L_{K,\mathbf{X}})\|^2} \\ & \leq \sqrt{1 + \mathcal{A} + \mathcal{A}^2} \leq 1 + \mathcal{A}, \end{aligned}$$

where $\mathcal{A} := \frac{2\kappa}{\sqrt{n\lambda}} \left\{ \frac{\kappa}{\sqrt{n\lambda}} + \sqrt{\mathcal{N}(\lambda)} \right\} \log(2/\delta)$ by Lemma 3. The first inequality holds by $\|A^s B^s\| \leq \|AB\|^s$ with $0 \leq s \leq 1$ for positive operators A and B on Hilbert spaces [39]. The second inequality can be derived by Eq. (11), $\|(L_{K,\mathbf{X}} + \lambda I)^{-1}\| \leq 1/\lambda$ and $\|(L_K + \lambda I)^{-1/2}\| \leq 1/\sqrt{\lambda}$. \square

Remark: Lemma 7.2 in [64] gives $\|(L_{K,\mathbf{X}} + \lambda I)^{-1/2}(L_K + \lambda I)^{1/2}\| \leq \sqrt{2}$ by assuming $\lambda > \frac{9}{n} \log \frac{n}{\delta}$; whereas our result does not require extra conditions on λ .

Based on the above lemmas, we are ready to prove Theorem 3.

Proof of Theorem 3. We first estimate $\|(L_{K,\mathbf{X}} + \lambda I)^{-1/2}(L_K + \lambda I)^{1/2}\|$ in Lemma 6 by taking $\lambda := \bar{c}n^{-\vartheta}$ and the capacity condition in Assumption 5: $\mathcal{N}(\lambda) \leq Q^2\lambda^{-\eta}$ with $\eta \in [0, 1]$. Accordingly, we have

$$\begin{aligned} & \|(L_{K,\mathbf{X}} + \lambda I)^{-1/2}(L_K + \lambda I)^{1/2}\| \leq 1 + \frac{2\kappa}{\sqrt{n\lambda}} \left\{ \frac{\kappa}{\sqrt{n\lambda}} + \sqrt{\mathcal{N}(\lambda)} \right\} \log \left(\frac{2}{\delta} \right) \\ & \leq 1 + \left(\frac{2\kappa^2}{\bar{c}} n^{-(1-\vartheta)} + 2\kappa\bar{c}^{-\left(\frac{1}{2} + \frac{\eta}{2}\right)} Q n^{-\frac{1-\vartheta-\vartheta\eta}{2}} \right) \log \left(\frac{2}{\delta} \right) \\ & \leq \left(1 + \frac{2\kappa(\kappa + Q)}{\bar{c}} n^{-\frac{1-\vartheta-\vartheta\eta}{2}} \right) \log \left(\frac{2}{\delta} \right), \end{aligned}$$

where we use $\log^r(2/\delta) \leq \log(2/\delta)$ due to $\log(2/\delta) > 1$ in the last inequality. Since $\|(L_{K,\mathbf{X}} + \lambda I)^{-1/2}(L_K + \lambda I)^{1/2}\|$ converges to zero when n is large enough, we require $\vartheta < \frac{1}{1+\eta}$ to ensure a positive convergence rate, which implies

$\vartheta \leq 1$. Then we bound $\|(L_K + \lambda I)^{-1/2}(L_K - L_{K,\mathbf{X}})\|^r$ by Lemma 3. By virtue of $(a + b)^r \leq a^r + b^r$ for any $r \in (0, 1]$ and $a, b \geq 0$, we have

$$\begin{aligned} \|(L_K + \lambda I)^{-1/2}(L_K - L_{K,\mathbf{X}})\|^r &\leq \left(\frac{2\kappa}{\sqrt{n}}\right)^r \left\{ \frac{\kappa}{(n\lambda)^{\frac{r}{2}}} + [\mathcal{N}(\lambda)]^{\frac{r}{2}} \right\} \log\left(\frac{2}{\delta}\right) \\ &\leq (2\kappa)^r (n\bar{c})^{-\frac{r}{2}} \left[\kappa n^{-\frac{r(1-\vartheta)}{2}} + \frac{Q}{\bar{c}^{\frac{\eta r}{2}}} n^{\frac{\vartheta \eta r}{2}} \right] \log\left(\frac{2}{\delta}\right) \\ &\leq \frac{2\kappa(Q + \kappa)}{\bar{c}} n^{-\frac{(1-\vartheta)r}{2}} \log\left(\frac{2}{\delta}\right), \end{aligned}$$

where the second one admits by the capacity condition in Assumption 5. Similarly, to bound $\|(L_K + \lambda I)^{-1/2}(L_K - L_{K,\mathbf{X}})(L_K + \lambda I)^{-1}\|^{1-r}$ by Lemma 5, we can derive that

$$\begin{aligned} \|(L_K + \lambda I)^{-1/2}(L_K - L_{K,\mathbf{X}})(L_K + \lambda I)^{-1}\|^{1-r} &\leq \lambda^{-\frac{1-r}{2}} \left(\kappa^2 (n\lambda)^{-(1-r)} + \kappa (n\lambda)^{-\frac{1-r}{2}} \right) \\ &\leq \frac{\kappa^2}{\bar{c}} \left(n^{\frac{3}{2}\vartheta-1)(1-r)} + n^{(\vartheta-\frac{1}{2})(1-r)} \right) \\ &\leq \frac{\kappa^2}{\bar{c}} n^{(\vartheta-\frac{1}{2})(1-r)}. \end{aligned}$$

Combining the above three inequalities, we have

$$\begin{aligned} \|f_{\mathbf{X},\lambda} - f_\lambda\| &\leq \frac{4R\kappa^3(Q + \kappa)^2}{\bar{c}^3} n^{-\frac{(1-\vartheta)\eta r + \vartheta}{2}} n^{(\vartheta-\frac{1}{2})(1-r)} \log^2\left(\frac{2}{\delta}\right) \\ &\leq \widetilde{C_{R,Q,\kappa,\bar{c}}} n^{-\frac{1-\vartheta(\eta r + 1 - 2r)}{2}} \log^2\left(\frac{2}{\delta}\right), \end{aligned}$$

where $\widetilde{C_{R,Q,\kappa,\bar{c}}} := 4R\kappa^3(Q + \kappa)^2/\bar{c}^3$ is independent of n and d .

Finally, the bias can be bounded by

$$\begin{aligned} \mathbf{B} &\leq 2\|f_{\mathbf{X},\lambda} - f_\lambda\|_{\mathcal{L}_{\rho_{\mathbf{X}}}}^2 + 2\|f_\lambda - f_\rho\|_{\mathcal{L}_{\rho_{\mathbf{X}}}}^2 \\ &\leq 2R^2 n^{-2\vartheta r} + \widetilde{C}_1 n^{-[1-\vartheta(\eta r + 1 - 2r)]} \log^4\left(\frac{2}{\delta}\right) \\ &\leq \widetilde{C} n^{-2\vartheta r} \log^4\left(\frac{2}{\delta}\right), \end{aligned}$$

where the third inequality holds by $2\vartheta r \leq 1 - \vartheta(\eta r + 1 - 2r)$ due to $\vartheta \leq \frac{1}{1+\eta}$, and $\widetilde{C}, \widetilde{C}_1$ are some constants independent of n and d . Accordingly, we can conclude the proof. \square

E Proof for the variance

Formally, we have the following theorem to bound the variance.

Theorem 4. (Variance) Under Assumptions 2, 3, then for $0 < \delta < 1$ with probability $1 - \delta - d^{-2}$, $\theta = \frac{1}{2} - \frac{2}{8+m}$, and d large enough, for any given $\varepsilon > 0$, we have

$$\mathbf{V} \lesssim \mathbf{V}_1 + \mathbf{V}_2,$$

where $\mathbf{V}_1 := \frac{\sigma^2 \beta}{d} \mathcal{N}_{\mathbf{X}}^{n\lambda+\gamma}$ and \mathbf{V}_2 is the residual term with

$$\mathbf{V}_2 := \begin{cases} \frac{\sigma^2 \log^{2+4\varepsilon} d}{(n\lambda + \gamma)^2 d^{4\theta-1}}, & \text{inner-product kernels} \\ \frac{\sigma^2}{(n\lambda + \gamma)^2} d^{-2\theta} \log^{1+\varepsilon} d, & \text{radial kernels.} \end{cases}$$

For inner-product kernels, our proof framework follows [18], and is briefly discussed in Section E.1. Nevertheless, error bound on radial kernels has not been investigated in [18] and is more subtle to handle (than that of inner-product kernels) due to the additionally introduced \mathbf{A} and $\mathbf{A} \odot \mathbf{A}$ in Table 2. Accordingly, we mainly focus on proofs for radial kernels.

E.1 Inner-product kernel matrices

In this subsection, we consider the inner-product kernel case with $k(\mathbf{x}, \mathbf{x}') = h(\langle \mathbf{x}, \mathbf{x}' \rangle / d)$. We briefly introduce our results that can be derived from proofs of Theorem 2 in [18] for completeness.

To prove Theorem 4, define

$$\widetilde{\mathbf{K}}^{\text{lin}}(\mathbf{X}, \mathbf{X}) := (n\lambda + \gamma)\mathbf{I} + \alpha\mathbf{1}\mathbf{1}^\top + \beta \frac{\mathbf{X}\mathbf{X}^\top}{d} \in \mathbb{R}^{n \times n}, \quad k^{\text{lin}}(\mathbf{x}, \mathbf{X}) := h(0)\mathbf{1} + \beta \frac{\mathbf{X}\mathbf{x}^\top}{d} \in \mathbb{R}^{n \times 1}, \quad (12)$$

and $k^{\text{lin}}(\mathbf{X}, \mathbf{x})$ is the transpose of $k^{\text{lin}}(\mathbf{x}, \mathbf{X})$. Note that γ in $\widetilde{\mathbf{K}}^{\text{lin}}$ corresponds to the *implicit* regularization and $n\lambda$ corresponds to the *explicit* regularization. Now we prove Theorem 4 for inner-product kernels.

Proof of Theorem 4 for inner-product kernels. According to the definition of \mathbf{V} , we have

$$\begin{aligned} \mathbf{V} &= \mathbb{E}_{\mathbf{x}, y} \text{tr} \left[k(\mathbf{x}, \mathbf{X})^\top (\mathbf{K} + n\lambda\mathbf{I})^{-1} \boldsymbol{\epsilon} \boldsymbol{\epsilon}^\top (\mathbf{K} + n\lambda\mathbf{I})^{-1} k(\mathbf{x}, \mathbf{X}) \right] = \mathbb{E}_{\mathbf{x}} \left\| (\mathbf{K} + n\lambda\mathbf{I})^{-1} k(\mathbf{x}, \mathbf{X}) \right\|_2^2 \mathbb{E}_{y|\mathbf{x}} \|\boldsymbol{\epsilon}\|_2^2 \\ &\leq \sigma^2 \mathbb{E}_{\mathbf{x}} \left\| (\mathbf{K} + n\lambda\mathbf{I})^{-1} k(\mathbf{x}, \mathbf{X}) \right\|_2^2 \\ &\leq \sigma^2 \left\| (\mathbf{K} + n\lambda\mathbf{I})^{-1} \widetilde{\mathbf{K}}^{\text{lin}} \right\|_2^2 \mathbb{E}_{\mathbf{x}} \left\| [\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} k^{\text{lin}}(\mathbf{x}, \mathbf{X}) \right\|_2^2 + \sigma^2 \left\| (\mathbf{K} + n\lambda\mathbf{I})^{-1} \right\|_2^2 \mathbb{E}_{\mathbf{x}} \left\| k(\mathbf{x}, \mathbf{X}) - k^{\text{lin}}(\mathbf{x}, \mathbf{X}) \right\|_2^2, \end{aligned} \quad (13)$$

where the first inequality comes from Assumption 2. To bound the terms in Eq. (13), we need

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} \left\| [\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} k^{\text{lin}}(\mathbf{x}, \mathbf{X}) \right\|_2^2 &= \mathbb{E}_{\mathbf{x}} \text{tr} \left[[\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} \left(\beta \frac{\mathbf{X}\mathbf{x}}{d} + h(0)\mathbf{1} \right) \left(\beta \frac{\mathbf{x}^\top \mathbf{X}^\top}{d} + h(0)\mathbf{1}^\top \right) [\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} \right] \\ &\leq \frac{1}{d} \|\boldsymbol{\Sigma}_d\|_2 \text{tr} \left([\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} \beta^2 \frac{\mathbf{X}\mathbf{X}^\top}{d} [\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} \right) + \frac{1}{d} \text{tr} \left([\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} h(0)^2 \mathbf{1}\mathbf{1}^\top [\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} \right) \\ &\leq \frac{\beta}{d} \|\boldsymbol{\Sigma}_d\|_2 \sum_{j=1}^n \frac{\lambda_j(\widetilde{\mathbf{X}})}{\left[n\lambda + \gamma + \lambda_j(\widetilde{\mathbf{X}}) \right]^2} + \frac{1}{d} \frac{h(0)^2 n}{\left[n\lambda + \gamma + \lambda_1(\widetilde{\mathbf{X}}) \right]^2} \\ &\asymp \frac{\beta}{d} \mathcal{N}_{\widetilde{\mathbf{X}}}^{n\lambda + \gamma} + \mathcal{O}\left(\frac{1}{nd}\right). \end{aligned} \quad (14)$$

To bound the remaining terms in Eq. (13), we also need the following results that can be obtained from [18]:

- (i) By Proposition A.2 in [18], with probability at least $1 - \delta - d^{-2}$, for $\theta = \frac{1}{2} - \frac{2}{8+m}$ and any given $\varepsilon > 0$, we have $\left\| \mathbf{K} + n\lambda\mathbf{I} - \widetilde{\mathbf{K}}^{\text{lin}} \right\|_2 \leq d^{-\theta} (\delta^{-1/2} + \log^{0.5+\varepsilon} d)$ and $\mathbb{E}_{\mathbf{x}} \left\| k(\mathbf{x}, \mathbf{X}) - k^{\text{lin}}(\mathbf{x}, \mathbf{X}) \right\|_2^2 \leq \widetilde{C}_1 d^{-(4\theta-1)} \log^{2+4\varepsilon} d$.
- (ii) $\left\| (\mathbf{K} + n\lambda\mathbf{I})^{-1} \right\|_2 \leq \frac{2}{n\lambda + \gamma}$ and $\left\| (\mathbf{K} + n\lambda\mathbf{I})^{-1} \widetilde{\mathbf{K}}^{\text{lin}} \right\|_2 \leq 2$ provided d is large enough such that $d^{-\theta} (\delta^{-1/2} + \log^{0.5+\varepsilon} d) \leq \gamma/2$.

Combining the above results, with probability at least $1 - \delta - d^{-2}$, for any given $\varepsilon > 0$, The error bound for the variance in Eq. (13) can be further given by

$$\begin{aligned} \mathbf{V} &\leq \sigma^2 \mathbb{E}_{\mathbf{x}} \left\| (\mathbf{K} + n\lambda\mathbf{I})^{-1} k(\mathbf{x}, \mathbf{X}) \right\|_2^2 \\ &\leq 2\sigma^2 \left\| (\mathbf{K} + n\lambda\mathbf{I})^{-1} \widetilde{\mathbf{K}}^{\text{lin}} \right\|_2^2 \mathbb{E}_{\mathbf{x}} \left\| [\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} k^{\text{lin}}(\mathbf{x}, \mathbf{X}) \right\|_2^2 + 2\sigma^2 \left\| \mathbf{K}^{-1} \right\|_2^2 \mathbb{E}_{\mathbf{x}} \left\| k(\mathbf{x}, \mathbf{X}) - k^{\text{lin}}(\mathbf{x}, \mathbf{X}) \right\|_2^2 \\ &\asymp \frac{8\sigma^2\beta}{d} \|\boldsymbol{\Sigma}_d\|_2 \sum_{j=1}^n \frac{\lambda_j(\widetilde{\mathbf{X}})}{\left[n\lambda + \gamma + \lambda_j(\widetilde{\mathbf{X}}) \right]^2} + \frac{8\sigma^2}{(n\lambda + \gamma)^2} \widetilde{C}_1 d^{-(4\theta-1)} \log^{2+4\varepsilon} d \\ &\asymp \frac{\sigma^2\beta}{d} \mathcal{N}_{\widetilde{\mathbf{X}}}^{n\lambda + \gamma} + \frac{\sigma^2}{(n\lambda + \gamma)^2} d^{-(4\theta-1)} \log^{2+4\varepsilon} d, \end{aligned}$$

which concludes the proof. \square

E.2 Radial kernel matrices

In this subsection, we consider the radial kernel case with $k(\mathbf{x}, \mathbf{x}') = h\left(\frac{1}{d}\|\mathbf{x} - \mathbf{x}'\|_2^2\right)$. Since the linearization of radial kernel matrices incurs in two additional terms \mathbf{A} and $\mathbf{A} \odot \mathbf{A}$, estimation for radial kernels is more technical than that of inner-product kernels. Accordingly, to prove Theorem 4 for radial kernels, we need to introduce the following notations and auxiliary results.

E.2.1 Auxiliary results

Recall $\tau := \text{tr}(\Sigma_d)/d$, define

$$\begin{aligned} \widetilde{\mathbf{K}}^{\text{lin}}(\mathbf{X}, \mathbf{X}) &:= (\gamma + n\lambda)\mathbf{I} + \alpha\mathbf{1}\mathbf{1}^\top + \beta\frac{\mathbf{X}\mathbf{X}^\top}{d} + h'(2\tau)\mathbf{A} + \frac{1}{2}h''(2\tau)\mathbf{A} \odot \mathbf{A} \\ k^{\text{lin}}(\mathbf{x}, \mathbf{X}) &:= h(2\tau)\mathbf{1} + \beta\frac{\mathbf{X}\mathbf{x}^\top}{d} - \frac{\beta}{2}\mathbf{A}(\mathbf{x}, \mathbf{X}) \in \mathbb{R}^{n \times 1}, \end{aligned} \quad (15)$$

where $\mathbf{A}(\mathbf{x}, \mathbf{X}) := \psi_{\mathbf{x}} + [\psi_1, \psi_2, \dots, \psi_n]^\top$ with $\psi_{\mathbf{x}} = \|\mathbf{x}\|_2^2/d - \tau$ and $\psi_i = \|\mathbf{x}_i\|_2^2/d - \tau$ for $i = 1, 2, \dots, n$. As discussed in Appendix B, we conclude that $\widetilde{\mathbf{K}}^{\text{lin}}$ admits the same eigenvalue decay as $\widetilde{\mathbf{X}}$ since \mathbf{A} is a rank-2 matrix. Accordingly, we have the following results.

Proposition 3. *Given $\mathbf{A}(\mathbf{x}, \mathbf{X})$ in Eq. (15), we have $\mathbb{E}_{\mathbf{x}}[\mathbf{X}\mathbf{x}\mathbf{A}(\mathbf{X}, \mathbf{x})] = \mu_3\mathbf{X}\Sigma_d^{1/2}\text{diag}(\Sigma_d)\mathbf{1}_n^\top$, where $\mu_3 := \mathbb{E}[\mathbf{t}(j)^3]$ does not depend on j because each entry in \mathbf{t} are independent for $j = 1, 2, \dots, d$. Further, $\mathbb{E}_{\mathbf{x}}[\mathbf{X}\mathbf{x}\mathbf{A}(\mathbf{X}, \mathbf{x})]$ is a rank-one matrix with its eigenvalue $\lambda_1(\mathbb{E}_{\mathbf{x}}[\mathbf{X}\mathbf{x}\mathbf{A}(\mathbf{X}, \mathbf{x})]) = \mathcal{O}(\sqrt{n/d})$.*

Proof of Proposition 3. According to the definition in Assumption 3, $\mathbf{x}_i = \Sigma_d^{1/2}\mathbf{t}_i$ with $\mathbb{E}[\mathbf{t}_i(j)] = 0$ and $\mathbb{V}[\mathbf{t}_i(j)] = 1$, we have the following expression

$$\mathbb{E}_{\mathbf{t}}[\mathbf{t}\mathbf{t}^\top\Sigma_d\mathbf{t}] = \mathbb{E}_{\mathbf{t}}\left[\mathbf{t}\sum_{i,j=1}^d \mathbf{t}(i)(\Sigma_d)_{ij}\mathbf{t}(j)\right] = \mu_3[(\Sigma_d)_{11}, (\Sigma_d)_{22}, \dots, (\Sigma_d)_{dd}]^\top,$$

where $\mu_3 := \mathbb{E}(\mathbf{t}_i^3)$. Accordingly, $\mathbb{E}_{\mathbf{x}}[\mathbf{X}\mathbf{x}\mathbf{A}(\mathbf{X}, \mathbf{x})]$ can be computed by

$$\begin{aligned} \mathbb{E}_{\mathbf{x}}[\mathbf{X}\mathbf{x}\mathbf{A}(\mathbf{X}, \mathbf{x})] &= \mathbb{E}_{\mathbf{x}}[\mathbf{X}\mathbf{x}(\psi_1 + \psi_{\mathbf{x}}), \mathbf{X}\mathbf{x}(\psi_2 + \psi_{\mathbf{x}}), \dots, \mathbf{X}\mathbf{x}(\psi_n + \psi_{\mathbf{x}})] \\ &= \mathbb{E}_{\mathbf{x}}[\mathbf{X}\mathbf{x}\psi_{\mathbf{x}}, \mathbf{X}\mathbf{x}\psi_{\mathbf{x}}, \dots, \mathbf{X}\mathbf{x}\psi_{\mathbf{x}}] \\ &= \mathbf{X}\Sigma_d^{1/2}\left[\frac{\mathbb{E}_{\mathbf{t}}[\mathbf{t}\mathbf{t}^\top\Sigma_d\mathbf{t}]}{d}, \frac{\mathbb{E}_{\mathbf{t}}[\mathbf{t}\mathbf{t}^\top\Sigma_d\mathbf{t}]}{d}, \dots, \frac{\mathbb{E}_{\mathbf{t}}[\mathbf{t}\mathbf{t}^\top\Sigma_d\mathbf{t}]}{d}\right] \\ &= \mu_3\mathbf{X}\Sigma_d^{1/2}\text{diag}(\Sigma_d)\mathbf{1}_n^\top. \end{aligned}$$

Note that, the matrix $\text{diag}(\Sigma_d)\mathbf{1}_n^\top$ is a rank-one matrix, which implies $\text{rank}(\mathbf{X}\Sigma_d^{1/2}\text{diag}(\Sigma_d)\mathbf{1}_n^\top) \leq 1$. Accordingly, its non-zero eigenvalue $\lambda_1(\mathbf{X}\Sigma_d^{1/2}\text{diag}(\Sigma_d)\mathbf{1}_n^\top)$ admits

$$\frac{1}{d}\lambda_1(\mathbf{X}\Sigma_d^{1/2}\text{diag}(\Sigma_d)\mathbf{1}_n^\top) = \frac{1}{d}\sum_{i=1}^n \mathbf{x}_i^\top \Sigma_d^{1/2} \text{diag}(\Sigma_d) = \frac{1}{d}\sum_{i=1}^n \mathbf{t}_i^\top \Sigma_d \text{diag}(\Sigma_d).$$

Due to $\mathbb{E}[\mathbf{t}_i^\top \Sigma_d \text{diag}(\Sigma_d)] = 0$ and $\mathbb{V}[\mathbf{t}_i^\top \Sigma_d \text{diag}(\Sigma_d)] = \|\Sigma_d \text{diag}(\Sigma_d)\|_2^2$, which, with a central limit theorem argument, implies $\sum_{i=1}^n \mathbf{t}_i^\top \Sigma_d \text{diag}(\Sigma_d) = \mathcal{O}(\sqrt{nd})$ due to $\|\Sigma_d \text{diag}(\Sigma_d)\|_2 \leq \|\Sigma_d\|_2 \|\text{diag}(\Sigma_d)\|_2 \leq \tilde{C}\|\text{diag}(\Sigma_d)\|_2$. Accordingly, we can conclude that $\frac{1}{d}\lambda_1(\mathbf{X}\Sigma_d^{1/2}\text{diag}(\Sigma_d)\mathbf{1}_n^\top) = \mathcal{O}(\sqrt{n/d})$. \square

Proposition 4. *Given $\mathbf{A}(\mathbf{x}, \mathbf{X})$ in Eq. (15), we have $\mathbb{E}_{\mathbf{x}}[\mathbf{A}(\mathbf{x}, \mathbf{X})\mathbf{A}(\mathbf{X}, \mathbf{x})] = \psi\psi^\top + \mathcal{O}(1/d)$. Further, it has only one non-zero eigenvalue that admits $\lambda_1(\mathbb{E}_{\mathbf{x}}[\mathbf{A}(\mathbf{x}, \mathbf{X})\mathbf{A}(\mathbf{X}, \mathbf{x})]) = \mathcal{O}(n)$.*

Proof of Proposition 4. By virtue of the following results [26]

$$\begin{aligned} \frac{1}{d}\mathbb{E}_{\mathbf{x}}\|\mathbf{x}\|_2^2 &= \frac{1}{d}\mathbb{E}_{\mathbf{t}}[\mathbf{t}^\top\Sigma_d\mathbf{t}] = \tau \\ \mathbb{V}_{\mathbf{x}}\left[\frac{\|\mathbf{x}\|_2^2}{d}\right] &= \frac{1}{d^2}\left((\mu_4 - 3)\sum_{i=1}^d ((\Sigma_d)_{ii})^2 + 2\text{tr}(\Sigma_d^2)\right) = \mathcal{O}\left(\frac{1}{d}\right), \end{aligned}$$

where $\mu_4 := \mathbb{E}[\mathbf{t}(i)^4]$ does not depend on i . Accordingly, each entry in $\mathbb{E}_{\mathbf{x}}[\mathbf{A}(\mathbf{x}, \mathbf{X})\mathbf{A}(\mathbf{X}, \mathbf{x})]$ can be computed as

$$\begin{aligned} \mathbb{E}_{\mathbf{x}}[\mathbf{A}(\mathbf{x}, \mathbf{X})\mathbf{A}(\mathbf{X}, \mathbf{x})]_{ij} &= \mathbb{E}_{\mathbf{x}}[(\psi_i + \psi_{\mathbf{x}})(\psi_j + \psi_{\mathbf{x}})] \\ &= \psi_i\psi_j + (\psi_i + \psi_j)\mathbb{E}_{\mathbf{x}}\psi_{\mathbf{x}} + \mathbb{E}_{\mathbf{x}}[\psi_{\mathbf{x}}^2] \\ &= \psi_i\psi_j + \mathbb{V}_{\mathbf{x}}\left[\frac{\|\mathbf{x}\|_2^2}{d}\right] \\ &= \psi_i\psi_j + \frac{\mu_4 - 3}{d^2} \text{tr}(\boldsymbol{\Sigma}_d \odot \boldsymbol{\Sigma}_d) + \frac{2 \text{tr}(\boldsymbol{\Sigma}_d^2)}{d^2}. \end{aligned}$$

Then we have

$$\mathbb{E}_{\mathbf{x}}[\mathbf{A}(\mathbf{x}, \mathbf{X})\mathbf{A}(\mathbf{X}, \mathbf{x})] = \boldsymbol{\psi}\boldsymbol{\psi}^\top + \mathcal{O}(1/d).$$

Therefore, $\boldsymbol{\psi}\boldsymbol{\psi}^\top$ is a rank-one matrix with $\lambda_1(\boldsymbol{\psi}\boldsymbol{\psi}^\top) = \|\boldsymbol{\psi}\|_2^2 = \mathcal{O}(n)$. Then $\lambda_1(\mathbb{E}_{\mathbf{x}}[\mathbf{A}(\mathbf{x}, \mathbf{X})\mathbf{A}(\mathbf{X}, \mathbf{x})])$ can be estimated by

$$\|\boldsymbol{\psi}\|_2^2 \leq \lambda_1(\mathbb{E}_{\mathbf{x}}[\mathbf{A}(\mathbf{x}, \mathbf{X})\mathbf{A}(\mathbf{X}, \mathbf{x})]) \leq \|\boldsymbol{\psi}\|_2^2 + n \underbrace{\left[\frac{\mu_4 - 3}{d^2} \text{tr}(\boldsymbol{\Sigma}_d \odot \boldsymbol{\Sigma}_d) + \frac{2 \text{tr}(\boldsymbol{\Sigma}_d^2)}{d^2} \right]}_{=\mathcal{O}(1/d)},$$

which implies $\lambda_1(\mathbb{E}_{\mathbf{x}}[\mathbf{A}(\mathbf{x}, \mathbf{X})\mathbf{A}(\mathbf{X}, \mathbf{x})]) = \mathcal{O}(n)$. \square

Lemma 7. *Given a radial kernel, under Assumption 3, for $\theta = \frac{1}{2} - \frac{2}{8+m}$, we have with probability at least $1 - d^{-2}$ with respect to the draw of \mathbf{X} , for d large enough, for any given $\varepsilon > 0$, we have*

$$\mathbb{E}_{\mathbf{x}} \left\| k(\mathbf{x}, \mathbf{X}) - k^{\text{lin}}(\mathbf{x}, \mathbf{X}) \right\|_2^2 \leq \widetilde{C}_1 d^{-2\theta} \log^{1+\varepsilon} d,$$

where \widetilde{C}_1 is some constant independent of n and d .

Remark: In fact, we only need the $(5 + m)$ -moment in Assumption 3 but we still follow with it for simplicity.

Proof of Lemma 7. We start with the entry-wise Taylor expansion for the smooth kernel at 2τ with $\tau := \text{tr}(\boldsymbol{\Sigma}_d)/d$

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}_j) &= h\left(\frac{1}{d}\|\mathbf{x} - \mathbf{x}_j\|_2^2\right) = h(2\tau) + h'(2\tau)\left(\frac{1}{d}\|\mathbf{x} - \mathbf{x}_j\|_2^2 - 2\tau\right) + \frac{h''(2\tau)}{2} \left(\frac{1}{d}\|\mathbf{x} - \mathbf{x}_j\|_2^2 - 2\tau\right)^2 + \mathcal{O}(d^{-3/2}) \\ &= h(2\tau) + h'(2\tau)\left(\psi_{\mathbf{x}} + \psi_j - \frac{2\mathbf{x}^\top \mathbf{x}_j}{d}\right) + \frac{h''(2\tau)}{2} \left(\psi_{\mathbf{x}} + \psi_j - \frac{2\mathbf{x}^\top \mathbf{x}_j}{d}\right)^2 + \mathcal{O}(d^{-3/2}), \end{aligned}$$

where $\psi_j = \|\mathbf{x}_j\|_2^2/d - \tau$ for $j = 1, 2, \dots, n$ as defined before. Accordingly, by virtue of $k^{\text{lin}}(\mathbf{x}, \mathbf{x}_j) = \frac{\beta \mathbf{x}^\top \mathbf{x}_j}{d} - \frac{\beta}{2}(\psi_{\mathbf{x}} + \psi_j)$ and Corollary 2 in [26], with probability at least $1 - d^{-2}$, for any $\varepsilon > 0$, we have

$$k(\mathbf{x}, \mathbf{x}_j) - k^{\text{lin}}(\mathbf{x}, \mathbf{x}_j) = \frac{h''(2\tau)}{2} \left(\frac{1}{d}\|\mathbf{x} - \mathbf{x}_j\|_2^2 - 2\tau\right)^2 \leq \widetilde{C} d^{-1+\frac{4}{m}} (\log d)^{\frac{1+\varepsilon}{2}},$$

where we only need $(5 + m)$ -moment. Therefore, with probability at least $1 - d^{-2}$, for any given $\varepsilon > 0$, we have

$$\left\| k(\mathbf{x}, \mathbf{X}) - k^{\text{lin}}(\mathbf{x}, \mathbf{X}) \right\|_2 \leq C_1 d^{-1/2+\frac{4}{m}} (\log d)^{\frac{1+\varepsilon}{2}} \leq \widetilde{C}_1 d^{-\theta} (\log d)^{\frac{1+\varepsilon}{2}},$$

which implies

$$\mathbb{E}_{\mathbf{x}} \left\| k(\mathbf{x}, \mathbf{X}) - k^{\text{lin}}(\mathbf{x}, \mathbf{X}) \right\|_2^2 \leq \widetilde{C}_2 d^{-2\theta} \log^{1+\varepsilon} d,$$

where \widetilde{C}_1 and \widetilde{C}_2 are some constant independent of n and d . \square

E.2.2 Proofs of Theorem 4 for radial kernels

Now we are ready to prove Theorem 4 for radial kernels.

Proof of Theorem 4 for radial kernels. Similar to Eq. (13), to estimate $\mathbf{v} \leq \sigma^2 \mathbb{E}_{\mathbf{x}} \|(\mathbf{K} + n\lambda \mathbf{I})^{-1} k(\mathbf{x}, \mathbf{X})\|_2^2$, we need to bound subsequently the following terms: $\|\mathbf{K} + n\lambda \mathbf{I} - \widetilde{\mathbf{K}}^{\text{lin}}(\mathbf{X}, \mathbf{X})\|_2$, $\|(\mathbf{K} + n\lambda \mathbf{I})^{-1}\|_2$, $\|(\mathbf{K} + n\lambda \mathbf{I})^{-1} \widetilde{\mathbf{K}}^{\text{lin}}(\mathbf{X}, \mathbf{X})\|_2$, $\mathbb{E}_{\mathbf{x}} \|[\widetilde{\mathbf{K}}^{\text{lin}}(\mathbf{X}, \mathbf{X})]^{-1} k^{\text{lin}}(\mathbf{x}, \mathbf{X})\|_2^2$, and $\mathbb{E}_{\mathbf{x}} \|k(\mathbf{x}, \mathbf{X}) - k^{\text{lin}}(\mathbf{x}, \mathbf{X})\|_2^2$.

In [26], the approximation error between radial kernel matrices and their linearization can be decomposed into three parts: the first-order term A_1 , the second-order term A_2 , and the third-order term A_3

$$\|\mathbf{K} + n\lambda \mathbf{I} - \widetilde{\mathbf{K}}^{\text{lin}}(\mathbf{X}, \mathbf{X})\|_2 := A_1 + A_2 + A_3,$$

where A_1 and A_3 admit $\|A_1\|_2 \leq d^{-\theta} \log^{2+4\varepsilon} d$ and $\|A_3\|_2 \leq d^{-\theta} \log^{2+4\varepsilon} d$. The second-order term A_2 admits $\Pr(\|A_2\|_2 \leq d^{-\theta} \delta^{-1/2}) \leq \delta$ by Proposition A.2 in [18] and [26]. Accordingly, with probability at least $1 - \delta - d^{-2}$, for $\theta = \frac{1}{2} - \frac{2}{8+m}$ and any given $\varepsilon > 0$, we have

$$\|\mathbf{K} + n\lambda \mathbf{I} - \widetilde{\mathbf{K}}^{\text{lin}}(\mathbf{X}, \mathbf{X})\|_2 \leq d^{-\theta} \left(\delta^{-1/2} + \log^{2+4\varepsilon} d \right).$$

According to Proposition 3 and 4, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}} \|[\widetilde{\mathbf{K}}^{\text{lin}}(\mathbf{X}, \mathbf{X})]^{-1} k^{\text{lin}}(\mathbf{X}, \mathbf{x})\|_2^2 \\ &= \beta^2 \mathbb{E}_{\mathbf{x}} \text{tr} \left[[\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} \left(\frac{\mathbf{X} \mathbf{x} \mathbf{x}^\top \mathbf{X}^\top}{d^2} - \frac{\mathbf{X} \mathbf{x} \mathbf{A}(\mathbf{X}, \mathbf{x})}{d} + \frac{1}{4} \mathbf{A}(\mathbf{x}, \mathbf{X}) \mathbf{A}(\mathbf{X}, \mathbf{x}) + h(2\tau)^2 \mathbf{1} \mathbf{1}^\top \right) [\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} \right] \\ &\leq \beta^2 \text{tr} \left([\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} \left(\frac{\mathbf{X} \mathbf{X}^\top \|\Sigma_d\|_2}{d^2} - \frac{\mu_3 \mathbf{X} \Sigma_d^{1/2} \text{diag}(\Sigma_d) \mathbf{1}_n^\top}{d} + \frac{1}{4} \mathbf{A}(\mathbf{x}, \mathbf{X}) \mathbf{A}(\mathbf{X}, \mathbf{x}) + h(2\tau)^2 \mathbf{1} \mathbf{1}^\top \right) [\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} \right) \\ &= \frac{\beta^2 \|\Sigma_d\|_2}{d} \sum_{i=1}^n \frac{\lambda_i(\mathbf{X} \mathbf{X}^\top / d)}{[\lambda_i(\widetilde{\mathbf{K}}^{\text{lin}})]^2} - \frac{\beta^2 \mu_3 \lambda_1(\mathbf{X} \Sigma_d^{1/2} \text{diag}(\Sigma_d) \mathbf{1}_n^\top)}{d [\lambda_1(\widetilde{\mathbf{K}}^{\text{lin}})]^2} + \beta^2 \frac{4\lambda_1(\mathbb{E}_{\mathbf{x}}[\mathbf{A}(\mathbf{x}, \mathbf{X}) \mathbf{A}(\mathbf{X}, \mathbf{x})]) + h(2\tau)^2 n}{[\lambda_1(\widetilde{\mathbf{K}}^{\text{lin}})]^2} \quad (16) \\ &\asymp \frac{\beta^2 \|\Sigma_d\|_2}{d} \sum_{i=1}^n \frac{\lambda_i(\mathbf{X} \mathbf{X}^\top / d)}{[\lambda_1(\widetilde{\mathbf{K}}^{\text{lin}})]^2} + \frac{\mathcal{O}(\sqrt{n/d})}{[\lambda_1(\widetilde{\mathbf{K}}^{\text{lin}})]^2} + \frac{\mathcal{O}(n)}{[\lambda_1(\widetilde{\mathbf{K}}^{\text{lin}})]^2} + \frac{\mathcal{O}(n)}{[\lambda_1(\widetilde{\mathbf{K}}^{\text{lin}})]^2} \\ &\asymp \frac{\beta}{d} \mathcal{N}_{\widetilde{\mathbf{X}}}^{m\lambda+\gamma} + \mathcal{O}\left(\frac{1}{n}\right). \end{aligned}$$

It can be found that, the above error bounds are the same as that of inner-product kernels, except two additional terms due to the considered \mathbf{A} and $\mathbf{A} \odot \mathbf{A}$ in the linearization, which can be shown small in the large n, d regime.

By virtue of $\|(\mathbf{K} + n\lambda \mathbf{I})^{-1}\|_2 \leq \frac{2}{n\lambda+\gamma}$ and $\|(\mathbf{K} + n\lambda \mathbf{I})^{-1} \widetilde{\mathbf{K}}^{\text{lin}}\|_2 \leq 2$ in [18], Lemma 7, and the above equations, with probability at least $1 - \delta - d^{-2}$, for any given $\varepsilon > 0$, we have

$$\begin{aligned} & \mathbf{v} \leq \sigma^2 \mathbb{E}_{\mathbf{x}} \|(\mathbf{K} + n\lambda \mathbf{I})^{-1} k(\mathbf{x}, \mathbf{X})\|_2^2 \\ &\leq 2\sigma^2 \left\| (\mathbf{K} + n\lambda \mathbf{I})^{-1} \widetilde{\mathbf{K}}^{\text{lin}} \right\|_2^2 \mathbb{E}_{\mathbf{x}} \|[\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} k^{\text{lin}}(\mathbf{X}, \mathbf{x})\|_2^2 + 2\sigma^2 \|\mathbf{K}^{-1}\|_2^2 \mathbb{E}_{\mathbf{x}} \|k(\mathbf{x}, \mathbf{X}) - k^{\text{lin}}(\mathbf{x}, \mathbf{X})\|_2^2 \\ &\leq 8\sigma^2 \mathbb{E}_{\mathbf{x}} \|[\widetilde{\mathbf{K}}^{\text{lin}}]^{-1} k^{\text{lin}}(\mathbf{X}, \mathbf{x})\|_2^2 + \frac{8\sigma^2}{(n\lambda + \gamma)^2} \widetilde{C}_1 d^{-2\theta} \log^{1+\varepsilon} d \quad (17) \\ &\asymp \frac{\sigma^2 \beta}{d} \mathcal{N}_{\widetilde{\mathbf{X}}}^{m\lambda+\gamma} + \frac{\sigma^2}{(n\lambda + \gamma)^2} d^{-2\theta} \log^{1+\varepsilon} d, \end{aligned}$$

where the second inequality admits by Lemma 7, and the last inequality follows by Eq. (16). Finally, we conclude the proof. \square

F Proof of Proposition 1

In this section, we discuss $\mathcal{N}_{\widetilde{\mathbf{X}}}^{n\lambda+\gamma}$ based on three eigenvalue decays: *harmonic decay*, *polynomial decay*, and *exponential decay* under two regimes $n < d$ and $n > d$.

F.1 $n < d$ case

Recall $b := n\lambda + \gamma > 0$, and $\mathcal{N}_{\tilde{\mathbf{X}}}^b := \sum_{i=1}^n \frac{\lambda_i(\tilde{\mathbf{X}})}{[b + \lambda_i(\tilde{\mathbf{X}})]^2}$, define $F(\lambda_i) := \frac{\lambda_i}{(b + \lambda_i)^2}$ where λ_i is short for $\lambda_i(\tilde{\mathbf{X}})$. We notice that, when $\lambda_i \leq b$, $F(\lambda_i)$ is an increasing function of λ_i , and thus a decreasing function of i when the above three eigenvalue decays are considered. Likewise, when $\lambda_i \geq b$, $F(\lambda_i)$ is a decreasing function of λ_i , and thus an increasing function of i . Without loss of generality, we assume that the first q eigenvalues satisfy $\lambda_i \geq b$ with $i = 1, 2, \dots, q$ and the remaining $n - q$ eigenvalues satisfy $\lambda_i \leq b$ with $i = m + 1, m + 2, \dots, n$. Clearly, the integer q can be chosen from 0 to n . Accordingly, denote $r_* := \text{rank}(\tilde{\mathbf{X}})$ which includes the rank-deficient case, $\mathcal{N}_{\tilde{\mathbf{X}}}^b$ can be upper bounded by the Riemann sum as follows.

Harmonic decay $\lambda_i(\tilde{\mathbf{X}}) \propto n/i$ for $i \in \{1, 2, \dots, r_*\}$ and $\lambda_i(\tilde{\mathbf{X}}) = 0$ for $i \in \{r_* + 1, \dots, n\}$

$$\begin{aligned} \frac{1}{d} \mathcal{N}_{\tilde{\mathbf{X}}}^b &= \frac{1}{d} \sum_{i=1}^{r_*} \frac{n/i}{(b + n/i)^2} = \frac{1}{d} \sum_{i=1}^q \frac{n/i}{(b + n/i)^2} + \frac{1}{d} \sum_{i=q+1}^{r_*} \frac{n/i}{(b + n/i)^2} \\ &\leq \frac{1}{nd} \int_1^{q+1} \frac{t}{(1 + \frac{bt}{n})^2} dt + \frac{1}{nd} \int_{q+1}^{r_*+1} \frac{t}{(1 + \frac{bt}{n})^2} dt \\ &= \frac{n}{b^2 d} \int_{\frac{b}{n}}^{\frac{(r_*+1)b}{n}} \frac{u}{(1+u)^2} du \text{ with the change of variable } u = tb/n \\ &= \frac{n}{b^2 d} \left[\ln \frac{n + (r_* + 1)b}{n + b} + \frac{n}{n + b + r_* b} - \frac{n}{n + b} \right] \\ &\leq \frac{n}{b^2 d} \ln \frac{n + (r_* + 1)b}{n + b} = \mathcal{O}\left(\frac{n}{b^2 d}\right). \end{aligned}$$

Polynomial decay: $\lambda_i(\tilde{\mathbf{X}}) \propto ni^{-2a}$ with $a > 1/2$ for $i \in \{1, 2, \dots, r_*\}$ and $\lambda_i(\tilde{\mathbf{X}}) = 0$ for $i \in \{r_* + 1, \dots, n\}$. Hence, we actually aim to bound

$$\begin{aligned} \frac{1}{d} \mathcal{N}_{\tilde{\mathbf{X}}}^b &= \frac{1}{d} \sum_{i=1}^{r_*} \frac{ni^{-2a}}{(b + ni^{-2a})^2} = \frac{1}{d} \sum_{i=1}^q \frac{ni^{-2a}}{(b + ni^{-2a})^2} + \frac{1}{d} \sum_{i=q+1}^{r_*+1} \frac{ni^{-2a}}{(b + ni^{-2a})^2} \\ &\leq \frac{1}{nd} \int_1^{r_*+1} \frac{t^{2a}}{(1 + \frac{t^{2a}b}{n})^2} dt \\ &= \frac{1}{2abd} \left(\frac{n}{b}\right)^{\frac{1}{2a}} \int_{b/n}^{(r_*+1)^{2a}b/n} \frac{u^{\frac{1}{2a}}}{(1+u)^2} du \text{ with the change of variable } u = t^{2a}b/n \\ &\leq \tilde{C} \frac{1}{2abd} \left(\frac{n}{b}\right)^{\frac{1}{2a}} \text{ since the integral is finite due to } 2a > 1 \end{aligned}$$

Exponential decay: $\lambda_i(\tilde{\mathbf{X}}) \propto ne^{-ai}$ with $a > 0$ for $i \in \{1, 2, \dots, r_*\}$ and $\lambda_i(\tilde{\mathbf{X}}) = 0$ for $i \in \{r_* + 1, \dots, n\}$.

We aim to bound the sum as

$$\begin{aligned} \frac{1}{d} \mathcal{N}_{\tilde{\mathbf{X}}}^b &= \frac{1}{d} \sum_{i=1}^{r_*} \frac{ne^{-ai}}{(b + ne^{-ai})^2} = \frac{1}{d} \sum_{i=1}^q \frac{ne^{-ai}}{(b + ne^{-ai})^2} + \frac{1}{d} \sum_{i=q+1}^{r_*} \frac{ne^{-ai}}{(b + ne^{-ai})^2} \\ &\leq \frac{1}{d} \int_1^{r_*+1} \frac{ne^{-at}}{(b + ne^{-at})^2} dt \\ &= \frac{1}{ad} \int_{ne^{-a(r_*+1)}}^{ne^{-a}} \frac{1}{(b + u)^2} du \text{ with the change of variable } u = ne^{-at} \\ &= \frac{1}{ad} \left(\frac{1}{b + ne^{-a(r_*+1)}} - \frac{1}{b + ne^{-a}} \right). \end{aligned}$$

Note that, the monotonicity of $\mathcal{N}_{\tilde{\mathbf{X}}}^b$ (also \mathbf{V}_1) with respect to n is relatively clear for *harmonic decay* and *polynomial decay* but is unclear in the case of *exponential decay*. Here we study the monotonicity in the *exponential decay*.

Denote the function $G(n) := \left(\frac{1}{b+ne^{-a(r_*+1)}} - \frac{1}{b+ne^{-a}} \right)$ with $b := n\lambda + \gamma$, taking $\lambda := \bar{c}n^{-\vartheta}$, its derivation is

$$G'(n) = \frac{-\bar{c}(1-\vartheta)n^{-\vartheta} - e^{-a(r_*+1)}}{[cn^{1-\vartheta} + \gamma + ne^{-a(r_*+1)}]^2} + \frac{\bar{c}(1-\vartheta)n^{-\vartheta} + e^{-a}}{[cn^{1-\vartheta} + \gamma + ne^{-a}]^2}, \quad (18)$$

which can be rewritten as

$$G'(n) = \frac{\bar{c}(1-\vartheta)n^{-\vartheta} + e^{-a}}{[cn^{1-\vartheta} + \gamma + ne^{-a(r_*+1)}]^2} \left(\underbrace{\frac{[\bar{c}n^{1-\vartheta} + \gamma + ne^{-a(r_*+1)}]^2}{[\bar{c}n^{1-\vartheta} + \gamma + ne^{-a}]^2}}_{\triangleq H_1(n)} - \underbrace{\frac{\bar{c}(1-\vartheta)n^{-\vartheta} + e^{-a(r_*+1)}}{\bar{c}(1-\vartheta)n^{-\vartheta} + e^{-a}}}_{\triangleq H_2(n)} \right).$$

It can be found that both $H_1(n)$ and $H_2(n)$ are decreasing functions with n . More specifically, their maximum and minimum can be achieved with

$$\max_n H_1(n) = H_1(1) = \left(\frac{\bar{c} + \gamma + e^{-a(r_*+1)}}{\bar{c} + \gamma + e^{-a}} \right)^2, \quad \min_n H_1(n) = \lim_{n \rightarrow \infty} H_1(n) = \left(\frac{e^{-a(r_*+1)}}{e^{-a}} \right)^2,$$

and

$$\max_n H_2(n) = H_2(1) = \frac{\bar{c}(1-\vartheta) + e^{-a(r_*+1)}}{\bar{c}(1-\vartheta) + e^{-a}}, \quad \min_n H_2(n) = \lim_{n \rightarrow \infty} H_2(n) = \frac{e^{-a(r_*+1)}}{e^{-a}}.$$

Accordingly, if $H_1(1) < H_2(1)$, we obtain a decreasing function $G(n)$ of n , which implies that $\mathcal{N}_{\tilde{\mathbf{X}}}^b$ will decrease with n . Here the condition $H_1(1) < H_2(1)$ indicates

$$\left(\frac{\bar{c} + \gamma + e^{-a(r_*+1)}}{\bar{c} + \gamma + e^{-a}} \right)^2 \leq \frac{\bar{c}(1-\vartheta) + e^{-a(r_*+1)}}{\bar{c}(1-\vartheta) + e^{-a}},$$

which is equivalent to

$$(\vartheta\bar{c} + \gamma)^2 \leq [e^{-a} + (1-\vartheta)\bar{c}] [e^{-a(r_*+1)} + (1-\vartheta)\bar{c}]. \quad (19)$$

Accordingly, if the above inequality holds, $\mathcal{N}_{\tilde{\mathbf{X}}}^b$ will decrease with n . In Section G.2, we will experimentally check whether this condition holds or not.

F.2 $n > d$ case and the large n limit

In this section, we consider the $n > d$ case, and further study the trend of \mathbf{V}_1 as $n \rightarrow \infty$. Note that, in this case, $\mathbf{X}\mathbf{X}^\top/d$ has at most $r_* \leq d$ non-zero eigenvalues. Accordingly, the Riemann sum is counted to r_* instead of n . Similar to the above description, we also consider the following three eigenvalue decays.

Harmonic decay $\lambda_i(\tilde{\mathbf{X}}) \propto n/i$, $i \in \{1, 2, \dots, d\}$

$$\begin{aligned} \frac{1}{d} \mathcal{N}_{\tilde{\mathbf{X}}}^b &= \frac{1}{d} \sum_{i=1}^{r_*} \frac{n/i}{(b+n/i)^2} = \frac{1}{d} \sum_{i=1}^q \frac{n/i}{(b+n/i)^2} + \frac{1}{d} \sum_{i=q+1}^{r_*} \frac{n/i}{(b+n/i)^2} \\ &\leq \frac{n}{b^2 d} \int_{\frac{b}{n}}^{\frac{(r_*+1)b}{n}} \frac{u}{(1+u)^2} du \\ &= \frac{n}{b^2 d} \left[\ln \frac{n + (r_*+1)b}{n+b} + \frac{n}{n+b+r_*b} - \frac{n}{n+b} \right]. \end{aligned}$$

In particular, taking the limit of $n \rightarrow \infty$, we have

$$\begin{aligned}
 \lim_{n \rightarrow \infty} \frac{1}{d} \mathcal{N}_{\tilde{\mathbf{X}}}^{rb} &= \lim_{n \rightarrow \infty} \frac{n}{b^2 d} \left[\ln \frac{n + (r_* + 1)b}{n + b} + \frac{n}{n + b + r_* b} - \frac{n}{n + b} \right] \\
 &= \lim_{n \rightarrow \infty} \frac{n}{b^2 d} \ln \frac{n + (r_* + 1)b}{n + b} + \lim_{n \rightarrow \infty} \frac{n}{b^2 d} \left(\frac{n}{n + b + r_* b} - \frac{n}{n + b} \right) \\
 &= \frac{r_*}{d} \left(\lim_{n \rightarrow \infty} \frac{1}{b} \frac{n}{n + b} - \lim_{n \rightarrow \infty} \frac{n^2}{b(n + b + r_*)(n + b)} \right) \\
 &\leq \lim_{n \rightarrow \infty} \frac{1}{b} \frac{n}{n + b} - \lim_{n \rightarrow \infty} \frac{n^2}{b(n + b + r_*)(n + b)} \\
 &= 0.
 \end{aligned}$$

Accordingly, by the squeeze theorem, we can conclude, given d , $\mathcal{N}_{\tilde{\mathbf{X}}}^{rb}$ tends to zero when $n \rightarrow \infty$.

Polynomial decay: $\lambda_i(\tilde{\mathbf{X}}) \propto ni^{-2a}$ with $a > 1/2$, $i \in \{1, 2, \dots, d\}$

$$\begin{aligned}
 \frac{1}{d} \mathcal{N}_{\tilde{\mathbf{X}}}^{rb} &= \frac{1}{d} \sum_{i=1}^{r_*} \frac{ni^{-2a}}{(b + ni^{-2a})^2} \leq \frac{1}{2abd} \left(\frac{n}{b}\right)^{\frac{1}{2a}} \int_{b/n}^{(r_*+1)^{2a}b/n} \frac{u^{\frac{1}{2a}}}{(1+u)^2} du \\
 &\leq \frac{1}{2abd} \left(\frac{n}{b}\right)^{\frac{1}{2a}} \int_0^\infty \frac{u^{\frac{1}{2a}}}{(1+u)^2} du \\
 &\leq \tilde{C} \frac{1}{2abd} \left(\frac{n}{b}\right)^{\frac{1}{2a}} \quad \text{since the integral is finite due to } 2a > 1
 \end{aligned}$$

Since the integral $\int \frac{u^{\frac{1}{2a}}}{(1+u)^2} du$ can behave rather differently for different choices of a , here we take $a = 1$ as an example. Taking the limit of $n \rightarrow \infty$, we have

$$\begin{aligned}
 \lim_{n \rightarrow \infty} \frac{1}{d} \mathcal{N}_{\tilde{\mathbf{X}}}^{rb} &= \lim_{n \rightarrow \infty} \frac{1}{2bd} \left(\frac{n}{b}\right)^{\frac{1}{2}} \int_{b/n}^{(r_*+1)^2b/n} \frac{u^{\frac{1}{2}}}{(1+u)^2} du \\
 &= \frac{1}{2bd} \lim_{n \rightarrow \infty} \left(\frac{n}{b}\right)^{\frac{1}{2}} \left(\arctan(\sqrt{u}) - \frac{\sqrt{u}}{u+1} \right) \Big|_{b/n}^{(r_*+1)^2b/n} \\
 &= \frac{1}{2bd} \lim_{n \rightarrow \infty} \sqrt{\frac{n}{b}} \left((r_* + 1)\sqrt{b/n} - \frac{(r_* + 1)\sqrt{b/n}}{(r_* + 1)^2b/n} - \sqrt{b/n} + \frac{\sqrt{b/n}}{b/n + 1} \right) \quad \text{using } \lim_{x \rightarrow 0} \frac{\arctan x}{x} = 1. \\
 &= 0.
 \end{aligned}$$

Exponential decay: $\lambda_i(\tilde{\mathbf{X}}) \propto ne^{-ai}$ with $a > 0$, $i \in \{1, 2, \dots, d\}$

$$\begin{aligned}
 \frac{1}{d} \mathcal{N}_{\tilde{\mathbf{X}}}^{rb} &= \frac{1}{d} \sum_{i=1}^{r_*} \frac{ne^{-ai}}{(b + ne^{-ai})^2} = \frac{1}{d} \sum_{i=1}^q \frac{ne^{-ai}}{(b + ne^{-ai})^2} + \frac{1}{d} \sum_{i=q+1}^{r_*} \frac{ne^{-ai}}{(b + ne^{-ai})^2} \\
 &\leq \frac{1}{ad} \int_{ne^{-a(r_*+1)}}^{ne^{-a}} \frac{1}{(b + u)^2} du \\
 &= \frac{1}{ad} \left(\frac{1}{b + ne^{-a(r_*+1)}} - \frac{1}{b + ne^{-a}} \right).
 \end{aligned}$$

Taking the limit of $n \rightarrow \infty$, we can directly have $\lim_{n \rightarrow \infty} \frac{1}{d} \mathcal{N}_{\tilde{\mathbf{X}}}^{rb} = 0$.

G Additional Experiments

In this section, we present additional experiments including the following parts:

- In Section G.1, we add the *MNIST* dataset [59] to verify the eigenvalue decay equivalence, and evaluate the effect by different orders in polynomial kernel.
- In Section G.2, our model works in a polynomial kernel setting under the *polynomial decay* and *exponential decay* of $\widetilde{\mathbf{X}}$ on the synthetic dataset.

G.1 Eigenvalue decay equivalence

Apart from the *YearPredictionMSD* dataset in the main text, we add the *MNIST* dataset [59] to verify the eigenvalue decay equivalence. We also compute eigenvalues of $\widetilde{\mathbf{X}} := \beta \mathbf{X} \mathbf{X}^\top / d + \alpha \mathbf{1} \mathbf{1}^\top$ for validation. Here the parameters α depends on the covariate Σ_d , which can be empirically estimated by the sample covariance $\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j)(\mathbf{x}_i - \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j)^\top$.

Results on the polynomial kernel with order 3 and the Gaussian kernel are presented in Figure 6 and 7, respectively. It can be observed that, the nonlinear kernel matrix \mathbf{K} admits almost the same eigenvalue as $\widetilde{\mathbf{X}} := \beta \mathbf{X} \mathbf{X}^\top / d + \alpha \mathbf{1} \mathbf{1}^\top$ with a constant shift γ , and accordingly exhibits the same eigenvalue decay with $\widetilde{\mathbf{X}}$ and $\mathbf{X} \mathbf{X}^\top / d$.

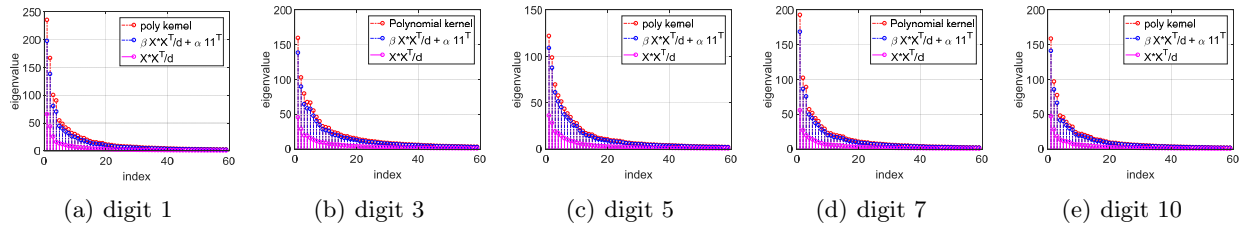


Figure 6: Top 60 eigenvalues of Polynomial kernel with order 3 and its linearization on the MNIST dataset. Note that the largest eigenvalue λ_1 is not plotted for better display.

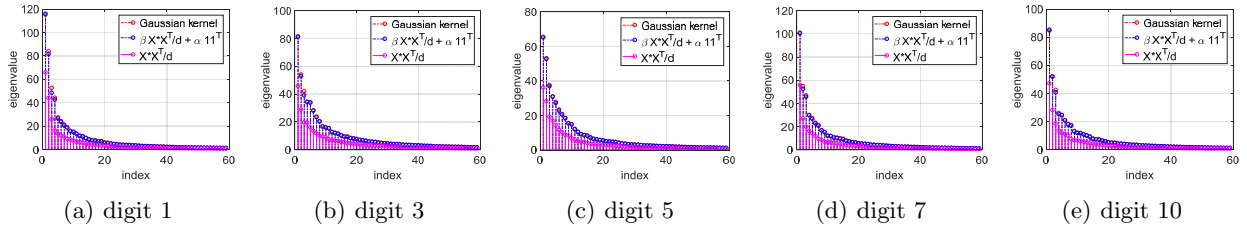


Figure 7: Top 60 eigenvalues of Gaussian kernel and its linearization on the MNIST dataset. Note that the largest eigenvalue λ_1 is not plotted for better display.

Besides, to study eigenvalue decay effected by the order in polynomial kernels, we present results of the order $p = 5$ and $p = 10$ in Figure 8. Experimental results show that, there is some gap between the original kernel and its linearization in higher orders. This is because, nonlinear kernel approximated by linear model here is based on Taylor expansion, which would incur in some residual errors as higher order in polynomial kernels brings in stronger non-linearity.

G.2 Results on the synthetic dataset

Here we evaluate our model with the polynomial kernel on the synthetic dataset under the polynomial/exponential decay of Σ_d . The data generation process follows with our experiments part in the main text such that $\widetilde{\mathbf{X}}$ admits the polynomial/exponential decay.

Results on the *polynomial decay* and the *exponential decay* are shown in Figure 9 and Figure 10, respectively. We find that, the bias achieves the certain $\mathcal{O}(n^{-2\theta r})$ convergence rate on both decays; while the variance shows different configurations on these two decays. To be specific, the tend of V_1 on the *polynomial decay* is unimodal,

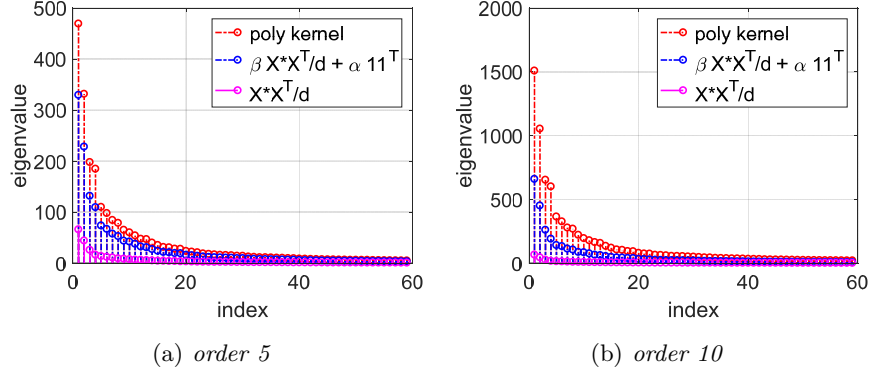


Figure 8: Top 60 eigenvalues of polynomial kernel matrices and their linearizations on the MNIST dataset (digit 1). Note that the largest eigenvalue λ_1 is not plotted for better display.

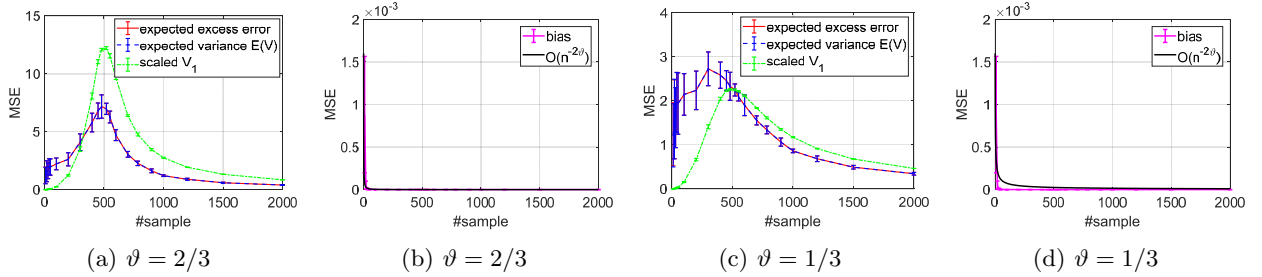


Figure 9: Polynomial decay of $\widetilde{\mathbf{X}}$ in the polynomial kernel case: MSE of the expected excess risk, the variance in Eq. (10), our derived \mathbf{V}_1 , the bias in Eq. (9), and our derived convergence rate $\mathcal{O}(n^{-2\vartheta r})$ with $r = 1$ in Theorem 2 under different ϑ .

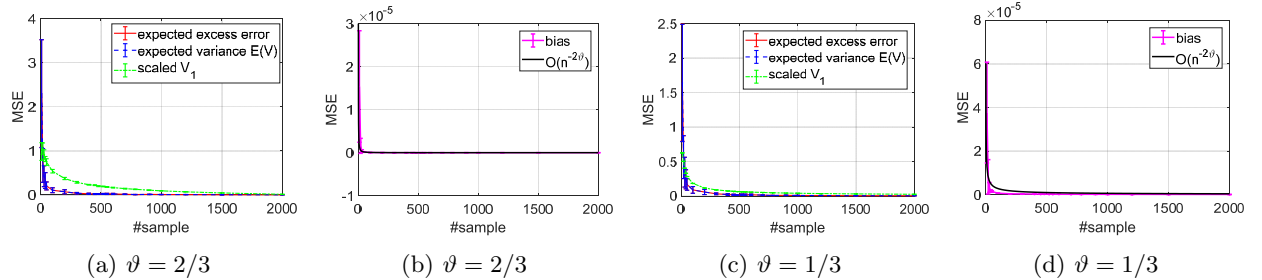


Figure 10: Exponential decay of $\widetilde{\mathbf{X}}$ in the polynomial kernel case: MSE of the expected excess risk, the variance in Eq. (10), our derived \mathbf{V}_1 , the bias in Eq. (9), and our derived convergence rate $\mathcal{O}(n^{-2\vartheta r})$ with $r = 1$ in Theorem 2 under different ϑ .

and thus the risk curve is bell-shaped. However, in Figure 10, \mathbf{V}_1 on the *exponential decay* monotonically decreases with n even if we set \bar{c} to 10^{-5} , 10^{-8} for a small regularization scheme.

Here we attempt to explain this phenomenon. In our setting, γ is set to zero. The condition in Eq. (19) can be reformulated as

$$(2\vartheta - 1)\bar{c} \leq e^{-a}(1 - \vartheta).$$

Clearly, if we choose $0 < \vartheta < 1/2$, the condition in Eq. (19) always holds. Hence, \mathbf{V}_1 will monotonically decrease with n . If $1/2 < \vartheta < 1$, we examine our result with $a = 1$ and $\vartheta = 2/3$. We conclude that the used $\bar{c} = 0.01 < e^{-1}$, so the tend of \mathbf{V}_1 is monotonically decreasing with n .