# Appendix

## A Glossary

Table 2: Glossary

| Name | Notation | Expression | Dimension |
|---|---|---|---|
| sampling distribution | $\rho$ | - | $\mathcal{X} \to \mathbb{R}^+$ |
| sampling size | $N$ | - | integer |
| input matrix | $\mathbf{X}$ | $(x_i)_{i=1}^N \underset{iid}{\sim} \rho_{\mathcal{X}}$ | $N \times d$ |
| output vector | $\mathbf{y}$ | $(y_i)_{i=1}^N$ | $N \times 1$ |
| sample | $\mathbf{Z}$ | $(\mathbf{X}, \mathbf{y})$ | $N \times (d+1)$ |
| noise | $\varepsilon$ | - | random scalar |
| noise variance | $\sigma^2$ | $\mathbb{E}[\varepsilon^2]$ | scalar |
| ridge | $\lambda$ | - | scalar |
| finite-rank kernel | $K$ | $\sum_{k=1}^M \lambda_k \psi_k(\cdot)\psi_k(\cdot)$ | $\mathcal{X} \times \mathcal{X} \to \mathbb{R}$ |
| kernel rank | $M$ | - | integer |
| $k$th eigenfunction | $\psi_k$ | - | $\mathcal{X} \to \mathbb{R}$ |
| $k$th value | $\lambda_k$ | - | scalar |
| - | $\boldsymbol{\psi}(x)$ | $[\psi_k(x)]_{k=1}^M$ | $M \times 1$ |
| - | $\boldsymbol{\Psi}$ | $[\psi_k(x_i)]_{k,i}$ | $M \times N$ |
| - | $\boldsymbol{\Lambda}$ | $\operatorname{diag}\left[\lambda_k\right]$ | $M \times M$ |
| kernel matrix | $\mathbf{K}$ | $[K(x_i, x_j)]_{i,j} = \boldsymbol{\Psi}^\top \boldsymbol{\Lambda} \boldsymbol{\Psi}$ | $N \times N$ |
| resolvent | $\mathbf{R}$ | $(\mathbf{K} + \lambda N \mathbf{I}_N)^{-1}$ | $N \times N$ |
| target function | $\tilde{f}$ | $\sum_{k=1}^M \tilde{\gamma}_k \psi_k + \tilde{\gamma}_{>M} \psi_{>M}$ | $\mathcal{X} \to \mathbb{R}$ |
| - | $\tilde{f}_{\leq M}$ | $\sum_{k=1}^M \tilde{\gamma}_k \psi_k$ | $\mathcal{X} \to \mathbb{R}$ |
| $k$th target coefficient | $\tilde{\gamma}_k$ | $\int_{\mathcal{X}} \tilde{f}(x)\psi_k(x)d\rho_{\mathcal{X}}(x)$ | scalar |
| - | $\boldsymbol{\gamma}$ | $[\gamma_k]$ | $M \times 1$ |
| orthonormal complement | $\psi_{>M}$ | - | $\mathcal{X} \to \mathbb{R}$ |
| complementary coefficient | $\tilde{\gamma}_{>M}$ | - | scalar |
| - | $\boldsymbol{\Psi}_{>M}$ | $[\psi_{>M}(x_i)]$ | $1 \times N$ |
| test error | $\mathcal{R}_{\mathbf{Z}, \lambda}$ | $\mathbb{E}_{x,\epsilon}\left[(f_{\mathbf{Z}, \lambda}(x) - \tilde{f}(x))^2\right]$ | scalar |
| bias | - | $\int_{\mathcal{X}}\left(f_{(\mathbf{X}, \tilde{f}(\mathbf{X})), \lambda}(x) - \tilde{f}(x)\right)^2 d\rho(x)$ | scalar |
| variance | - | $\mathcal{R}_{\mathbf{Z}, \lambda} - \operatorname{bias} = \mathbb{E}_{x,\varepsilon}\left(\mathbf{K}_x^\top \mathbf{R} \boldsymbol{\varepsilon}\right)^2$ | scalar |
| fluctuation matrix | $\boldsymbol{\Delta}$ | $\frac{1}{N}\boldsymbol{\Psi}\boldsymbol{\Psi}^\top - \mathbf{I}_M$ | $M \times M$ |
| fluctuation | $\delta$ | $\|\boldsymbol{\Delta}\|_{\operatorname{op}}$ | scalar |
| error vector | $\boldsymbol{E}$ | $[\eta_k]$ | $M \times 1$ |
| - | $\eta_k$ | $\frac{1}{N}\sum_{i=1}^N \psi_k(x_i)\psi_{>M}(x_i)$ | scalar |
| - | $\mathbf{B}$ | $(\mathbf{I}_M + \boldsymbol{\Delta} + \lambda \boldsymbol{\Lambda}^{-1})^{-1}$ | $M \times M$ |
| - | $\bar{\mathbf{P}}$ | $\operatorname{diag}\left[\frac{\lambda_k}{\lambda_k + \lambda}\right]$ | $M \times M$ |

## B Classical KRR Theory

In an effort to keep our manuscript as self-contained as possible, we recall the Mercer decomposition, representer theorem for kernel ridge regression as well as the form of the bias-variance tradeoff in the KRR context.

### B.1 Mercer Decomposition

We begin with a general kernel $K^{(\infty)} : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$.

**Proposition B.1.** *[12] Fix a sample distribution $\rho$. Let $K^{(\infty)} : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ be a reproducing kernel with corresponding RKHS $\mathcal{H}^{(\infty)}$. There exists a decreasing sequence of real numbers $\lambda_1 \geq \lambda_2 \geq ...$, called the eigenvalues of the kernel $K^{(\infty)}$; and a sequence of pairwise-orthonormal functions $\{\psi_k\}_{k=1}^{\infty} \subset L_\rho^2$, called the eigenfunctions of $K^{(\infty)}$, such that for all $x, x' \in \mathcal{X}$, we have*

$$K^{(\infty)}(x, x') = \sum_{k=1}^{\infty} \lambda_k \psi_k(x) \psi_k(x') \tag{14}$$

In particular, we assume $\lambda_k = 0$, $\forall k > M$. In this case, we say the kernel $K(x, x') = \sum_{k=1}^{M} \lambda_k \psi_k(x) \psi_k(x')$ is of finite rank $M$ with corresponding (finite-dimensional) RKHS $\mathcal{H}$, re-covering equation (2).

The first of these results, allows us to explicitly express the finite-rank kernel ridge regressor $f_{\mathbf{Z},\lambda}$.

**Proposition B.2** (Representer Theorem - [38, Chapter 12])*. Let $\mathbf{R} \stackrel{\text{def.}}{=} (\mathbf{K} + \lambda N \mathbf{I}_N)^{-1} \in \mathbb{R}^{N \times N}$ be the resolvent matrix and recall the kernel ridge regressor $f_{\mathbf{Z},\lambda}$ given by equation (3):*

$$f_{\mathbf{Z},\lambda} \stackrel{\text{def.}}{=} \arg\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^{N} (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}}^2$$

*Then, for every $x \in \mathcal{X}$, we have the expression*

$$f_{\mathbf{Z},\lambda}(x) = \mathbf{y}^\top \mathbf{R} \mathbf{K}_x, \ \forall x \in \mathcal{X}, \tag{15}$$

*where $\mathbf{K}_x \stackrel{\text{def.}}{=} [K(x_i, x)]_{i=1}^{N} \in \mathbb{R}^{N \times 1}$.*

### B.2 Compact Matrix Expression

First, let $\mathbf{\Psi} \stackrel{\text{def.}}{=} (\psi_k(x_i))_{k=1,i=1}^{M,N}$ be the random $M \times N$ matrix defined by evaluating the $M$ eigenfunctions on all input training instances $\mathbf{X} \stackrel{\text{def.}}{=} (x_i)_{i=1}^{N}$, $\mathbf{\Lambda} \stackrel{\text{def.}}{=} \mathrm{diag}[\lambda_k] \in \mathbb{R}^{M \times M}$, and $\psi(x) \stackrel{\text{def.}}{=} [\psi_k(x)]_{k=1}^{M} \in \mathbb{R}^{M \times 1}$. The advantage of this notation is that we can rewrite the equations in a more compact form. For equation (15):

$$f_{\mathbf{Z},\lambda}(x) = \mathbf{y}^\top \underbrace{(\mathbf{\Psi}^\top \mathbf{\Lambda} \mathbf{\Psi} + \lambda N \mathbf{I}_M)^{-1}}_{\mathbf{R}} \mathbf{\Psi}^\top \mathbf{\Lambda} \psi(x); \tag{16}$$

for equation (4):

$$\tilde{f}(x) = \tilde{\boldsymbol{\gamma}}^\top \psi(x) + \tilde{\gamma}_{>M} \psi_{>M}(x). \tag{17}$$

Last but not least, we define some important quantities for later analysis.

**Definition B.3** (Fluctuation matrix)*. The fluctuation matrix is the random $M \times M$-matrix given by $\mathbf{\Delta} \stackrel{\text{def.}}{=} \frac{1}{N} \mathbf{\Psi} \mathbf{\Psi}^\top - \mathbf{I}_M$. Our analysis will often involve the operator norm of $\mathbf{\Delta}$, which we denote by $\delta \stackrel{\text{def.}}{=} \|\mathbf{\Delta}\|_{op}$.*

The fluctuation matrix $\mathbf{\Delta}$ measures the first source of randomness in the KRR's test error. Namely it encodes the degree of non-orthonormality between the vectors obtained by evaluating of the $M$ eigenfunctions $\psi_1, \ldots, \psi_M$ on the input $\mathbf{X}$.

The second source of randomness in the KRR's test error comes from the empirical evaluation of the dot product of the eigenfunction $\psi_k$'s and the orthogonal complement $\psi_{>M}$:

**Definition B.4** (Error Vector)*. $\boldsymbol{E} \stackrel{\text{def.}}{=} \frac{1}{N} \mathbf{\Psi} \psi_{>M}(\mathbf{X})$ is called the error vector.*

The random matrix $\mathbf{\Delta}$ and the random vector $\boldsymbol{E}$ are centered; i.e. $\mathbb{E}_{\mathbf{X}}[\mathbf{\Delta}] = 0$ and $\mathbb{E}_{\mathbf{X}}[\boldsymbol{E}] = 0$.

## B.3 Bias-Variance Decomposition

396 The derivation of several contemporary KRR generalization bounds [6, 26, 27] involves the classical
397 Bias-Variance Trade-off:

398 **Proposition B.5** (Bias-Variance Trade-off). *Fix a sample* $\mathbf{Z}$*. Recall the definition 3.4 of test error*
399 $\mathcal{R}_{\mathbf{Z},\lambda}$*, bias, and variance:*

$$\mathcal{R}_{\mathbf{Z},\lambda} \overset{\text{def.}}{=} \mathbb{E}_{x,\epsilon}\left[(f_{\mathbf{Z},\lambda}(x) - \tilde{f}(x))^2\right] = \mathbb{E}_\epsilon\left[\int_{\mathcal{X}}\left(f_{\mathbf{Z},\lambda}(x) - \tilde{f}(x)\right)^2 d\rho(x)\right];$$

$$bias \overset{\text{def.}}{=} \int_{\mathcal{X}}\left(f_{(\mathbf{X},\tilde{f}(\mathbf{X})),\lambda}(x) - \tilde{f}(x)\right)^2 d\rho(x);$$

$$variance \overset{\text{def.}}{=} \mathcal{R}_{\mathbf{Z},\lambda} - bias.$$

400 *Then, we can write* $variance_{test} = \mathbb{E}_{x,\varepsilon}\left(\mathbf{K}_x^\top \mathbf{R}\boldsymbol{\varepsilon}\right)^2$ *and hence the test error* $\mathcal{R}_{\mathbf{Z},\lambda}$ *admits a decompo-*
401 *sition:*

$$R_{\mathbf{Z},\lambda} = bias + \mathbb{E}_{x,\varepsilon}\left(\mathbf{K}_x^\top \mathbf{R}\boldsymbol{\varepsilon}\right)^2.$$

402 *Proof.* See the proof of Theorem C.8. □

403 # C  Proofs

404 In this section, we will derive the essential lemmata and propositions for proving the main theorems.

405 ## C.1  Formula Derivation

406 We begin with writing the test error in convenient forms.

407 ### C.1.1  Bias

408 We first derive, from the definition of the bias, a convenient expression to proceed:

409 **Proposition C.1** (Bias Expression). *Let* $\boldsymbol{\Psi}_{>M} \overset{\text{def.}}{=} [\psi_{>M}(x_i)]_{i=1}^N$ *as an* $1 \times N$*- row vector,* $\left(\begin{smallmatrix}\boldsymbol{\Psi}\\\boldsymbol{\Psi}_{>M}\end{smallmatrix}\right)$
410 *as an* $(M+1) \times N$ *matrix. Denote* $\mathbf{P} \overset{\text{def.}}{=} \left(\mathbf{P}_{\leq M} \quad \mathbf{P}_{>M}\right) = \boldsymbol{\Lambda}\boldsymbol{\Psi}\mathbf{R}\left(\boldsymbol{\Psi}^\top \quad \boldsymbol{\Psi}_{>M}^\top\right) \in \mathbb{R}^{M \times (M+1)}$
411 *,* $\mathbf{P}_{\leq M} \in \mathbb{R}^{M \times M}$ *and* $\mathbf{P}_{>M} \in \mathbb{R}^{M \times 1}$*. Then the bias admits the following expression:*

$$bias = \underbrace{\tilde{\gamma}_{>M}^2}_{\text{Finite Rank Error}} + \underbrace{\|\tilde{\boldsymbol{\gamma}} - \mathbf{P}_{\leq M}\tilde{\boldsymbol{\gamma}} - \tilde{\gamma}_{>M}\mathbf{P}_{>M}\|_2^2}_{\text{Fitting Error}}.$$

412 *Proof.* Recall that, by equations (16) and (17), we can write

$$\tilde{f}(x) = \tilde{\boldsymbol{\gamma}}^\top \boldsymbol{\psi}(x) + \tilde{\gamma}_{>M}\psi_{>M}(x),$$

$$f_{(\mathbf{X},\tilde{f}(\mathbf{X})),\lambda}(x) = (\tilde{\boldsymbol{\gamma}}^\top \boldsymbol{\Psi} + \tilde{\gamma}_{>M}\boldsymbol{\Psi}_{>M}^\top)\mathbf{R}\boldsymbol{\Psi}^\top \boldsymbol{\Lambda}\boldsymbol{\Psi}(x).$$

413 Hence

$$
\begin{aligned}
\text{bias} &= \mathbb{E}_x\left[\left(\tilde{\boldsymbol{\gamma}}^\top \boldsymbol{\psi}(x) + \tilde{\gamma}_{>M}\psi_{>M}(x) - (\tilde{\boldsymbol{\gamma}}^\top \boldsymbol{\Psi} + \tilde{\gamma}_{>M}\boldsymbol{\Psi}_{>M}^\top)\mathbf{R}\boldsymbol{\Psi}^\top \boldsymbol{\Lambda}\boldsymbol{\Psi}(x)\right)^2\right]\\
&= \left\|\begin{pmatrix}\tilde{\boldsymbol{\gamma}}\\\tilde{\gamma}_{>M}\end{pmatrix} - \begin{pmatrix}\mathbf{P}\begin{pmatrix}\tilde{\boldsymbol{\gamma}}\\\tilde{\gamma}_{>M}\end{pmatrix}\\0\end{pmatrix}\right\|_2^2 \quad\quad (18)\\
&= \underbrace{\tilde{\gamma}_{>M}^2}_{\text{Finite Rank Error}} + \underbrace{\|\tilde{\boldsymbol{\gamma}} - \mathbf{P}_{\leq M}\tilde{\boldsymbol{\gamma}} - \tilde{\gamma}_{>M}\mathbf{P}_{>M}\|_2^2}_{\text{Fitting Error}},
\end{aligned}
$$

414 in line (18), we use Parseval's identity. □

415 We proceed by reformulating the projection matrix $\mathbf{P}$, first with the left matrix $\mathbf{P}_{\leq M}$:

**Lemma C.2.** *Recall the following notations*

$$\mathbf{K} \overset{\text{def.}}{=} \boldsymbol{\Psi}^\top \boldsymbol{\Lambda} \boldsymbol{\Psi}, \quad \mathbf{R} \overset{\text{def.}}{=} (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1}, \quad \mathbf{P}_{\leq M} \overset{\text{def.}}{=} \boldsymbol{\Lambda} \boldsymbol{\Psi} \mathbf{R} \boldsymbol{\Psi}^\top.$$

*Define the symmetric random matrix* $\mathbf{B} \overset{\text{def.}}{=} (\mathbf{I}_M + \boldsymbol{\Delta} + \lambda \boldsymbol{\Lambda}^{-1})^{-1}$. *It holds that*

$$\mathbf{P}_{\leq M} = \mathbf{I}_M - \lambda \mathbf{B} \boldsymbol{\Lambda}^{-1}.$$

*Proof.* We first observe that

$$\boldsymbol{\Psi} \boldsymbol{\Psi}^\top \mathbf{P}_{\leq M} = \boldsymbol{\Psi} \boldsymbol{\Psi}^\top \boldsymbol{\Lambda} \boldsymbol{\Psi} \mathbf{R} \boldsymbol{\Psi}^\top \tag{19}$$

$$= \boldsymbol{\Psi} \mathbf{K} (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1} \boldsymbol{\Psi}^\top$$

$$= \boldsymbol{\Psi} \left( \mathbf{I}_M - \lambda N (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1} \right) \boldsymbol{\Psi}^\top$$

$$= \boldsymbol{\Psi} \boldsymbol{\Psi}^\top - \lambda N \boldsymbol{\Psi} (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1} \boldsymbol{\Psi}^\top. \tag{20}$$

From lines (19)- (20) and the definition of the fluctuation matrix $\boldsymbol{\Delta}$ we deduce

$$\frac{1}{N} \boldsymbol{\Psi} \boldsymbol{\Psi}^\top (\mathbf{I}_M - \mathbf{P}_{\leq M}) = \lambda \boldsymbol{\Psi} (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1} \boldsymbol{\Psi}^\top$$

$$(\mathbf{I}_M + \boldsymbol{\Delta})(\mathbf{I}_M - \mathbf{P}_{\leq M}) = \lambda \boldsymbol{\Psi} \mathbf{R} \boldsymbol{\Psi}^\top$$

$$(\mathbf{I}_M + \boldsymbol{\Delta})(\mathbf{I}_M - \mathbf{P}_{\leq M}) = \lambda \boldsymbol{\Lambda}^{-1} \mathbf{P}_{\leq M}$$

$$(\boldsymbol{\Lambda} + \boldsymbol{\Lambda} \boldsymbol{\Delta})(\mathbf{I}_M - \mathbf{P}_{\leq M}) = \lambda \mathbf{P}_{\leq M}$$

$$\boldsymbol{\Lambda} + \boldsymbol{\Lambda} \boldsymbol{\Delta} = (\boldsymbol{\Lambda} + \boldsymbol{\Lambda} \boldsymbol{\Delta} + \lambda \mathbf{I}_M) \mathbf{P}_{\leq M}. \tag{21}$$

Rearranging (21) and applying the definition of $\boldsymbol{B}$ we find that

$$\mathbf{P}_{\leq M} = (\boldsymbol{\Lambda} + \boldsymbol{\Lambda} \boldsymbol{\Delta} + \lambda \mathbf{I}_M)^{-1} (\boldsymbol{\Lambda} + \boldsymbol{\Lambda} \boldsymbol{\Delta}) \tag{22}$$

$$= \mathbf{I}_M - \lambda (\boldsymbol{\Lambda} + \boldsymbol{\Lambda} \boldsymbol{\Delta} + \lambda \mathbf{I}_M)^{-1}$$

$$= \mathbf{I}_M - \lambda (\boldsymbol{\Lambda} + \boldsymbol{\Lambda} \boldsymbol{\Delta} + \lambda \mathbf{I}_M)^{-1} \boldsymbol{\Lambda} \boldsymbol{\Lambda}^{-1}$$

$$= \mathbf{I}_M - \lambda \mathbf{B} \boldsymbol{\Lambda}^{-1}. \tag{23}$$

$\square$

Arguing analogously for the right matrix $\mathbf{P}_{>M}$, we draw the subsequent similar conclusion.

**Lemma C.3.** *Recall the following notations*

$$\mathbf{K} \overset{\text{def.}}{=} \boldsymbol{\Psi}^\top \boldsymbol{\Lambda} \boldsymbol{\Psi}, \quad \mathbf{R} \overset{\text{def.}}{=} (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1}, \quad \mathbf{P}_{>M} \overset{\text{def.}}{=} \boldsymbol{\Lambda} \boldsymbol{\Psi} \mathbf{R} \boldsymbol{\Psi}^\top_{>M},$$

$$\boldsymbol{E} \overset{\text{def.}}{=} \frac{1}{N} \boldsymbol{\Psi} \boldsymbol{\Psi}^\top_{>M}, \quad \mathbf{B} \overset{\text{def.}}{=} (\mathbf{I}_M + \boldsymbol{\Delta} + \lambda \boldsymbol{\Lambda}^{-1})^{-1}.$$

*We have that* $\mathbf{P}_{>M} = \mathbf{B} \boldsymbol{E}$.

*Proof.* Similarly to (19)- (20) we note that

$$\boldsymbol{\Psi} \boldsymbol{\Psi}^\top \mathbf{P}_{>M} = \boldsymbol{\Psi} \boldsymbol{\Psi}^\top \boldsymbol{\Lambda} \boldsymbol{\Psi} \mathbf{R} \boldsymbol{\Psi}^\top_{>M}$$

$$= \boldsymbol{\Psi} \mathbf{K} (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1} \boldsymbol{\Psi}^\top_{>M}$$

$$= \boldsymbol{\Psi} \left( \mathbf{I}_M - \lambda N (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1} \right) \boldsymbol{\Psi}^\top_{>M}$$

$$= \boldsymbol{\Psi} \boldsymbol{\Psi}^\top_{>M} - \lambda N \boldsymbol{\Psi} (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1} \boldsymbol{\Psi}^\top_{>M}.$$

Analogously to the computations in (22)-(23)

$$(\mathbf{I}_M + \boldsymbol{\Delta}) \mathbf{P}_{>M} = \boldsymbol{E} - \lambda \boldsymbol{\Psi} (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1} \boldsymbol{\Psi}^\top_{>M}$$

$$(\mathbf{I}_M + \boldsymbol{\Delta}) \mathbf{P}_{>M} = \boldsymbol{E} - \lambda \boldsymbol{\Lambda}^{-1} \boldsymbol{\Lambda} \boldsymbol{\Psi} (\mathbf{K} + \lambda N \mathbf{I}_M)^{-1} \boldsymbol{\Psi}^\top_{>M}$$

$$(\mathbf{I}_M + \boldsymbol{\Delta}) \mathbf{P}_{>M} = \boldsymbol{E} - \lambda \boldsymbol{\Lambda}^{-1} \mathbf{P}_{>M}$$

$$(\boldsymbol{\Lambda} + \boldsymbol{\Lambda} \boldsymbol{\Delta}) \mathbf{P}_{>M} = \boldsymbol{\Lambda} \boldsymbol{E} - \lambda \mathbf{P}_{>M}$$

$$(\boldsymbol{\Lambda} + \boldsymbol{\Lambda} \boldsymbol{\Delta} + \lambda \mathbf{I}_M) \mathbf{P}_{>M} = \boldsymbol{\Lambda} \boldsymbol{E}$$

$$\mathbf{P}_{>M} = (\boldsymbol{\Lambda} + \boldsymbol{\Lambda} \boldsymbol{\Delta} + \lambda \mathbf{I}_M)^{-1} \boldsymbol{\Lambda} \boldsymbol{E}$$

$$\mathbf{P}_{>M} = \mathbf{B} \boldsymbol{E}.$$

$\square$

**Lemma C.4** (Fitting Error). *Recall the notation*

$$\text{fitting error} = \|\tilde{\gamma} - \mathbf{P}_{\leq M}\tilde{\gamma} - \tilde{\gamma}_{>M}\mathbf{P}_{>M}\|_2^2,$$

$$\mathbf{B} \overset{\text{def.}}{=} (\mathbf{I}_M + \boldsymbol{\Delta} + \lambda\boldsymbol{\Lambda}^{-1})^{-1}.$$

*We have fitting error* $= \left\|\mathbf{B}\left(\lambda\boldsymbol{\Lambda}^{-1}\tilde{\gamma} - \boldsymbol{E}\tilde{\gamma}_{>M}\right)\right\|_2^2.$

*Proof.* By lemmata C.2 and C.3,

$$\|\tilde{\gamma} - \mathbf{P}_{\leq M}\tilde{\gamma} - \tilde{\gamma}_{>M}\mathbf{P}_{>M}\|_2^2 = \left\|\tilde{\gamma} - (\mathbf{I}_M - \lambda\mathbf{B}\boldsymbol{\Lambda}^{-1})\tilde{\gamma}\tilde{\gamma} - \mathbf{B}\boldsymbol{E}\tilde{\gamma}_{>M}\right\|_2^2$$
$$= \left\|\mathbf{B}\left(\lambda\boldsymbol{\Lambda}^{-1}\tilde{\gamma} - \boldsymbol{E}\tilde{\gamma}_{>M}\right)\right\|_2^2.$$

$\square$

Hence we come up with a new expression of the bias:

**Proposition C.5** (Bias). *Recall that* $\mathbf{B} \overset{\text{def.}}{=} (\mathbf{I}_M + \boldsymbol{\Delta} + \lambda\boldsymbol{\Lambda}^{-1})^{-1}$. *The bias* $\mathbb{E}_x\left(f_{\mathbf{X}}^{\lambda}(x) - \tilde{f}(x)\right)^2$ *has the following expression:*

$$\text{bias} = \tilde{\gamma}_{>M}^2 + \left\|\mathbf{B}\left(\lambda\boldsymbol{\Lambda}^{-1}\tilde{\gamma} - \tilde{\gamma}_{>M}\boldsymbol{E}\right)\right\|_2^2.$$

*Proof.* We apply Proposition C.1 and Lemma C.4 to obtain the result. $\square$

### C.1.2 Variance

If we consider noise in the label, we have to compute the variance part of the test error.

**Proposition C.6** (Variance Expression). *Define*

$$\mathbf{M} \overset{\text{def.}}{=} \mathbb{E}_x[\mathbf{K}_x\mathbf{K}_x^\top]$$
$$= \mathbb{E}_x[\boldsymbol{\Psi}^\top\boldsymbol{\Lambda}\boldsymbol{\psi}(x)\boldsymbol{\psi}(x)^\top\boldsymbol{\Lambda}\boldsymbol{\Psi}]$$
$$= \boldsymbol{\Psi}^\top\boldsymbol{\Lambda}\mathbb{E}_x[\boldsymbol{\psi}(x)\boldsymbol{\psi}(x)^\top]\boldsymbol{\Lambda}\boldsymbol{\Psi}$$
$$= \boldsymbol{\Psi}^\top\boldsymbol{\Lambda}\mathbf{I}_M\boldsymbol{\Lambda}\boldsymbol{\Psi}$$
$$= \boldsymbol{\Psi}^\top\boldsymbol{\Lambda}^2\boldsymbol{\Psi}.$$

*We can further simplify the variance part:*

$$\text{variance} \overset{\text{def.}}{=} \mathbb{E}_{x,\varepsilon}\left[\left(\mathbf{K}_x^\top\mathbf{R}\boldsymbol{\varepsilon}\right)^2\right]$$
$$= \mathbb{E}_{x,\varepsilon}\left[\boldsymbol{\varepsilon}^\top\mathbf{R}\mathbf{K}_x\mathbf{K}_x^\top\mathbf{R}\boldsymbol{\varepsilon}\right]$$
$$= \mathbb{E}_{\varepsilon}\left[\boldsymbol{\varepsilon}^\top\mathbf{R}\mathbf{M}\mathbf{R}\boldsymbol{\varepsilon}\right]$$
$$= \sigma^2 \cdot \text{Tr}[\mathbf{R}\mathbf{M}\mathbf{R}].$$

**Theorem C.7** (Variance). *Recall that* $\mathbf{B} \overset{\text{def.}}{=} (\mathbf{I}_M + \boldsymbol{\Delta} + \lambda\boldsymbol{\Lambda}^{-1})^{-1}$. *The variance part, variance, can be expressed as:*

$$\text{variance} = \frac{\sigma^2}{N}\text{Tr}\left[\mathbf{B}^2(\mathbf{I}_M + \boldsymbol{\Delta})\right].$$

*Proof.* We argue similarly as in lemma C.2. Since

$$\boldsymbol{\Psi}\boldsymbol{\Psi}^\top\boldsymbol{\Lambda}\boldsymbol{\Psi}\mathbf{R} = \boldsymbol{\Psi}\mathbf{K}(\mathbf{K} + \lambda N\mathbf{I}_M)^{-1}$$
$$= \boldsymbol{\Psi}(\mathbf{I}_M - \lambda N\mathbf{R})$$
$$= \boldsymbol{\Psi} - \lambda N\boldsymbol{\Psi}\mathbf{R},$$

therefore, we deduce that

$$(\mathbf{I}_M + \boldsymbol{\Delta})\boldsymbol{\Lambda}\boldsymbol{\Psi}\mathbf{R} = \frac{1}{N}\boldsymbol{\Psi} - \lambda\boldsymbol{\Psi}\mathbf{R} \tag{24}$$

$$(\mathbf{I}_M + \boldsymbol{\Delta})\boldsymbol{\Lambda}\boldsymbol{\Psi}\mathbf{R} = \frac{1}{N}\boldsymbol{\Psi} - \lambda\boldsymbol{\Lambda}^{-1}\boldsymbol{\Lambda}\boldsymbol{\Psi}\mathbf{R}$$

$$(\mathbf{I}_M + \boldsymbol{\Delta} + \lambda\boldsymbol{\Lambda}^{-1})\boldsymbol{\Lambda}\boldsymbol{\Psi}\mathbf{R} = \frac{1}{N}\boldsymbol{\Psi}$$

$$\boldsymbol{\Lambda}\boldsymbol{\Psi}\mathbf{R} = \frac{1}{N}(\mathbf{I}_M + \boldsymbol{\Delta} + \lambda\boldsymbol{\Lambda}^{-1})^{-1}\boldsymbol{\Psi}$$

$$\boldsymbol{\Lambda}\boldsymbol{\Psi}\mathbf{R} = \frac{1}{N}\mathbf{B}\boldsymbol{\Psi}. \tag{25}$$

By leveraging the identity $\mathbf{M} = \boldsymbol{\Psi}^\top\boldsymbol{\Lambda}^2\boldsymbol{\Psi}$ and elementary properties of the trace map, the computations in (24)-(25) imply that

$$\mathrm{Tr}[\mathbf{R}\mathbf{M}\mathbf{R}] = \mathrm{Tr}[\mathbf{R}\boldsymbol{\Psi}^\top\boldsymbol{\Lambda}^2\boldsymbol{\Psi}\mathbf{R}] \tag{26}$$

$$= \mathrm{Tr}\left[(\boldsymbol{\Lambda}\boldsymbol{\Psi}\mathbf{R})^\top(\boldsymbol{\Lambda}\boldsymbol{\Psi}\mathbf{R})\right] \tag{27}$$

$$= \mathrm{Tr}\left[\left(\frac{1}{N}\mathbf{B}\boldsymbol{\Psi}\right)^\top\left(\frac{1}{N}\mathbf{B}\boldsymbol{\Psi}\right)\right] \tag{28}$$

$$= \frac{1}{N}\mathrm{Tr}\left[\frac{1}{N}\boldsymbol{\Psi}^\top\mathbf{B}^\top\mathbf{B}\boldsymbol{\Psi}\right]$$

$$= \frac{1}{N}\mathrm{Tr}\left[\mathbf{B}^\top\mathbf{B}\cdot\frac{1}{N}\boldsymbol{\Psi}\boldsymbol{\Psi}^\top\right] \tag{29}$$

$$= \frac{1}{N}\mathrm{Tr}\left[\mathbf{B}^\top\mathbf{B}(\mathbf{I}_M + \boldsymbol{\Delta})\right] \tag{30}$$

$$= \frac{1}{N}\mathrm{Tr}\left[\mathbf{B}^2(\mathbf{I}_M + \boldsymbol{\Delta})\right]; \tag{31}$$

in more detail: in line (26), we use the definition of $\mathbf{M}$; in line (27), we use the fact that both $\boldsymbol{\Lambda}$ and $\mathbf{R}$ are symmetric; in line (28), we use line (25); in line (29), we use the cyclicity of the trace; in line (30), we use the definition of $\boldsymbol{\Delta}$; in line (31), we use the symmetry of $\mathbf{B}$. We obtain the result upon applying Lemma C.6. $\qquad\square$

### C.1.3 Test Error

The Bias-Variance trade-off (see Proposition B.5) decomposed the KRR's test error into two terms, the bias and variance. Since Propositions C.5 and C.7 give us exact expressions for the bias and variance, respectively, we deduce the following exact expression for the KRR's test error.

**Theorem C.8** (Exact Formula for KRR's Test Error). *The test error $\mathcal{R}_{\mathbf{Z},\lambda}$ of KRR equals*

$$\mathcal{R}_{\mathbf{Z},\lambda} = \underbrace{\overbrace{\left\|\mathbf{B}\left(\lambda\boldsymbol{\Lambda}^{-1}\tilde{\boldsymbol{\gamma}} - \tilde{\gamma}_{>M}\boldsymbol{E}_M\right)\right\|_2^2}^{\textit{fitting error}} + \overbrace{\tilde{\gamma}_{>M}^2}^{\textit{finite rank error}}}_{\textit{bias}} + \underbrace{\frac{\sigma_{\textit{noise}}^2}{N}\mathrm{Tr}\left[\mathbf{B}^2(\mathbf{I}_M + \boldsymbol{\Delta})\right]}_{\textit{variance}},$$

*where* $\mathbf{B} \overset{\text{def.}}{=} (\mathbf{I}_M + \boldsymbol{\Delta} + \lambda\boldsymbol{\Lambda}^{-1})^{-1}$.

*Proof.* We begin with the bias/variance decomposition:

$$R_{\mathbf{Z}}^\lambda \overset{\text{def.}}{=} \mathbb{E}_y\|f_{\mathbf{Z}}^\lambda - \tilde{f}\|_{L_{\rho_\mathcal{X}}^2}^2$$

$$= \mathbb{E}_{x,y}\left(\mathbf{K}_x^\top\mathbf{R}\mathbf{y} - \tilde{f}(x)\right)^2$$

$$= \mathbb{E}_{\varepsilon,x}\left(\mathbf{K}_x^\top\mathbf{R}(\tilde{f}(\mathbf{X}) + \varepsilon) - \tilde{f}(x)\right)^2$$

$$= \mathbb{E}_x\left(f_{\mathbf{X}}^\lambda(x) - \tilde{f}(x)\right)^2 + \mathbb{E}_{x,\varepsilon}\left[\left(\mathbf{K}_x^\top\mathbf{R}\varepsilon\right)^2\right]$$

$$= \text{bias} + \text{variance},$$

457  then we apply Propositions C.5 and C.7. □

458  For the validation of the Theorem C.8, please see Appendix D for details.

459  The matrix $\mathbf{B}$ plays an important role in the expression since it encodes most information of the KRR.
460  Therefore, the following subsection will discuss the approximation of the matrix $\mathbf{B}$.

## C.2  Matrix Approximation

462  Recall that the matrix $\mathbf{B} \overset{\text{def.}}{=} (\mathbf{I}_M + \boldsymbol{\Delta} + \lambda\boldsymbol{\Lambda}^{-1})^{-1}$ is the inverse of a random matrix. The following
463  lemma helps to approximate $\mathbf{B}$. Informally, it says that: given that $\delta \overset{\text{def.}}{=} \|\boldsymbol{\Delta}\|_{\text{op}} < 1$. We have

$$\mathbf{B} = \sum_{s=0}^{\infty} (-\bar{\mathbf{P}}\boldsymbol{\Delta})^s \bar{\mathbf{P}}$$

464  in operator norm $\|\cdot\|_{op}$ for an $M \times M$ matrix $\bar{P}$ depending only on the $M$ eigenvalues $\{\lambda_k\}_{k=1}^{M}$
465  and on the ridge $\lambda > 0$. More precisely we have the following.

466  **Lemma C.9** (B-Expansion). *Given that $\delta \overset{\text{def.}}{=} \|\boldsymbol{\Delta}\|_{op} < 1$. It holds that*

$$\lim_{n \uparrow \infty} \left\| \mathbf{B} - \sum_{s=0}^{n} (-\bar{\mathbf{P}}\boldsymbol{\Delta})^s \bar{\mathbf{P}} \right\|_{op} = 0$$

467  *where $\bar{\mathbf{P}} \overset{\text{def.}}{=} \operatorname{diag}\left[\frac{\lambda_k}{\lambda_k + \lambda}\right]_k = \boldsymbol{\Lambda}(\boldsymbol{\Lambda} + \lambda\mathbf{I}_M)^{-1} \in \mathbb{R}^{M \times M}$.*

468  *Proof.* Set $\mathbf{A} = \mathbf{I}_M + \lambda\boldsymbol{\Lambda}^{-1}$ and repeatedly use the formula $(\mathbf{A} + \boldsymbol{\Delta})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\boldsymbol{\Delta}(\mathbf{A} + \boldsymbol{\Delta})^{-1}$
469  from [31], we have

$$\begin{aligned}
\mathbf{B} &\overset{\text{def.}}{=} (\mathbf{I}_M + \boldsymbol{\Delta} + \lambda\boldsymbol{\Lambda}^{-1})^{-1} \\
&= (\mathbf{A} + \boldsymbol{\Delta})^{-1} \\
&= \mathbf{A}^{-1} - \mathbf{A}^{-1}\boldsymbol{\Delta}(\mathbf{A} + \boldsymbol{\Delta})^{-1} \\
&= \mathbf{A}^{-1} - \mathbf{A}^{-1}\boldsymbol{\Delta}\left(\mathbf{A}^{-1} - \mathbf{A}^{-1}\boldsymbol{\Delta}(\mathbf{A} + \boldsymbol{\Delta})^{-1}\right) \\
&= \mathbf{A}^{-1} - \mathbf{A}^{-1}\boldsymbol{\Delta}\mathbf{A}^{-1} + (\mathbf{A}^{-1}\boldsymbol{\Delta})^2(\mathbf{A} + \boldsymbol{\Delta})^{-1} \\
&= \sum_{s=0}^{n} (-\mathbf{A}^{-1}\boldsymbol{\Delta})^s \mathbf{A}^{-1} + (-\mathbf{A}^{-1}\boldsymbol{\Delta})^{n+1}(\mathbf{A} + \boldsymbol{\Delta})^{-1}
\end{aligned}$$

470  Note that $\mathbf{A}^{-1} = (\mathbf{I}_M + \lambda\boldsymbol{\Lambda}^{-1})^{-1} = \boldsymbol{\Lambda}(\boldsymbol{\Lambda} + \lambda\mathbf{I}_M)^{-1} = \bar{\mathbf{P}}$ with operator norm $\frac{\lambda_1}{\lambda_1 + \lambda} < 1$, hence
471  we have $(\mathbf{A}^{-1}\boldsymbol{\Delta})^{n+1} = (-\bar{\mathbf{P}}\boldsymbol{\Delta})^{n+1} \to 0$ in operator norm as $n \to \infty$. Hence

$$\begin{aligned}
\mathbf{B} &= \sum_{s=0}^{\infty} (-\mathbf{A}^{-1}\boldsymbol{\Delta})^s \mathbf{A}^{-1} \\
&= \sum_{s=0}^{\infty} (-\bar{\mathbf{P}}\boldsymbol{\Delta})^s \bar{\mathbf{P}}
\end{aligned}$$

472  in operator norm. □

473  Due to the convergence result in lemma C.9, it is natural to define:

474  **Definition C.10.** *For any $n \in \mathbb{N} \cup \{\infty\}$, write $\mathbf{B}^{(n)} = \sum_{s=0}^{n} (-\bar{\mathbf{P}}\boldsymbol{\Delta})^s \bar{\mathbf{P}}$. For example, We have*

$$\mathbf{B}^{(0)} = \bar{\mathbf{P}}$$
$$\mathbf{B}^{(1)} = \bar{\mathbf{P}} - \bar{\mathbf{P}}\boldsymbol{\Delta}\bar{\mathbf{P}}$$
$$\mathbf{B}^{(\infty)} = \mathbf{B}$$

Although lemma C.9 is valid when $\delta < 1$, we need a slightly stronger condition that $\delta$ is upper bounded by an arbitrary constant strictly small than 1. For simplicity, we assume this constant to be $\frac{1}{2}$ in the following lemma:

**Lemma C.11** (B-Approximation). *Assume that $\delta \overset{\text{def.}}{=} \|\boldsymbol{\Delta}\|_{op} < \frac{1}{2}$. Let $\mathbf{B}^{(n)} = \sum_{s=0}^{n}(-\bar{\mathbf{P}}\boldsymbol{\Delta})^s\bar{\mathbf{P}}$ be the $n$th-order approximation of the matrix $\mathbf{B}$ as in definition C.10. Then we have*

$$\left\|\mathbf{B} - \mathbf{B}^{(n)}\right\|_{op} < 2\delta^{n+1}.$$

*Proof.* We first bound the operator norm of the matrix $\mathbf{B}$: since the minimum singular value of the matrix $\bar{\mathbf{P}}^{-1} + \boldsymbol{\Delta}$ is at least

$$\frac{\lambda_k + \lambda}{\lambda_k} - \|\Delta\|_{op} \geq 1 + \frac{\lambda}{\lambda_k} - \frac{1}{2} > \frac{1}{2},$$

and hence

$$\|\mathbf{B}\|_{op} = \left\|\left(\bar{\mathbf{P}}^{-1} + \boldsymbol{\Delta}\right)^{-1}\right\|_{op} < 2.$$

Also, we have

$$\begin{aligned}
\mathbf{B} - \mathbf{B}^{(n)} &= \sum_{s=n+1}^{\infty}(-\bar{\mathbf{P}}\boldsymbol{\Delta})^s\bar{\mathbf{P}} \\
&= (-\bar{\mathbf{P}}\boldsymbol{\Delta})^{n+1}\sum_{s=0}^{\infty}(-\bar{\mathbf{P}}\boldsymbol{\Delta})^s\bar{\mathbf{P}} \\
&= (-\bar{\mathbf{P}}\boldsymbol{\Delta})^{n+1}\mathbf{B}.
\end{aligned}$$

Hence $\left\|\mathbf{B} - \mathbf{B}^{(n)}\right\|_{op} \leq \left\|\bar{\mathbf{P}}\boldsymbol{\Delta}\right\|_{op}^{n+1}\|\mathbf{B}\|_{op} < \|\boldsymbol{\Delta}\|_{op}^{n+1} \cdot 2 = 2\delta^{n+1}$, since we have $\left\|\bar{\mathbf{P}}\right\|_{op} = \frac{\lambda_1}{\lambda_1+\lambda} < 1$. $\square$

Note that the upper bound $\frac{1}{2}$ of $\delta$ can be replaced by any constant strictly small than 1 to get a similar conclusion.

**Remark C.12.** *Using the concentration result from random matrix theory, for $M < N$, one can show with high probability that the operator norm $\delta$ of the fluctuation matrix $\boldsymbol{\Delta}$ is less than 1.* [4]

See subsection C.3 for details. Then we can use the the above lemmata C.9 and C.11 to approximate the test error of KRR:

**Proposition C.13** (Bias Approximation). *Fix a sample $\mathbf{Z}$ of $\rho$ such that $\delta \overset{\text{def.}}{=} \|\boldsymbol{\Delta}\|_{op} < \frac{1}{2}$. Then the $bias_{test}$ term is bounded above and below by*

$$\left|bias - \left(\left\|\bar{\mathbf{P}}w\right\|_2^2 + \tilde{\gamma}_{>M}^2\right)\right| \leq 2\delta\left\|\bar{\mathbf{P}}w\right\|_2^2 + \|w\|_2^2\,\delta^2 p(\delta),$$

*where $\bar{\mathbf{P}} \overset{\text{def.}}{=} \boldsymbol{\Lambda}(\boldsymbol{\Lambda} + \lambda\mathbf{I}_M)^{-1}$, $w = \lambda\boldsymbol{\Lambda}^{-1}\tilde{\boldsymbol{\gamma}} - \tilde{\gamma}_{>M}E$, and $p(\delta) \overset{\text{def.}}{=} 5 + 4\delta + 4\delta^2$. By writing $E = (\eta_k)_{k=1}^{M}$, the upper-bound simplifies to*

$$bias \leq (1 + 2\delta)\sum_{k=1}^{M}\frac{(\lambda\tilde{\gamma}_k - \tilde{\gamma}_{>M}\eta_k\lambda_k)^2}{(\lambda_k + \lambda)^2} + \tilde{\gamma}_{>M}^2 + \|w\|_2^2\,\delta^2 p(\delta).$$

*Analogously, the lower bound can be derived.*

---

[4]From there, we differentiate the approach from Bach [6]: From Propositions C.5 and C.7, it is inevitable to approximate the matrix $\mathbf{B}$, and we have $\mathbf{I}_M$ as support of the inverse. Bach instead uses RHKS basis to express the fluctuation matrix and is hence forced to use $\lambda\mathbf{I}_M$ as the support. As a result, he would need to require that the fluctuation is less than $\lambda$ and hence his requirement on $N$ is antiproportional to $\lambda$ in Theorem 5.1.

495  *Proof.* Let $w = \lambda \mathbf{\Lambda}^{-1} \tilde{\boldsymbol{\gamma}} - \tilde{\gamma}_{>M} \boldsymbol{E}$. We apply lemma C.4 followed by the 1st-order approximation
496  $\mathbf{B}^{(1)}$ of the matrix $\mathbf{B}$ in lemma C.11:

$$\text{fitting error} = \|\mathbf{B}w\|_2^2 = \left\| \mathbf{B}^{(1)}w + \left( \mathbf{B} - \mathbf{B}^{(1)} \right) w \right\|_2^2$$

$$= \left\| \mathbf{B}^{(1)}w \right\|_2^2 + w^\top \left( \mathbf{B} - \mathbf{B}^{(1)} \right) \mathbf{B}^{(1)}w + w^\top \mathbf{B}^{(1)} \left( \mathbf{B} - \mathbf{B}^{(1)} \right) w + \left\| \left( \mathbf{B} - \mathbf{B}^{(1)} \right) w \right\|_2^2$$

$$\leq \left\| \mathbf{B}^{(1)}w \right\|_2^2 + 2 \left\| \mathbf{B}^{(1)} \right\|_{\text{op}} \left\| \mathbf{B} - \mathbf{B}^{(1)} \right\|_{\text{op}} \|w\|_2^2 + \left\| \mathbf{B} - \mathbf{B}^{(1)} \right\|_{\text{op}}^2 \|w\|_2^2$$

$$\leq \left\| \mathbf{B}^{(1)}w \right\|_2^2 + 2 \cdot (1 + \delta) \cdot 2\delta^2 \|w\|_2^2 + 4\delta^4 \|w\|_2^2$$

$$\leq \left\| \mathbf{B}^{(1)}w \right\|_2^2 + 4 \|w\|_2^2 \delta^2 (1 + \delta + \delta^2)$$

$$\leq \left\| \left( \bar{\mathbf{P}} - \bar{\mathbf{P}} \boldsymbol{\Delta} \bar{\mathbf{P}} \right) w \right\|_2^2 + 4 \|w\|_2^2 \delta^2 (1 + \delta + \delta^2)$$

$$\leq \left\| \mathbf{I}_M - \bar{\mathbf{P}} \boldsymbol{\Delta} \right\|_{\text{op}}^2 \left\| \bar{\mathbf{P}} w \right\|_2^2 + 4 \|w\|_2^2 \delta^2 (1 + \delta + \delta^2)$$

$$\leq \left( 1 + 2 \left\| \bar{\mathbf{P}} \boldsymbol{\Delta} \right\|_{\text{op}} + \left\| \bar{\mathbf{P}} \boldsymbol{\Delta} \right\|_{\text{op}}^2 \right) \left\| \bar{\mathbf{P}} w \right\|_2^2 + 4 \|w\|_2^2 \delta^2 (1 + \delta + \delta^2)$$

$$\leq \left( 1 + 2 \left\| \bar{\mathbf{P}} \boldsymbol{\Delta} \right\|_{\text{op}} \right) \left\| \bar{\mathbf{P}} w \right\|_2^2 + \|w\|_2^2 \delta^2 (5 + 4\delta + 4\delta^2)$$

$$\leq (1 + 2\delta) \left\| \bar{\mathbf{P}} w \right\|_2^2 + \|w\|_2^2 \delta^2 (5 + 4\delta + 4\delta^2).$$

497  Hence we have the upper bound:

$$\text{bias} \leq \tilde{\gamma}_{>M}^2 + (1 + 2\delta) \left\| \bar{\mathbf{P}} w \right\|_2^2 + \|w\|_2^2 \delta^2 p(\delta).$$

498  We argue similarly for the lower bound using: $\|\mathbf{A}\|_{\text{op}} \|v\|_2^2 \geq v^\top \mathbf{A} v \geq - \|\mathbf{A}\|_{\text{op}} \|v\|_2^2$ for any
499  $\mathbf{A} \in \mathbb{R}^{M \times M}$, $v \in \mathbb{R}^{M \times 1}$. $\qquad\square$

500  We argue similarly for variance.

501  **Proposition C.14** (Variance Approximation)**.** *Fix a sampling $\mathbf{Z}$ such that $\delta \overset{\text{def.}}{=} \|\boldsymbol{\Delta}\|_{op} < \frac{1}{2}$. Then we*
502  *have*

$$\left| \text{variance} - \frac{\sigma^2}{N} \sum_{k=1}^{M} \frac{\lambda_k^2}{(\lambda_k + \lambda)^2} \right| \leq \delta \frac{\sigma^2}{N} \sum_{k=1}^{M} \frac{\lambda_k^2}{(\lambda_k + \lambda)^2} + M \frac{\sigma^2}{N} (1 + \delta) \delta^2 p(\delta),$$

503  *where $p(\delta) \overset{\text{def.}}{=} 5 + 4\delta + 4\delta^2$, and $\sigma^2 \overset{\text{def.}}{=} \mathbb{E}[\epsilon^2]$ is the noise variance.*

504  *Proof.* Note that $\text{Tr}\, \mathbf{A} \leq M \|\mathbf{A}\|_{\text{op}}$ for any matrix $\mathbf{A} \in \mathbb{R}^{M \times M}$. Since $\mathbf{B}^2(\mathbf{I}_M + \boldsymbol{\Delta}) =$
505  $(\mathbf{B}^{(1)})^2(\mathbf{I}_M + \boldsymbol{\Delta}) + 2\mathbf{B}^{(1)} \left( \mathbf{B} - \mathbf{B}^{(1)} \right) (\mathbf{I}_M + \boldsymbol{\Delta}) + \left( \mathbf{B} - \mathbf{B}^{(1)} \right)^2 (\mathbf{I}_M + \boldsymbol{\Delta})$, we can bound the
506  residue term by $\delta$:

$$\text{Tr} \left[ 2\mathbf{B}^{(1)} \left( \mathbf{B} - \mathbf{B}^{(1)} \right) (\mathbf{I}_M + \boldsymbol{\Delta}) + \left( \mathbf{B} - \mathbf{B}^{(1)} \right)^2 (\mathbf{I}_M + \boldsymbol{\Delta}) \right]$$

$$\leq M(1 + \delta) \left\| \mathbf{B} - \mathbf{B}^{(1)} \right\|_{\text{op}} \left( 2 \left\| \mathbf{B}^{(1)} \right\|_{\text{op}} + \left\| \mathbf{B} - \mathbf{B}^{(1)} \right\|_{\text{op}} \right)$$

$$\leq M(1 + \delta) \cdot 2\delta^2 (2(1 + \delta) + 2\delta^2)$$

$$\leq 4M\delta^2 (1 + \delta)(1 + \delta + \delta^2),$$

507  For the main terms, we have

$$\text{Tr}[(\mathbf{B}^{(1)})^2(\mathbf{I}_M + \boldsymbol{\Delta})] \leq \text{Tr}[\bar{\mathbf{P}}^2] \cdot \left\| (\mathbf{I}_M - \boldsymbol{\Delta}\bar{\mathbf{P}})^2 (\mathbf{I}_M + \boldsymbol{\Delta}) \right\|_{\text{op}}$$

$$= \text{Tr}[\bar{\mathbf{P}}^2] \left\| \mathbf{I}_M + \boldsymbol{\Delta}(\mathbf{I}_M - 2\bar{\mathbf{P}}) + (\boldsymbol{\Delta}\bar{\mathbf{P}})^2 - 2\boldsymbol{\Delta}\bar{\mathbf{P}}\boldsymbol{\Delta} + (\boldsymbol{\Delta}\bar{\mathbf{P}})^2 \boldsymbol{\Delta} \right\|_{\text{op}}$$

$$\leq \text{Tr}[\bar{\mathbf{P}}^2] \left\| \mathbf{I}_M + \boldsymbol{\Delta}(\mathbf{I}_M - 2\bar{\mathbf{P}}) \right\|_{\text{op}} + M \left\| (\boldsymbol{\Delta}\bar{\mathbf{P}})^2 - 2\boldsymbol{\Delta}\bar{\mathbf{P}}\boldsymbol{\Delta} + (\boldsymbol{\Delta}\bar{\mathbf{P}})^2 \boldsymbol{\Delta} \right\|_{\text{op}}$$

$$\leq \text{Tr}[\bar{\mathbf{P}}^2](1 + \delta) + M\delta^2 (1 + \delta).$$

We apply Theorem C.7 to yield a bound on variance:

$$\left| \text{variance} - \frac{\sigma^2}{N} \sum_{k=1}^{M} \frac{\lambda_k^2}{(\lambda_k + \lambda)^2} \right| \leq \delta \frac{\sigma^2}{N} \sum_{k=1}^{M} \frac{\lambda_k^2}{(\lambda_k + \lambda)^2} + M \frac{\sigma^2}{N} (1 + \delta) \delta^2 p(\delta).$$

$\square$

Note that the above propositions C.13 and C.14 give absolute (non-probabilistic) bounds on the test error, once $\delta$ is controlled.

## C.3 Concentration Results

In this subsection, we focus on bounding the operator norm $\delta$ of the fluctuation matrix $\mathbf{\Delta}$. We can assume the data-generating distribution $\rho$ and eigenfunctions $\psi_k$ are well-behaved in the sense that:

**Assumption C.15** (Sub-Gaussian-ness). *We assume probability distribution of the random variable $\psi_k(x)$, where $x \in \rho$, has sub-Gaussian norm bounded by a positive constant $G > 0$, for all $k \in \{1, ..., M\} \cup \{> M\}$[5].*

In particular, if the random variable $\psi_k(x)$ is bounded, the assumption C.15 is fulfilled. First, we establish some concentration results.

**Lemma C.16** (Theorem 3.59 in [36]). *Let $\mathbf{A}$ be an $n \times N$ matrix with independent isotropic sub-Gaussian columns in $\mathbb{R}^n$ which sub-gaussian norm is bounded by a positive constant $G > 0$. Then for all $t \geq 0$, with probability at least $1 - 2\exp(-\frac{1}{3}t^2)$, we have*

$$\left\| \frac{1}{N} \mathbf{A}\mathbf{A}^\top - \mathbf{I}_n \right\|_{op} \leq \max(a, a^2), \tag{32}$$

*where $a \overset{\text{def.}}{=} C\sqrt{\frac{n}{N}} + \frac{t}{\sqrt{N}}$, for all constant $C \geq 12G^2$ .*

*Proof.* Let $a \overset{\text{def.}}{=} C\sqrt{\frac{n}{N}} + \frac{t}{\sqrt{N}}$ with $C > 0$ to be determined, and $\epsilon \overset{\text{def.}}{=} \max\{a, a^2\}$. The first step to show that :

$$\max_{x \in \mathcal{N}} \left| \frac{1}{N} \left\| \mathbf{A}^\top x \right\|_2^2 - 1 \right| \leq \epsilon$$

for some $\frac{1}{4}$-net $\mathcal{N}$ on the sphere $\mathbb{S}^{n-1} \subset \mathbb{R}^n$. Choose such a net $\mathcal{N}$ with $|\mathcal{N}| < \left( 1 + \frac{2}{1/4} \right)^n = 9^n$. Let $\mathbf{A}_i$ be the $i$th column of the matrix $\mathbf{A}$ and let $Z_i \overset{\text{def.}}{=} \mathbf{A}_i^\top x$ be a random variable. By definition of $\mathbf{A}$, $Z_i$ is centered with unit variance with sub-Gaussian norm upper bounded by $G$. Note that $G \geq \frac{1}{\sqrt{2}} \mathbb{E}[Z_i^2]^{1/2} = \frac{1}{\sqrt{2}}$, and the random variable $Z_i^2 - 1$ is centered and has sub-exponential norm upper bounded by $4G^2$. Hence by an exponential deviation inquality [6], we have, for any $x \in \mathbb{S}^{n-1}$:

$$\mathbb{P}\left\{ \left| \frac{1}{N} \left\| \mathbf{A}^\top x \right\|_2^2 - 1 \right| \geq \frac{\epsilon}{2} \right\} = \mathbb{P}\left\{ \left| \frac{1}{N} \sum_{i=1}^{N} Z_i^2 - 1 \right| \geq \frac{\epsilon}{2} \right\}$$

$$\leq 2\exp\left( -\frac{1}{2} e^{-1} G^{-4} \min\{\epsilon, \epsilon^2\} \right)$$

$$= \leq 2\exp\left( -\frac{1}{2} e^{-1} G^{-4} a^2 \right)$$

$$\leq 2\exp\left( -\frac{1}{2} e^{-1} G^{-4} (C^2 n + t^2) \right).$$

---

[5]it means the orthonormal complement $\psi_{>M}$ is also mentioned in the assumption.

[6]This inequality is Corollary 5.17 from [36].

Then by union bound, we have

$$\mathbb{P}\left\{\max_{x\in\mathcal{N}}\left|\frac{1}{N}\left\|\mathbf{A}^\top x\right\|_2^2-1\right|\geq\frac{\epsilon}{2}\right\}\leq 9^n\cdot 2\exp\left(-\frac{1}{2}e^{-1}G^{-4}(C^2n+t^2)\right)$$

$$\leq 2\exp\left(-\frac{1}{2}e^{-1}G^{-4}t^2\right),$$

for $C\geq\sqrt{2e\log 9}G^2$. Since $12>\sqrt{2e\log 9}$, for simplicity, we assume $C>12G^2$. Moreover, since $G\geq\frac{1}{\sqrt{2}}$, we have $\frac{1}{2}e^{-1}G^{-4}\leq\frac{1}{3}$, we have

$$\mathbb{P}\left\{\max_{x\in\mathcal{N}}\left|\frac{1}{N}\left\|\mathbf{A}^\top x\right\|_2^2-1\right|\geq\frac{\epsilon}{2}\right\}\leq 2\exp\left(-\frac{1}{3}t^2\right).$$

Then by the $\frac{1}{4}$-net argument, with probability at least $1-2\exp\left(-\frac{1}{3}t^2\right)$, we have

$$\left\|\frac{1}{N}\mathbf{A}\mathbf{A}^\top-\mathbf{I}_n\right\|_{\mathrm{op}}\leq\frac{4}{2}\max_{x\in\mathcal{N}}\left|\frac{1}{N}\left\|\mathbf{A}^\top x\right\|_2^2-1\right|$$

$$\leq\epsilon=\max\{a,a^2\}.$$

$\square$

**Lemma C.17.** *Assume Assumption C.15 holds and that $N>\exp(4(12G^2)^2(M+1))$. Then with a probability of at least $1-2/N$, we have*

$$\max\left\{\delta,\|\boldsymbol{E}_M\|_2\right\}\leq\sqrt{\frac{\log N}{N}}.$$

*Proof.* Set $n=M+1$, $\mathbf{A}=\left(\begin{smallmatrix}\boldsymbol{\Psi}_{\leq M}\\\psi_{>M}(\mathbf{X})^\top\end{smallmatrix}\right)\in\mathbb{R}^{(M+1)\times N}$. Then

$$\frac{1}{N}\mathbf{A}\mathbf{A}^\top-\mathbf{I}_n=\begin{pmatrix}\frac{1}{N}\boldsymbol{\Psi}_{\leq M}\boldsymbol{\Psi}_{\leq M}^\top & \boldsymbol{E}_M\\ \boldsymbol{E}_M^\top & \eta_{>M}+1\end{pmatrix}-\mathbf{I}_n=\begin{pmatrix}\boldsymbol{\Delta}_M & \boldsymbol{E}_M\\ \boldsymbol{E}_M^\top & \eta_{>M}\end{pmatrix}.$$

where $\eta_{>M}\overset{\text{def.}}{=}\frac{1}{N}\sum_{i=1}^N\psi_{>M}(x_i)^2-1$. On one hand, the operator norm of the above matrix bounds $\delta$ and $\|\boldsymbol{E}_M\|_2$ from above:

$$\left\|\begin{pmatrix}\boldsymbol{\Delta}_M & \boldsymbol{E}_M\\ \boldsymbol{E}_M^\top & \eta_{>M}\end{pmatrix}\right\|_{\mathrm{op}}=\max_{\|\mathbf{u}\|_2^2+v^2=1}\left\|\begin{pmatrix}\boldsymbol{\Delta}_M & \boldsymbol{E}_M\\ \boldsymbol{E}_M^\top & \eta_{>M}\end{pmatrix}\begin{pmatrix}\mathbf{u}\\ v\end{pmatrix}\right\|_2$$

$$=\max_{\|\mathbf{u}\|_2^2+v^2=1}\left\|\begin{pmatrix}\boldsymbol{\Delta}_M\mathbf{u}+v\boldsymbol{E}_M\\ \boldsymbol{E}_M^\top\mathbf{u}+\eta_{>M}v\end{pmatrix}\right\|_2$$

$$\geq\max_{\|\mathbf{u}\|_2^2+v^2=1}\|\boldsymbol{\Delta}_M\mathbf{u}+v\boldsymbol{E}_M\|_2$$

$$\geq\max_{\|\mathbf{u}\|_2^2=1,v=0}\|\boldsymbol{\Delta}_M\mathbf{u}+v\boldsymbol{E}_M\|_2$$

$$\geq\max_{\|\mathbf{u}\|_2^2=1}\|\boldsymbol{\Delta}_M\mathbf{u}\|_2=\delta,$$

and

$$\left\|\begin{pmatrix}\boldsymbol{\Delta}_M & \boldsymbol{E}_M\\ \boldsymbol{E}_M^\top & \eta_{>M}\end{pmatrix}\right\|_{\mathrm{op}}\geq\max_{\|\mathbf{u}\|_2^2+v^2=1}\|\boldsymbol{\Delta}_M\mathbf{u}+v\boldsymbol{E}_M\|_2\geq\max_{\|\mathbf{u}\|_2^2=0,|v|=1}\|\boldsymbol{\Delta}_M\mathbf{u}+v\boldsymbol{E}_M\|_2=\|\boldsymbol{E}_M\|_2.$$

On the other hand, set $t=\frac{1}{2}\sqrt{\log N}$, $C=12G^2$, since $N>\exp(4C^2(M+1))$, we have

$$a=C\sqrt{\frac{n}{N}}+\frac{t}{\sqrt{N}}=12G^2\sqrt{\frac{M+1}{N}}+\frac{1}{2}\sqrt{\frac{\log N}{N}}\leq\sqrt{\frac{\log N}{N}}<1.$$

By Lemma C.17, then with probability of at least $1-2\exp(-\frac{1}{3}t^2)=1-2\exp(-\frac{1}{12})/N>1-2/N$, we have

$$\left\|\begin{pmatrix}\boldsymbol{\Delta}_M & \boldsymbol{E}_M\\ \boldsymbol{E}_M^\top & \eta_{>M}\end{pmatrix}\right\|_{\mathrm{op}}\leq\max\{a,a^2\}=a\leq\sqrt{\frac{\log N}{N}}.$$

Combine the both results and we conclude the upper bounds. $\square$

23

In particular, as $N \to \infty$, $\delta$ vanishes almost surely. In empirical calculation, if the requirement $N > \exp(4(12G^2)^2(M+1))$ exponential in $M$ is too demanding for a large integer $M$, we can take $t = N^s$ for any positive number $s \in \left(0, \frac{1}{2}\right)$ instead of $t = \frac{1}{2}\log N$. In this way, we decrease the requirement to $N$ polynomial in $M$ in sacrificing the decay from $\mathcal{O}\left(\sqrt{\frac{\log N}{N}}\right)$ to $\mathcal{O}\left(N^{s-1/2}\right)$. For simplicity purpose, we do not list out the result with this decay in this paper.

## C.4 Refined Test Error Analysis

We can apply the above concentration results to refine the following bounds on the finite-rank KRR test error.

**Definition C.18.** *To ease the notation, we denote $\underline{r} \overset{\text{def.}}{=} \min_k\{|\tilde{\gamma}_k/\lambda_k|\}$ and $\overline{r} \overset{\text{def.}}{=} \max_k\{|\tilde{\gamma}_k/\lambda_k|\}$.*

### C.4.1 Refined Bounds on Bias

Recall that Proposition C.13 bounding the bias in terms of $\delta$ and $\eta_k$. For the former one, we can choose: for $N > \max\{\exp(4(12G^2)^2(M+1)), 9\}$, by Lemma C.17, with probability of at least $1 - 2/N$, we have $\delta \le \sqrt{\frac{\log N}{N}} < \sqrt{\frac{\log 9}{9}} < \frac{1}{2}$. For the latter one, we have to control the vector $w$:

**Lemma C.19.** *Let $w = \lambda \mathbf{\Lambda}^{-1}\tilde{\gamma} - \tilde{\gamma}_{>M}\mathbf{E}$. We have*

$$\lambda^2 \underline{r}^2 M - 2\lambda \underline{r}|\tilde{\gamma}_{>M}|\sqrt{M}\|E\|_2 \le \|w\|_2^2 \le \left(\lambda\overline{r}\sqrt{M} + \tilde{\gamma}_{>M}\|E\|_2\right)^2;$$

$$\frac{\lambda^2 \lambda_M}{(\lambda_M + \lambda)^2}\|\tilde{f}_{\le M}\|_{\mathcal{H}}^2 - \frac{1}{2}|\tilde{\gamma}_{>M}|\|\tilde{f}_{\le M}\|_{L_\rho^2}\|\mathbf{E}\|_2 \le \|\bar{\mathbf{P}}w\|_2^2 \le \lambda\|\tilde{f}_{\le M}\|_{\mathcal{H}}^2 + \frac{1}{2}|\tilde{\gamma}_{>M}|\|\tilde{f}_{\le M}\|_{L_\rho^2}\|\mathbf{E}\|_2 + \tilde{\gamma}_{>M}^2\|\mathbf{E}\|_2^2.$$

*Proof.* Since $\lambda^2 \underline{r}^2 M \le \left\|\lambda\mathbf{\Lambda}^{-1}\tilde{\gamma}\right\|_2^2 \le \lambda^2 \overline{r}^2 M$ and $\|\tilde{\gamma}_{>M}\mathbf{E}\|_2^2 = \tilde{\gamma}_{>M}^2\|\mathbf{E}\|_2^2$, we have

$$\lambda^2 \underline{r}^2 M - 2\lambda \underline{r}|\tilde{\gamma}_{>M}|\sqrt{M}\|E\|_2 \le \|w\|_2^2 \le \left(\lambda\overline{r}\sqrt{M} + \tilde{\gamma}_{>M}\|E\|_2\right)^2.$$

Similarly, we can bound $\|\bar{\mathbf{P}}w\|_2$. Observe that:

$$\|\bar{\mathbf{P}}w\|_2^2 = \lambda^2 \underbrace{\sum_{k=1}^M \frac{\tilde{\gamma}_k^2}{(\lambda_k + \lambda)^2}}_{I} - 2\lambda\tilde{\gamma}_{>M}\underbrace{\sum_{k=1}^M \frac{\tilde{\gamma}_k \lambda_k \eta_k}{(\lambda_k + \lambda)^2}}_{II} + \tilde{\gamma}_{>M}^2 \underbrace{\sum_{k=1}^M \frac{\lambda_k^2 \eta_k^2}{(\lambda_k + \lambda)^2}}_{III}.$$

Since $1 \ge \frac{\lambda}{\lambda_M + \lambda} \ge \frac{\lambda}{\lambda_k + \lambda}$, we have the upper bound:

$$I = \lambda^2 \sum_{k=1}^M \frac{\tilde{\gamma}_k^2}{(\lambda_k + \lambda)^2} \le \lambda \sum_{k=1}^M \frac{\lambda}{\lambda_k + \lambda}\frac{\tilde{\gamma}_k^2}{\lambda_k + \lambda} \le \lambda \sum_{k=1}^M \frac{\tilde{\gamma}^2}{\lambda_k} = \lambda\|\tilde{f}_{\le M}\|_{\mathcal{H}}^2. \tag{33}$$

where $\tilde{f}_{\le M} \overset{\text{def.}}{=} \sum_{k=1}^M \tilde{\gamma}_k \psi_k = \tilde{f} - \tilde{\gamma}_{>M}\psi_{>M}$. For the lower bound, we have:

$$I = \lambda^2 \sum_{k=1}^M \frac{\tilde{\gamma}_k^2}{(\lambda_k + \lambda)^2} \ge \lambda^2 \sum_{k=1}^M \frac{\lambda_k}{(\lambda_k + \lambda)^2}\frac{\tilde{\gamma}_k^2}{\lambda_k} \ge \lambda^2 \frac{\lambda_M}{(\lambda_M + \lambda)^2}\|\tilde{f}_{\le M}\|_{\mathcal{H}}^2 \tag{34}$$

Similarly, since $4\lambda\lambda_k \le (\lambda_k + \lambda)^2$,

$$|II| = 2\lambda|\tilde{\gamma}_{>M}|\sum_{k=1}^M \frac{|\tilde{\gamma}_k||\lambda_k||\eta_k|}{(\lambda_k + \lambda)^2} \le \frac{1}{2}|\tilde{\gamma}_{>M}|\sum_{k+1}^M |\tilde{\gamma}_k \eta_k| \le \frac{1}{2}|\tilde{\gamma}_{>M}|\sqrt{\sum_{k+1}^M \tilde{\gamma}_k^2 \sum_{k=1}^M \eta_k^2} \le \frac{1}{2}|\tilde{\gamma}_{>M}|\|\tilde{f}_{\le M}\|_{L_\rho^2}\|\mathbf{E}\|_2.$$

And

$$III = \tilde{\gamma}_{>M}^2 \sum_{k=1}^M \frac{\lambda_k^2 \eta_k^2}{(\lambda_k + \lambda)^2} \le \tilde{\gamma}_{>M}^2 \sum_{k=1}^M \eta_k^2 = \tilde{\gamma}_{>M}^2 \|\mathbf{E}\|_2^2.$$

$\square$

Combining the above result, we state the following theorem:

**Theorem C.20.** *For $N > \max\left\{\exp(4(12G^2)^2(M+1)), 9\right\}$ and for any constant $C_1 > 8\left(\lambda\bar{r}\sqrt{M} + \frac{1}{2}|\tilde{\gamma}_{>M}|\right)^2 + \frac{5}{2}\|\tilde{f}\|_{L^2_\rho}^2$ (independent to N), with a probability of at least $1 - 2/N$, we have the upper and lower bounds of bias:*

$$bias \leq \tilde{\gamma}_{>M}^2 + \lambda\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2 + \left(\frac{1}{4}\|\tilde{f}\|_{L^2_\rho}^2 + 2\lambda\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2\right)\sqrt{\frac{\log N}{N}} + C_1\frac{\log N}{N};$$

$$bias \geq \tilde{\gamma}_{>M}^2 + \frac{\lambda^2\lambda_M}{(\lambda_M+\lambda)^2}\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2 - \left(\frac{1}{4}\|\tilde{f}\|_{L^2_\rho}^2 + \frac{2\lambda^2}{\lambda_1+\lambda}\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2\right)\sqrt{\frac{\log N}{N}} - C_1\frac{\log N}{N}.$$

*For $\lambda \to 0$, we have a simpler bound: with a probability of at least $1 - 2/N$, we have*

$$\lim_{\lambda\to 0} bias \leq \left(1 + \frac{\log N}{N}\right)\tilde{\gamma}_{>M}^2 + 6\tilde{\gamma}_{>M}^2\left(\frac{\log N}{N}\right)^{\frac{3}{2}};$$

$$\lim_{\lambda\to 0} bias \geq \left(1 - \frac{\log N}{N}\right)\tilde{\gamma}_{>M}^2 - 6\tilde{\gamma}_{>M}^2\left(\frac{\log N}{N}\right)^{\frac{3}{2}}. \tag{35}$$

*For $\tilde{\gamma}_{>M}^2 = 0$, that is $\tilde{f} \in \mathcal{H}$, we have a simpler upper bound on bias: with a probability of at least $1 - 2/N$, we have*

$$bias \leq \lambda\|\tilde{f}\|_{\mathcal{H}}^2\left(1 + 2\sqrt{\frac{\log N}{N}}\right) + C_1\frac{\log N}{N};$$

$$bias \geq \frac{\lambda^2\lambda_M}{(\lambda_M+\lambda)^2}\|\tilde{f}\|_{\mathcal{H}}^2\left(1 - 2\sqrt{\frac{\log N}{N}}\right) - C_1\frac{\log N}{N}. \tag{36}$$

*Proof.* By Proposition C.13 and Lemma C.19,

$$\text{fitting error} \leq (1+2\delta)\left\|\bar{\mathbf{P}}w\right\|_2^2 + \|w\|_2^2\,\delta^2 p(\delta) \tag{37}$$

$$\leq (1+2\delta)(\lambda\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2 + \frac{1}{2}|\tilde{\gamma}_{>M}|\|\tilde{f}_{\leq M}\|_{L^2_\rho}\|\boldsymbol{E}\|_2 + \tilde{\gamma}_{>M}^2\|\boldsymbol{E}\|_2^2) + \|w\|_2^2\,\delta^2 p(\delta) \tag{38}$$

$$\leq (1+2\delta)(\lambda\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2 + \frac{1}{4}\|\tilde{f}\|_{L^2_\rho}^2\|\boldsymbol{E}\|_2 + \tilde{\gamma}_{>M}^2\|\boldsymbol{E}\|_2^2) + \|w\|_2^2\,\delta^2 p(\delta). \tag{39}$$

where in line (37), we use Proposition C.13; in line (38), we use Lemma C.19; in line (38), we use the fact that $2ab \leq a^2 + b^2$ where $a = |\tilde{\gamma}_{>M}|, b = \|\tilde{f}_{\leq M}\|_{L^2_\rho}$.

Now we apply the concentration result in Lemma C.17: with a probability of at least $1 - 2/N$:

$$\text{fitting error} \leq \left(1 + 2\sqrt{\frac{\log N}{N}}\right)\lambda\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2 + \frac{1}{4}\|\tilde{f}\|_{L^2_\rho}^2\sqrt{\frac{\log N}{N}}$$

$$+ \frac{\log N}{N}\left(\|w\|_2^2\,p(\delta) + (1+2\delta)\tilde{\gamma}_{>M}^2 + \frac{1}{2}\|\tilde{f}\|_{L^2_\rho}^2\right)$$

$$\leq \lambda\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2 + \left(\frac{1}{4}\|\tilde{f}\|_{L^2_\rho}^2 + 2\lambda\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2\right)\sqrt{\frac{\log N}{N}} + C_1\frac{\log N}{N},$$

where we choose $C_1 > 0$ to be such that:

$$\|w\|_2^2\,p(\delta) + (1+2\delta)\tilde{\gamma}_{>M}^2 + \frac{1}{2}\|\tilde{f}\|_{L^2_\rho}^2 \leq \|w\|_2^2\,p(\delta) + \left(1 + 2\delta + \frac{1}{2}\right)\|\tilde{f}\|_{L^2_\rho}^2$$

$$\leq \|w\|_2^2\,p\left(\frac{1}{2}\right) + \left(1 + 2\cdot\frac{1}{2} + \frac{1}{2}\right)\|\tilde{f}\|_{L^2_\rho}^2$$

$$\leq 8\left(\lambda\bar{r}\sqrt{M} + |\tilde{\gamma}_{>M}|\|E\|_2\right)^2 + \frac{5}{2}\|\tilde{f}\|_{L^2_\rho}^2$$

$$\leq 8\left(\lambda\bar{r}\sqrt{M} + \frac{1}{2}|\tilde{\gamma}_{>M}|\right)^2 + \frac{5}{2}\|\tilde{f}\|_{L^2_\rho}^2 < C_1.$$

579 Hence we have an upper bound for the bias. We argue similarly for the lower bound:

$$\text{fitting error} \geq \left(1 - 2\sqrt{\frac{\log N}{N}}\right) \frac{\lambda^2 \lambda_M}{(\lambda_M + \lambda)^2} \|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2 - \frac{1}{4}\|\tilde{f}\|_{L_\rho^2}^2 \sqrt{\frac{\log N}{N}}$$

$$- \frac{\log N}{N} \left(\|w\|_2^2 \, p(\delta) + (1 + 2\delta)\tilde{\gamma}_{>M}^2 + |\tilde{\gamma}_{>M}|\|\tilde{f}_{\leq M}\|_{L_\rho^2}\right)$$

$$\geq \frac{\lambda^2 \lambda_M}{(\lambda_M + \lambda)^2}\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2 - \left(\frac{1}{4}\|\tilde{f}\|_{L_\rho^2}^2 + \frac{2\lambda^2}{\lambda_1 + \lambda}\|\tilde{f}_{\leq M}\|_{\mathcal{H}}^2\right)\sqrt{\frac{\log N}{N}} - C_1 \frac{\log N}{N}.$$

580 For $\lambda \to 0$, note that $w \to -\tilde{\gamma}_{>M}\boldsymbol{E}$. This yields

$$\lim_{\lambda \to 0} \text{fitting error} \leq \lim_{\lambda \to 0}\left\{(1 + 2\delta)\left\|\bar{\mathbf{P}}w\right\|_2^2 + \|w\|_2^2\,\delta^2 p(\delta)\right\}$$

$$= (1 + 2\delta)\left\|-\tilde{\gamma}_{>M}\boldsymbol{E}\right\|_2^2 + \left\|-\tilde{\gamma}_{>M}\boldsymbol{E}\right\|_2^2\,\delta^2 p(\delta)$$

$$= \tilde{\gamma}_{>M}^2\left\|\boldsymbol{E}\right\|_2^2\left(1 + \delta(2 + \delta p(\delta))\right).$$

581 Hence, by plugging in $\delta < \frac{1}{2}$, with probability of at least $1 - 2/N$,

$$\lim_{\lambda \to 0} \text{fitting error} \leq \tilde{\gamma}_{>M}^2 \frac{\log N}{N}\left(1 + 6\sqrt{\frac{\log N}{N}}\right)$$

$$\lim_{\lambda \to 0} \text{bias} \leq \left(1 + \frac{\log N}{N}\right)\tilde{\gamma}_{>M}^2 + 6\tilde{\gamma}_{>M}^2\left(\frac{\log N}{N}\right)^{\frac{3}{2}}.$$

582 For lower bound, it follows similarly:

$$\lim_{\lambda \to 0} \text{bias} \geq \left(1 - \frac{\log N}{N}\right)\tilde{\gamma}_{>M}^2 - 6\tilde{\gamma}_{>M}^2\left(\frac{\log N}{N}\right)^{\frac{3}{2}},$$

583 and we obtain line (35). For the case where $\tilde{\gamma}_{>M} = 0$, recalculate and simplify line (38) to obtain
584 line (36). □

### C.4.2 Refined Bounds on Variance

586 Similarly, we can refine Theorem C.14 to get a bound on the variance:

587 **Theorem C.21.** *For $N > \max\left\{(12G)^4(M+1)^2, 9\right\}$, and set $C_2 = 12$ (independent to $N$), with a*
588 *probability of at least $1 - 2/N$, we have the upper and lower bounds of variance:*

$$\text{variance} \leq \sigma^2 \frac{M}{N}\left(1 + \sqrt{\frac{\log N}{N}} + C_2 \frac{\log N}{N}\right);$$

$$\text{variance} \geq \frac{\lambda_M^2}{(\lambda_M + \lambda)^2}\sigma^2 \frac{M}{N}\left(1 - \sqrt{\frac{\log N}{N}}\right) - C_2\sigma^2 \frac{M}{N}\frac{\log N}{N}.$$

589 *Proof.* We argue analogously as in Theorem C.20: by Proposition C.14 and Lemma C.17, we have

$$\text{variance} \leq (1 + \delta)\frac{\sigma^2}{N}\sum_{k=1}^{M}\frac{\lambda_k^2}{(\lambda_k + \lambda)^2} + M\frac{\sigma^2}{N}(1 + \delta)\delta^2 p(\delta)$$

$$\leq (1 + \delta)\sigma^2\frac{M}{N} + \sigma^2\frac{M}{N}(1 + \delta)\delta^2 p(\delta)$$

$$\leq \left(1 + \sqrt{\frac{\log N}{N}}\right)\sigma^2\frac{M}{N} + \sigma^2\frac{M}{N}\frac{\log N}{N}\left(1 + \frac{1}{2}\right)p\left(\frac{1}{2}\right)$$

$$\leq \left(1 + \sqrt{\frac{\log N}{N}}\right)\sigma^2\frac{M}{N} + 12\sigma^2\frac{M}{N}\frac{\log N}{N}$$

590 Hence we can choose $C_2 = 12$. For the lower bound, since $\frac{\lambda_k^2}{(\lambda_k + \lambda)^2} > \frac{\lambda_M^2}{(\lambda_M + \lambda)^2}$, we have

$$\text{variance} \geq (1 - \delta)\frac{\sigma^2}{N} \sum_{k=1}^{M} \frac{\lambda_k^2}{(\lambda_k + \lambda)^2} - M\frac{\sigma^2}{N}(1 + \delta)\delta^2 p(\delta)$$

$$\geq \frac{\lambda_M^2}{(\lambda_M + \lambda)^2}\sigma^2 \frac{M}{N}\left(1 - \sqrt{\frac{\log N}{N}}\right) - 12\sigma^2 \frac{M}{N}\frac{\log N}{N}.$$

591 □

592 Note that in both Theorems C.20 and C.21, the constants $C_1, C_2 > 0$ is not optimized.

# D Numerical Validation

594 In this section, we illustrate our result for KRR with two different finite rank kernels.

## D.1 Truncated NTK

595 First, we need to define a finite-rank kernel $K : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$. We set $\mathcal{X} = \mathbb{S}^1 \subset \mathbb{R}^2$. By reparametrization, we write $\mathbb{S}^1 \cong [0, 2\pi]/_{0 \sim 2\pi}$. We assume the data are drawn uniformly on the circle, that is $\rho_{\mathcal{X}} = \text{unif}[\mathbb{S}^1]$. We can use the Fourier functions $\cos(k\cdot), \sin(k\cdot)$ as the orthogonal eigenfunctions of the kernel. We define the NTK

$$K^{(\infty)}(\theta, \theta') \stackrel{\text{def.}}{=} \frac{\cos(\theta - \theta')(\pi - |\theta - \theta'|)}{2\pi}$$

596 for all $\theta, \theta' \in [0, 2\pi]$. 3) We choose a rank-$M$ truncation $K(\theta, \theta') = \sum_{k=1}^{M} \lambda_k \psi_k(\theta)\psi_k(\theta')$ for all $\theta, \theta' \in [0, 2\pi]$. For the first few eigenvalues of the kernel, please see Table 3 for example. Before

| $k$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | $\infty$ |
|---|---|---|---|---|---|---|---|---|
| $\lambda_k$ | $\frac{1}{\pi^2}$ | $\frac{1}{8}$ | $\frac{1}{8}$ | $\frac{5}{9\pi^2}$ | $\frac{5}{9\pi^2}$ | $\frac{17}{225\pi^2}$ | $\frac{17}{225\pi^2}$ | - |
| $\psi_k(\theta)$ | 1 | $\sqrt{2}\cos(\theta)$ | $\sqrt{2}\sin(\theta)$ | $\sqrt{2}\cos(2\theta)$ | $\sqrt{2}\sin(2\theta)$ | $\sqrt{2}\cos(4\theta)$ | $\sqrt{2}\sin(4\theta)$ | - |
| $\sum_{k'=0}^{k}\lambda_{k'}$ | 0.1013 | 0.2263 | 0.3513 | 0.4076 | 0.4639 | 0.4716 | 0.4792 | 0.5 |

Table 3: The first few eigenvalues of the NTK

597
598 proceeding to test error computation, we present a training example, Figure 4, to give readers an intuition on the truncated NTK (tNTK).
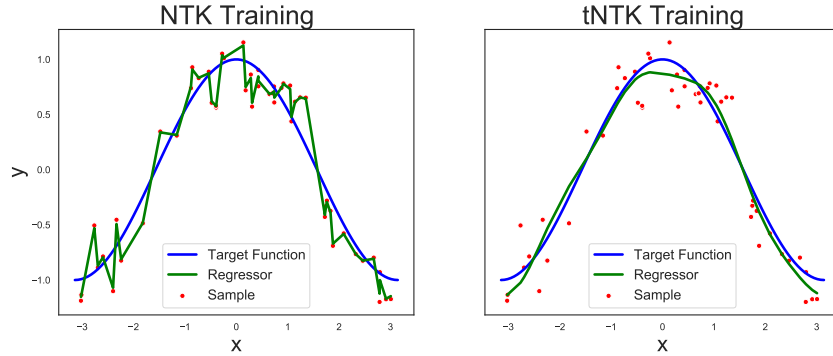


Figure 4: (left): NTK training; (right): tNTK training where $N = 50, M = 7$. $\sigma^2 = 0.05, \lambda = \sigma^2/N$.

599

ocr

## D.2 Test Error Computations

In the following tNTK training, we set the hyperparameters as follows:

**Target function**   We choose a simple target function $\tilde{f}(x) = \cos x = \frac{1}{\sqrt{2}}\psi_2(x)$. Throughout the experiment, we set the noise variance $\sigma^2 = 0.05$.

**Ridge**   We choose $\lambda = \frac{\sigma^2}{N}$. In Figure 5 (left), we set $N = 50, \lambda = 0.05/50$ for tNTK training; (right) we set set $\lambda = 0.05/50$ for varying $N$ from 10 to 200.

**Error bars**   In Figure 7 (right), for each value of $N$, we run over 10 iterations of random samples and compute the test error. The error bars are shown as the difference between the upper and the lower quartiles.



Figure 5: (left): tNTK training; (right): the decay of test error as $N$ varies.

**Lower bound**   See the subsection below.

## D.3 Bound Comparison

We continue with the experiment on the tNTK this time with varying $N$ and compare our upper bound with [6].

**Upper bounds**   In Figure 6, the expression of Bach's and our upper bounds are directly computed:

$$\text{Bach's upper bound} = 4\lambda\|\tilde{f}\|_{\mathcal{H}}^2 + \frac{8\sigma^2 R^2}{\lambda N}(1 + 2\log N)$$

$$\text{Our upper bound without residue} = \lambda\|\tilde{f}\|_{\mathcal{H}}^2\left(1 + 2\sqrt{\frac{\log N}{N}}\right),$$

where the constants $\|\tilde{f}\|_{\mathcal{H}}^2$ and $R^2$ can be computed directed from the choice of kernel and target function. For simplicity reason, we drop the residue term $C_1\frac{\log N}{N}$ since it is overshadowed by the other terms and the constant $C_1$ is not optimized.

## D.4 Legendre Kernel

To illustrate the bounds with another finite-rank, we choose a simple legendre kernel (LK):

$$K(x, z) = \sum_{k=0}^{M} \lambda_k P_k(x) P_k(z)$$

where $P_k$ is the Legendre polynomial of degree $k$, and $\lambda_k > 0$ are the eigenvalues.
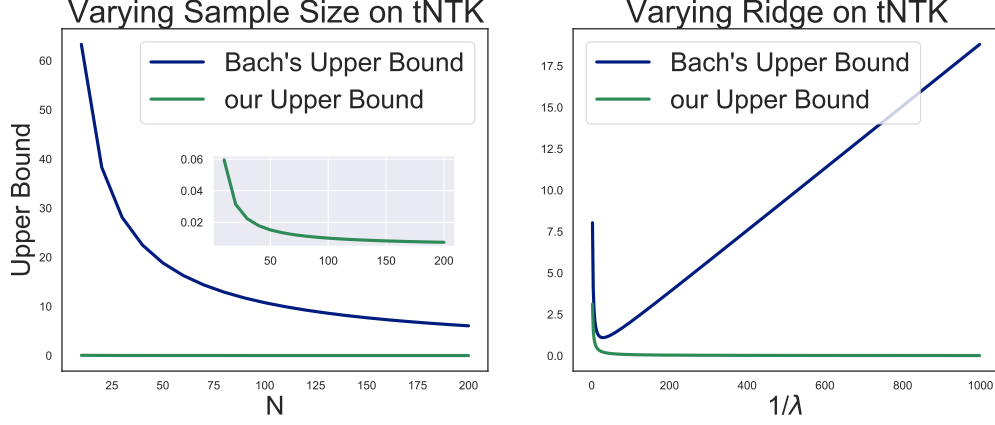
28

Figure 6: Bound improvement on tNTK. Residue term is dropped in our bound.

**Eigenvalues** To better compare the Legendre kernel $K$ with the NTK, we choose $\lambda_k = C \cdot (k+1)^{-2}$ of quadratic decay such that the spectral sums are the same: $\sum_{k=0}^{\infty} \lambda_k = 0.5$. Hence we choose $C = 0.5 / \sum_{k=1}^{\infty} k^{-2} = \frac{3}{\pi^2}$.

**Target function** We choose a simple target function $\tilde{f}(x) = x^2 = \frac{1}{3}P_0(x) + \frac{2}{3}P_2(x)$. Throughout the experiment, we set the noise variance $\sigma^2 = 0.05$.

### D.5 Test Error Computation

**Ridge** As before, our bound suggests that, to balance the bias and the variance with a fixed $N$, we can choose $\lambda = \frac{\sigma^2}{N}$. In Figure 7 (left), we set $N = 50, \lambda = 0.05/50$ for KRR training; (right) we set set $\lambda = 0.05/50$ for varying $N$ from 10 to 200.

**Error bars** In Figure 7 (right), for each value of $N$, we run over 10 iterations of random samples and compute the test error. The error bars are shown as the different between the upper and the lower quartiles. The median is taken as average.
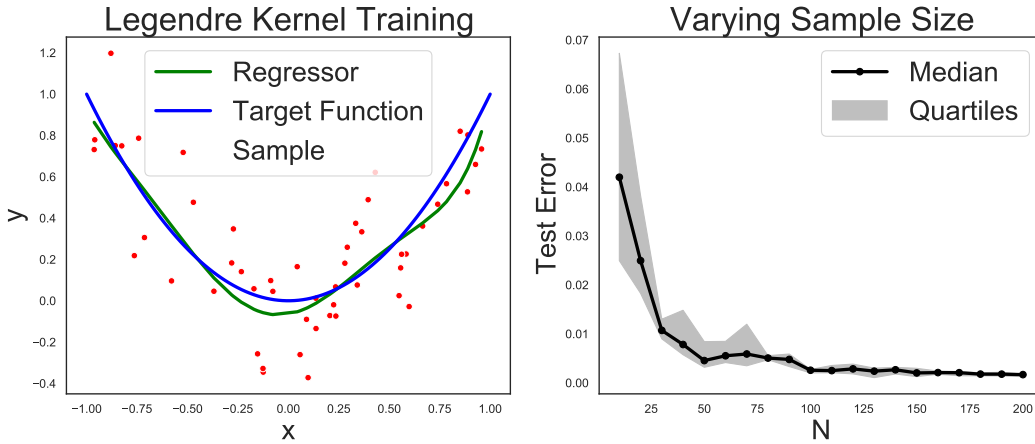


Figure 7: (left): LK training; (right): the decay of test error as $N$ varies. Same as Figure 1.

632 **Upper bounds**  In Figure 8, the expression of Bach's and our upper bounds are directly computed:

$$\text{Bach's upper bound} = 4\lambda\|\tilde{f}\|_{\mathcal{H}}^2 + \frac{8\sigma^2 R^2}{\lambda N}(1 + 2\log N)$$

$$\text{Our upper bound without residue} = \lambda\|\tilde{f}\|_{\mathcal{H}}^2\left(1 + 2\sqrt{\frac{\log N}{N}}\right),$$

633 where the constants $\|\tilde{f}\|_{\mathcal{H}}^2$ and $R^2$ can be computed directed from the choice of kernel and target
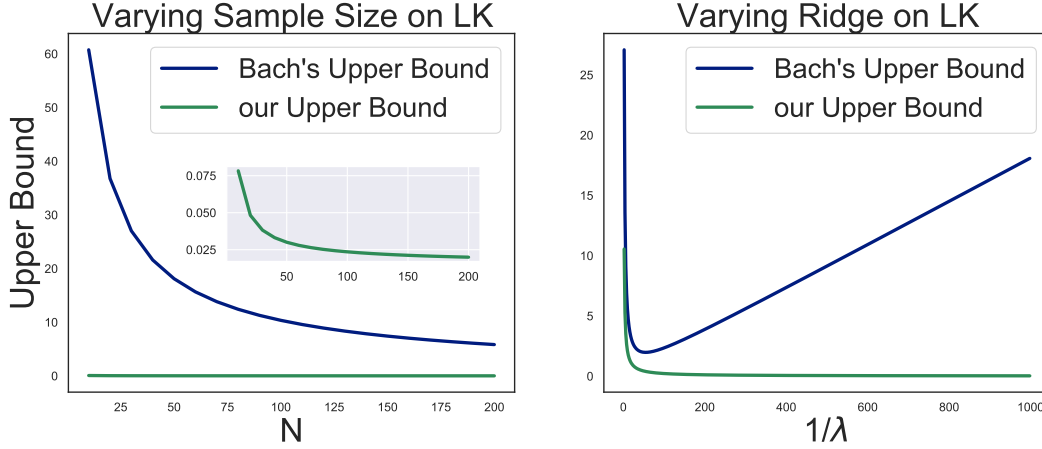634 function.



Figure 8: Bound improvement on LK. Same as Figure 2.

635 **Lower Bound**  Last but not least, we need to show our lower bound is valid. To see this clearly, we
636 need to write the bound in exact sums instead of in HKRS norm square $\|\tilde{f}\|_{\mathcal{H}}^2$: namely, we compute $I$

$$\frac{\lambda^2\lambda_M}{(\lambda_M + \lambda)^2}\|\tilde{f}\|_{\mathcal{H}}^2 \leq I = \lambda^2\sum_{k=1}^{M}\frac{\tilde{\gamma}_k^2}{(\lambda_k + \lambda)^2} \leq \lambda\|\tilde{f}\|_{\mathcal{H}}^2, \tag{40}$$

637 instead of using the inequality (40) in Lemma C.19; and

$$M\frac{\lambda_M^2}{(\lambda_M + \lambda)^2} \leq \sum_{k=1}^{M}\frac{\lambda_k^2}{(\lambda_k + \lambda)^2} \leq M, \tag{41}$$

638 instead of using the inequality (40) in Theorem C.21. Then we can compute our bounds as:

$$\text{Our upper bound without residue} = \lambda^2 I\left(1 + 2\sqrt{\frac{\log N}{N}}\right) + \frac{\sigma^2}{N}\sum_{k=1}^{M}\frac{\lambda_k^2}{(\lambda_k + \lambda)^2}\left(1 + \sqrt{\frac{\log N}{N}}\right),$$

$$\text{Our lower bound without residue} = \lambda^2 I\left(1 - 2\sqrt{\frac{\log N}{N}}\right) + \frac{\sigma^2}{N}\sum_{k=1}^{M}\frac{\lambda_k^2}{(\lambda_k + \lambda)^2}\left(1 - \sqrt{\frac{\log N}{N}}\right),$$

639 and we drop the residue terms $C_1\frac{\log N}{N}$ and $C_2\frac{\sigma^2}{N}M\frac{\log N}{N}$ by the same reason as before. From Figure
640 3, we can see that our bounds precisely describe the decay of the test error. Our bounds are not
641 'bounding' the test errors in smaller instances due to the absence of the residue terms, which increases
642 the interval of confidence of our approximation. But for larger instances, say $N > 100$, all upper and
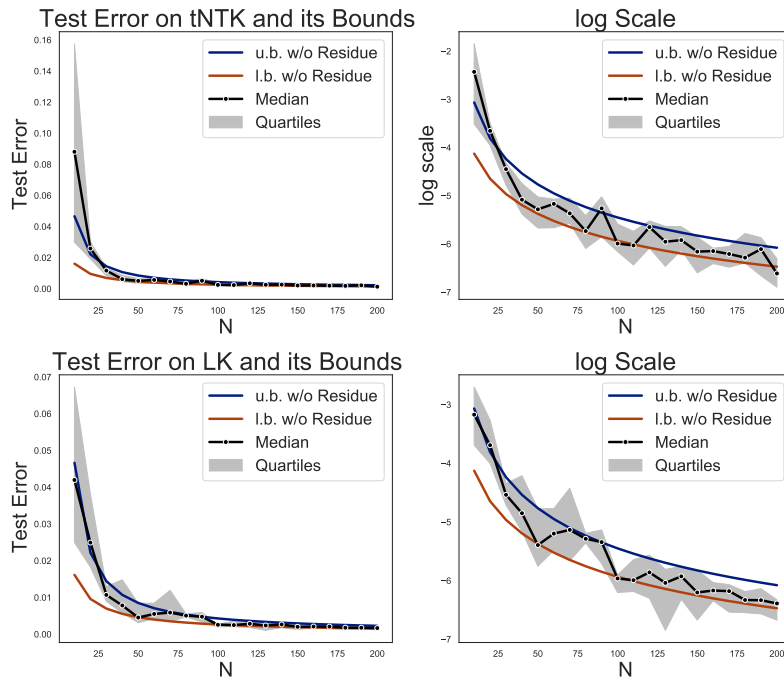643 lower bounds, and the averaged test error converge to the same limit.

30

Figure 9: Our bounds comparing to the averaged test error with varying $N$, over 10 iterations. Same as Figure 3.