

Convergence Rates for Learning Linear Operators from Noisy Data*

Maarten V. de Hoop[†], Nikola B. Kovachki[‡], Nicholas H. Nelsen[‡], and Andrew M. Stuart[‡]

Abstract. We study the Bayesian inverse problem of learning a linear operator on a Hilbert space from its noisy pointwise evaluations on random input data. Our framework assumes that this target operator is self-adjoint and diagonal in a basis shared with the Gaussian prior and noise covariance operators arising from the imposed statistical model and is able to handle target operators that are compact, bounded, or even unbounded. We establish posterior contraction rates with respect to a family of Bochner norms as the number of data tend to infinity and derive related lower bounds on the estimation error. In the large data limit, we also provide asymptotic convergence rates of suitably defined excess risk and generalization gap functionals associated with the posterior mean point estimator. In doing so, we connect the posterior consistency results to nonparametric learning theory. Furthermore, these convergence rates highlight and quantify the difficulty of learning unbounded linear operators in comparison with the learning of bounded or compact ones. Numerical experiments confirm the theory and demonstrate that similar conclusions may be expected in more general problem settings.

Key words. operator regression, linear inverse problems, Bayesian inference, posterior consistency, learning theory

AMS subject classifications. 62G20, 62C10, 68T05, 47A62

1. Introduction. Learning operators between Hilbert spaces provides a natural framework for the application of tools from supervised learning to the development of surrogate models that accelerate scientific computation by approximating existing expensive models and to the discovery of new models consistent with observed data when no model exists. In order to develop a deeper understanding of operator learning, this paper is concerned with nonparametric regression under random design on a separable infinite-dimensional Hilbert space H . We consider the problem of learning L , an unknown (possibly unbounded and in general densely defined) self-adjoint linear operator on H , from data pairs $\{(x_n, y_n)\}_{n=1}^N$ related by

$$(1.1) \quad y_n = Lx_n + \eta_n, \quad n \in \{1, \dots, N\}.$$

Here the $\{\eta_n\}_{n=1}^N$ represent noise and $N \in \mathbb{N}$ is referred to as the sample size.

The estimation of L from the data (1.1) is in general an ill-posed linear inverse problem [24]. We say that L is more difficult to learn than another operator L' if the estimation of L

*Submitted to the editors August 27, 2021.

Funding: MVdH is supported by the Simons Foundation under the MATH + X program, U.S. Department of Energy, Office of Basic Energy Sciences, Chemical Sciences, Geosciences, and Biosciences Division under grant number DE-SC0020345, the National Science Foundation (NSF) under grant DMS-1815143, and the corporate members of the Geo-Mathematical Imaging Group at Rice University. NHN is supported by the NSF Graduate Research Fellowship Program under grant DGE-1745301. AMS is supported by NSF (grant DMS-1818977). NBK, NHN, and AMS are supported by NSF (grant AGS-1835860) and ONR (grant N00014-19-1-2408).

[†]Department of Computational and Applied Mathematics, Rice University, Houston, TX 77005 USA (mdehoop@rice.edu).

[‡]Division of Engineering and Applied Science, California Institute of Technology, Pasadena, CA 91125 USA (nkovachki@caltech.edu, nnelsen@caltech.edu, astuart@caltech.edu).

requires larger sample sizes N to achieve a fixed error tolerance relative to that of L' , that is, L has worse sample complexity. Broadly, our work aims to provide an answer to the question:

What is the relative difficulty of learning different types of linear operators?

The case in which H comprises (e.g., real-valued) functions over a domain $D \subset \mathbb{R}^d$ is a particular focus; for example, it is of interest to understand the relative difficulty of learning forward and inverse operators arising from partial differential equations (PDEs). Indeed, there is an emerging body of work on supervised learning between Banach spaces focused primarily on forward, typically nonlinear, PDE solution operators [2, 11, 36, 38, 39, 44, 46, 52]. In the context of dynamical systems, there is also literature focused on learning the Koopman operator or its generator, both linear operators, from data generated by the underlying dynamics for system identification or forecasting [14, 27, 32, 33, 47]. Finally, there is interest in using surrogate forward models to speed up (e.g., Bayesian) inversion techniques [38] and in directly learning regularizers for inversion, or even the regularized inverse solution operator itself [9].

The study of linear function-to-function models within functional data analysis (FDA) [49] is also a well-established area; see [20, 30, 51, 58] and references therein. Much of this work concerns the setting $H = L^2((0, 1); \mathbb{R})$ and linear models based on kernel integral operators under colored noise. The operator estimation is then reduced to learning the kernel function, usually in a reproducing kernel Hilbert space (RKHS) framework. Linear operator learning has also been considered in machine learning [1], particularly in the context of conditional expectation operators [43] and conditional mean embeddings in RKHS [29, 31, 53].

The authors in [30, 51] study functional linear regression between two Hilbert spaces and work directly with a spectral representation of the estimator which allows them to obtain consistency of the prediction error assuming only boundedness of the true operator [30], rather than compactness as assumed in much of the FDA literature; convergence rates are established in [51]. While unbounded operators are not considered in these works, their approaches could likely be modified to handle them. Relatedly, the authors of [56] and [13] have similar motivations to us; the former establishes sample complexities for learning Schatten-class compact operators (motivated by inverse problem solution maps) while the latter for learning compact operators associated to Green's functions of elliptic PDEs (motivated by PDE discovery). Our theory also treats these types of operators but goes further by proving sample complexities for the direct learning of *unbounded operators*, which are of primary interest in these papers (the inverse map in the former and the partial differential operator in the latter).

We consider (1.1) in an idealized setting in which the eigenbases of the target ground truth operator and data covariance operator are assumed known. Although quite strong, these assumptions may be realized in practice when, for example, prior knowledge is available that the covariance commutes with the target (hence simultaneously diagonalizable in the same eigenbasis) or that the target obeys certain physical principles (e.g., commutes with translation operators). This allows us to work in coordinates with respect to the eigenbasis and hence reduce inference of the operator to that of its eigenvalue sequence; similar ideas were recently applied to learn a differential operator arising in an advection-diffusion model [48]. Our proof techniques in this linear diagonal setting follow the program set forth in [35], which is concerned with the Bayesian approach to ill-posed inverse problems in a diagonalized problem setting. Although the authors of [17, 50] established optimal convergence rates for regularized

least squares algorithms in statistical direct and inverse learning problems with both infinite-dimensional input *and* output spaces, these results do not apply to our simpler linear operator regression setting (1.1). This is because their requirement that the appropriately defined point evaluation map (here acting on linear operators) is Hilbert–Schmidt never holds with H infinite-dimensional, since the implied operator-valued kernel is not trace-class; the limitations of this strong assumption were also acknowledged in more general RKHS settings [29, 43].

1.1. Our Contributions. It is clear that much of the work regarding (1.1) has focused on bounded or compact L . Instead, we provide a unified framework accommodating the learning of compact, bounded, and unbounded operators that allows us to compare the difficulty of learning operators defined by both forward and inverse problems. In particular, we show that unbounded linear operators can be learned in a stable and consistent manner, but with associated convergence rates in the large data limit that are worse relative to those of their continuous (bounded or compact) counterparts (see Figure 1).

The primary contributions we make in this paper are now listed:

- (i) we formulate linear operator learning both as a nonparametric Bayesian inverse problem with non-compact forward map and as a statistical learning optimization problem;
- (ii) in the large sample limit, we prove convergence of the posterior estimator using dimension/discretization-independent asymptotic upper and lower error bounds in a family of suitable norms, in expectation and with high probability over input data;
- (iii) our theory demonstrates that, as a point estimator, the posterior mean exhibits no loss in statistical performance when compared to the full posterior solution;
- (iv) thus, as a byproduct of our posterior mean analysis, we analogously establish asymptotic convergence rates of the excess risk and generalization gap;
- (v) we perform numerical experiments to illustrate the learning of compact, bounded, and unbounded linear operators from noisy data, and the results both support the theory and confirm our conclusions beyond the confines of our idealized theorem setting.

The remainder of the paper is organized as follows. In Section 2, we describe the theoretical setting in which we work; the formulation in (i) is addressed in Subsection 2.2, 2.3, and Proposition 2.2. Our main theoretical results are in Section 3, where we provide asymptotic convergence rates of the inverse problem solution to the truth; contributions (ii)–(iii) are found in Theorems 3.1, 3.3, 3.5, and 3.6 and (iv) in Theorems 3.7 to 3.9. We discuss the theory in Section 4, and numerical results (v) that illustrate, support, and extend beyond the convergence theory are provided in Section 5. Our concluding remarks follow in Section 6. Appendix A is devoted to proofs of the main results, with supporting lemmas in Appendix B.

2. Problem Formulation. After overviewing some notation in Subsection 2.1, we recast (1.1) as a random Bayesian inverse problem in Subsection 2.2. Subsection 2.3 gives an optimization perspective and defines expected risk and generalization gap in the infinite-dimensional setting. In Subsection 2.4, we describe our main assumptions and explicitly characterize the conditional posterior probability measure under these assumptions.

2.1. Preliminaries. Let $(H, \langle \cdot, \cdot \rangle, \|\cdot\|)$ be a real, separable, infinite-dimensional Hilbert space. For any self-adjoint positive definite linear operator A on H , we define

$$\langle \cdot, \cdot \rangle_A := \langle A^{-1/2} \cdot, A^{-1/2} \cdot \rangle, \quad \|\cdot\|_A := \|A^{-1/2} \cdot\|.$$

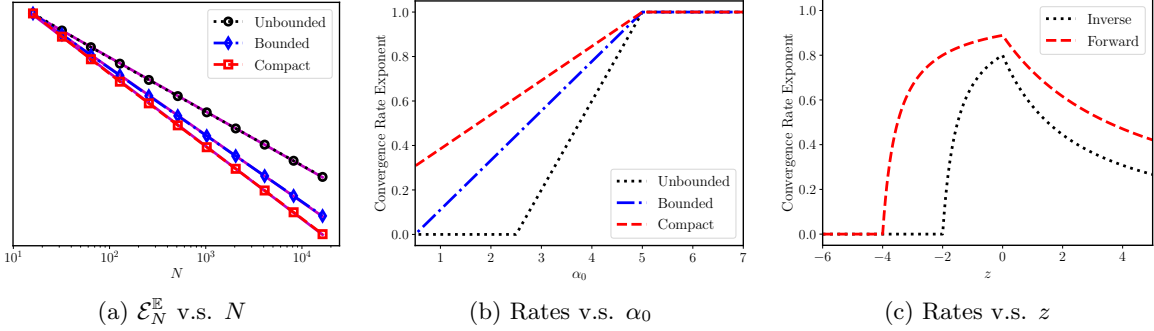


Figure 1. The difficulty of learning unbounded linear operators (as discussed in Section 4): Figure 1a displays the observed convergence rates (same as Table 1, Column 4) for the expected excess risk (Theorem 3.7) corresponding to unbounded $(-\Delta)$, bounded (Id) , and compact $((-\Delta)^{-1})$ true target operators, where the vertical axis is rescaled to allow comparison of the three differing slopes. In Figure 1b, the rate approaches zero the soonest for unbounded operator $-\Delta$ as the out-of-distribution test measure smoothness α_0 is decreased (Theorem 3.1). Figure 1c shows that the rate for learning the unbounded “inverse map” $-\Delta$ (with training measure smoothness $\alpha_1 = 4.5$) is always worse than that for learning the compact “forward map” $(-\Delta)^{-1}$ (with $\alpha_1 = 2.5$) as the shift z in the prior regularity, with respect to the regularity of the truth, is varied (Theorem 3.7).

The set $\mathcal{L}(H_1; H_2)$ is the space of bounded linear operators mapping Hilbert space H_1 into H_2 , and when $H_1 = H_2 = H$, we write $\mathcal{L}(H)$. Similarly, the separable Hilbert space of Hilbert–Schmidt operators from H_1 to H_2 is denoted by $(\text{HS}(H_1; H_2), \langle \cdot, \cdot \rangle_{\text{HS}(H_1; H_2)}, \|\cdot\|_{\text{HS}(H_1; H_2)}) \subset \mathcal{L}(H_1; H_2)$, and when $H_1 = H_2 = H$, we write $(\text{HS}(H), \langle \cdot, \cdot \rangle_{\text{HS}}, \|\cdot\|_{\text{HS}})$. For any $a \in H_2, b \in H_1$, $a \otimes_{H_1} b \in \text{HS}(H_1; H_2)$ denotes the outer product defined by $(a \otimes_{H_1} b)c := \langle b, c \rangle_{H_1} a$ for any $c \in H_1$, and we also use the shorthand $a \otimes b \in \text{HS}(H)$. Additionally, when dealing with a possibly unbounded linear operator T on H , we denoted its domain by the subspace $\mathcal{D}(T) \subseteq H$.

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a complete probability space. We use Bochner integrals throughout the article and when no confusion is caused by doing so, use the symbol \mathbb{E} with no other superscripts to denote averaging over all sources of randomness. We implicitly justify all exchanges of expectation and series with the Fubini–Tonelli theorem. The chi-square distribution with $n \in \mathbb{N}$ degrees of freedom is denoted by $\chi^2(n)$. We primarily work with centered Gaussian measures $\mathcal{N}(0, \mathcal{C})$ on measurable space $(H, \mathcal{B}(H))$, where $\mathcal{C} \in \mathcal{L}(H)$ is its symmetric, nonnegative covariance operator and $\mathcal{B}(H)$ is the standard Borel σ -algebra on H , understood to be defined with respect to common probability space $(\Omega, \mathcal{F}, \mathbb{P})$. Such a \mathcal{C} admits a sequence $\{\psi_j\}_{j \in \mathbb{N}}$ of eigenvectors forming an orthonormal basis of H and nonnegative eigenvalues $\{\theta_j^2\}_{j \in \mathbb{N}}$. Without loss of generality, we may assume that the eigenvalues are ordered to be nonincreasing; operator \mathcal{C} is necessarily trace-class on H , hence compact, so that $\{\theta_j^2\}_{j \in \mathbb{N}}$ is summable. A sample $h \sim \mathcal{N}(0, \mathcal{C})$ may be obtained from the Karhunen–Loève (KL) expansion $h = \sum_{j=1}^{\infty} \theta_j \xi_j \psi_j$, where $\{\xi_j\}_{j \in \mathbb{N}}$ is an independent and identical distributed (i.i.d.) sequence with $\xi_1 \sim \mathcal{N}(0, 1)$ [54, 55]. Thus, the coordinates $\langle \psi_j, h \rangle$ of h in eigenbasis $\{\psi_j\}$ are independent and distributed as $\mathcal{N}(0, \theta_j^2)$, and $\mathbb{E}\|h\|^2 = \text{tr}(\mathcal{C}) < \infty$. We also consider generalizations of Gaussian measures defined on spaces larger than H and on spaces of operators.

For $p, q \in \mathbb{R}$, we write $p \wedge q := \min\{p, q\}$. For two sequences $\{a_j\}$ and $\{b_j\}$ of nonnegative real numbers, we write $a_j \asymp b_j$ if the sequence $\{a_j/b_j\}$ is bounded away from zero and infinity

(uniformly over the index j), $a_j \lesssim b_j$ if $\{a_j/b_j\}$ is bounded, and $a_j \gtrsim b_j$ if $b_j \lesssim a_j$. Similarly, we use standard asymptotic notation, writing $a_n = O(b_n)$ as $n \rightarrow \infty$ if there exist $n', C' > 0$ such that $a_n \leq C'b_n$ for all $n \geq n'$, $a_n = \Omega(b_n)$ as $n \rightarrow \infty$ if $b_n = O(a_n)$, $a_n = \Theta(b_n)$ as $n \rightarrow \infty$ if both $a_n = O(b_n)$ and $a_n = \Omega(b_n)$, and $a_n = o(b_n)$ as $n \rightarrow \infty$ if $\lim_{n \rightarrow \infty} a_n/b_n = 0$. We make frequent use of the Sobolev-like sequence Hilbert spaces

$$\mathcal{H}^s(\mathbb{N}; \mathbb{R}) := \left\{ v: \mathbb{N} \rightarrow \mathbb{R} \mid \sum_{j=1}^{\infty} j^{2s} |v_j|^2 < \infty \right\}$$

for any $s \in \mathbb{R}$, equipped with the natural $\{j^s\}$ -weighted $\ell^2(\mathbb{N}; \mathbb{R})$ inner product and norm.

2.2. Bayesian Inversion. Consider the linear operator equation

$$(2.1) \quad y = Lx + \eta,$$

where $L: \mathcal{D}(L) \subseteq H \rightarrow H$ is a densely defined self-adjoint linear operator, $x \in H$, and $\eta \sim \pi_0$ is a random variable modeling observational noise. Rather than the standard inverse problem setting in which x is viewed as the unknown, we instead view x itself as a given *pointwise evaluation linear map on operators*: $L \mapsto Lx$; thus (2.1), or multiple such identities, defines an equation for *operator unknown* L . We consider the setting where x and η are Gaussian and independent, and draw N pairs $\{(x_n, \eta_n)\}_{n=1}^N$ to define a likelihood; we then assume that L is *a priori* Gaussian which makes the likelihood conjugate. To this end, it is useful to denote the N -fold product of H by H^N and define $H_{\mathcal{K}} := \text{Im}(\mathcal{K}^{1/2}) = \mathcal{K}^{1/2}H \subset H$, where $\mathcal{K} \in \mathcal{L}(H)$ is symmetric, positive definite, and trace-class on H ; equipped with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{K}}$, $H_{\mathcal{K}}$ is a RKHS and may also be viewed as the Cameron–Martin space of $\mathcal{N}(0, \mathcal{K})$ [55]. In our setting, an important role of \mathcal{K} is to ensure that the support of the prior on L is large enough to encompass unbounded operators. We make the following assumptions about (2.1).

Assumption 2.1. *It holds that $X := (x_1, \dots, x_N) \in H^N$, comprised of i.i.d. $x_n \sim \nu := \mathcal{N}(0, \mathcal{C}_1)$, and $E := (\eta_1, \dots, \eta_N)$, comprised of i.i.d. $\eta_n \sim \pi_0 := \mathcal{N}(0, \mathcal{C}_2)$, are independent centered Gaussians. The data covariance $\mathcal{C}_1: H \rightarrow H$ is a symmetric, positive definite linear operator that is trace-class on H , and the noise covariance operator $\mathcal{C}_2 \in \mathcal{L}(H)$ is symmetric and positive definite. Furthermore, for some $\mathcal{K} \in \mathcal{L}(H)$ symmetric, positive definite, and trace-class on H such that $H_{\mathcal{C}_1} \subsetneq H_{\mathcal{K}}$, there exists a Borel measurable linear operator $L^\dagger \in \text{HS}(H_{\mathcal{K}}; H)$ that is self-adjoint with domain $\mathcal{D}(L^\dagger) := \{h \in H: L^\dagger h \in H\} \supseteq H_{\mathcal{K}}$ dense in H , commutes with \mathcal{K} on $\mathcal{D}(L^\dagger)$, and generates the data $\{(x_n, y_n)\}_{n=1}^N$, $N \in \mathbb{N}$, where*

$$(2.2) \quad y_n = L^\dagger x_n + \eta_n, \quad n \in \{1, \dots, N\}.$$

We note that, although x is a Gaussian random variable in H (since \mathcal{C}_1 is assumed trace-class on H), η need not be (since \mathcal{C}_2 is not assumed trace-class on H) and must be interpreted weakly (see [3, Sec. 5.1], [5, 18, 35]). Viewing L as a random variable, for our prior model we place a centered Gaussian prior $\mu_0 := \mathcal{N}(0, \mathcal{C}_3)$ on L independent of x and η . The sense in which μ_0 is a proper Gaussian measure requires some care. The most natural Hilbert space of linear operators on H is the set of Hilbert–Schmidt operators $\text{HS}(H) \subset \mathcal{L}(H)$, and one could try to interpret μ_0 as a measure supported on $\text{HS}(H)$. However, this is unsatisfactory as it does not allow for prior knowledge that L^\dagger may fail to be continuous and thus not an

element of $\text{HS}(H)$. Instead, we take μ_0 as a proper Gaussian measure on the *larger space* $\text{HS}(H_{\mathcal{K}}; H) \supset \mathcal{L}(H)$ so that $\mathcal{C}_3 \in \mathcal{L}(\text{HS}(H_{\mathcal{K}}; H))$ is assumed symmetric, positive definite, and trace-class on $\text{HS}(H_{\mathcal{K}}; H)$, but not necessarily trace-class on $\text{HS}(H)$.

Recovery of the ground truth element L^\dagger from the given data (2.2) may be formulated in the notation of a traditional statistical inverse problem as follows. Concatenating on the index n , we write $Y := (y_1, \dots, y_N)$ and equip the product space H^N with the inner product $\langle U, V \rangle_{H^N} = \sum_{n=1}^N \langle u_n, v_n \rangle$ for any $U, V \in H^N$; this makes H^N a Hilbert space. Next, for $Z \in H_{\mathcal{K}}^N$ (the N -fold copy of $H_{\mathcal{K}}$), we define the linear map $K_Z \in \mathcal{L}(\text{HS}(H_{\mathcal{K}}; H); H^N)$ by $T \mapsto K_Z T := (Tz_1, \dots, Tz_N)$. We thus have the given data (2.2) in the standard form

$$(2.3) \quad Y = K_X L^\dagger + E,$$

where $X \sim \nu^{\otimes N}$ and $E \sim \pi_0^{\otimes N}$. K_X is analogous to the *forward map* in the terminology of inverse problems, except here it is a random variable. The following proposition shows that although this forward map is injective, contrary to usual inverse problems it is not compact.

Proposition 2.2. *K_X is injective ν -almost surely and for any $Z \in H_{\mathcal{K}}^N$; K_Z is not compact.*

Proof. For any nonzero $T \in \text{HS}(H_{\mathcal{K}}; H)$, as $H_{\mathcal{C}_1} \subset H_{\mathcal{K}}$ by Assumption 2.1, $T \in \mathcal{L}(H_{\mathcal{C}_1}; H)$ so $Tx_n \sim T_\# \nu$ i.i.d., $n \in \{1, \dots, N\}$. But $\nu(\{x \in H : Tx = 0\}) = (T_\# \nu)(\{y : y = 0\}) = 0$, so that $\mathbb{P}\{\omega \in \Omega : Tx_n(\omega) = 0 \ \forall n \leq N\} \leq \mathbb{P}\{\omega \in \Omega : Tx_1(\omega) = 0\} = 0$ proves the first assertion.

For the second, notice that $\text{HS}(H_{\mathcal{K}}; H)$ is a vector-valued RKHS [42] with operator-valued kernel $(z, z') \mapsto K_z K_{z'}^* = \langle z, z' \rangle_{\mathcal{K}} \text{Id}$, where $K_z : T \mapsto Tz$ and $K_z^* : h \mapsto h \otimes_{H_{\mathcal{K}}} z$. However, $\text{Id} \in \mathcal{L}(H)$ is not compact on H , which implies $K_z \in \mathcal{L}(\text{HS}(H_{\mathcal{K}}; H); H)$ is not compact either. Extending the argument from $z \in H_{\mathcal{K}}$ (for K_z) to $Z \in H_{\mathcal{K}}^N$ (for K_Z) completes the proof. ■

We assume that L, X , and E are *a priori* independent and consider the Bayesian inverse problem of finding the distribution of L , given X and Y , assuming that the random variables are related by (2.3) (with L^\dagger replaced by L). A rigorous justification of the posterior formulae to follow, namely that it is Gaussian with mean and covariance (2.5) and that it is absolutely continuous with respect to the prior μ_0 , with Radon–Nikodym derivative (2.6), may be obtained by application of [22, Theorems 32, 13, 37]. Because X and L are *a priori* independent, the posterior measure on L given X and Y , denoted by $\mu_N^{Y,X}$, is the same as that obtained when L is conditioned on Y with X fixed, almost surely (a.s.); this measure is Gaussian:

$$(2.4) \quad \mu_N^{Y,X} = \mathcal{N}(\bar{L}^{(N)}, \mathcal{C}^{(N)}).$$

The mean is the linear operator $\bar{L}^{(N)} = \mathbb{E}^{L \sim \mu_N^{Y,X}} L \in \text{HS}(H_{\mathcal{K}}; H)$ and $\mathcal{C}^{(N)} : \text{HS}(H_{\mathcal{K}}; H) \rightarrow \text{HS}(H_{\mathcal{K}}; H)$ is the covariance operator. If we define $I_{N \times N} \otimes \mathcal{C}_2$ to be the Kronecker product of $I_{N \times N} \in \mathbb{R}^{N \times N}$ and $\mathcal{C}_2 \in \mathcal{L}(H)$, which is a N -by- N block diagonal operator with diagonal entries \mathcal{C}_2 , then the posterior mean and covariance are given by

$$(2.5a) \quad \bar{L}^{(N)} = \mathcal{C}_3 K_X^* (I_{N \times N} \otimes \mathcal{C}_2 + K_X \mathcal{C}_3 K_X^*)^{-1} Y,$$

$$(2.5b) \quad \mathcal{C}^{(N)} = \mathcal{C}_3 - \mathcal{C}_3 K_X^* (I_{N \times N} \otimes \mathcal{C}_2 + K_X \mathcal{C}_3 K_X^*)^{-1} K_X \mathcal{C}_3$$

in our infinite-dimensional setting (see [35, Prop. 3.1], [37, 40]), and additionally,

$$(2.6) \quad \frac{d\mu_N^{Y,X}}{d\mu_0}(L) = \frac{1}{Z_N^{Y,X}} \exp\left(-\sum_{n=1}^N \frac{1}{2} \|Lx_n\|_{\mathcal{C}_2}^2 + \sum_{n=1}^N \langle Lx_n, y_n \rangle_{\mathcal{C}_2}\right),$$

where $Z_N^{Y,X}$ is the required normalization constant. The likelihood formula (2.6) motivates the learning theory formulation that follows in the next subsection.

Remark 2.3. Our analysis is not limited to the assumption that ν is a Gaussian measure on H ; the results in this paper would be virtually unchanged if ν was instead assumed to be given by a Karhunen–Loève random series expansion with non-Gaussian (but still independent) coefficients (see, e.g., [54]), provided they are sub-Gaussian. Indeed, since $x \sim \nu$ enters (2.3) only through the forward operator K_X , the inverse problem would still be linear Gaussian and hence the posterior would have the same form. However, the concentration inequalities and moment bounds appearing in Lemmas B.5 to B.7 within Appendix B are specific to the chi-square distribution, which arises because ν is Gaussian; to generalize, these would need to be replaced with results pertaining to more general sub-exponential distributions. \diamond

We now define two natural statistical approximations to L^\dagger : the *posterior mean estimator* $\bar{L}^{(N)}$ and the *posterior sample estimator* $L^{(N)} \sim \mu_N^{Y,X}$. The sense in which these estimators are close to L^\dagger , as a function of data volume N , measures the amount of data required for given level of estimation accuracy; this is the focus of our analysis. The usual norms (e.g., operator, Hilbert–Schmidt, and trace norms) do not usefully metrize estimation errors if L^\dagger is unbounded. Instead, the following weaker *Bochner norm* is proposed to measure such error. Let ν' be a centered probability measure supported on H with finite fourth moment, $\mathbb{E}^{x' \sim \nu'} \|x'\|^4 < \infty$, and denote its covariance operator by $\mathcal{C}_0 := \mathbb{E}^{x' \sim \nu'} [x' \otimes x']$. The identity $\langle v, Tu \rangle = \text{tr}(Tu \otimes v)$ for any linear $T: \mathcal{D}(T) \subseteq H \rightarrow H$ and for all $u \in \mathcal{D}(T)$, $v \in H$, yields

$$(2.7) \quad \|T\|_{L_{\nu'}^2(H;H)}^2 := \mathbb{E}^{x' \sim \nu'} \|Tx'\|^2 = \text{tr}(T\mathcal{C}_0^{1/2}(T\mathcal{C}_0^{1/2})^*) = \|T\mathcal{C}_0^{1/2}\|_{\text{HS}}^2 = \|T\|_{\text{HS}(H_{\mathcal{C}_0};H)}^2,$$

assuming \mathcal{C}_0 is trace-class. This motivates the next fact.

Proposition 2.4. Under Assumption 2.1, $L^\dagger \in \text{HS}(H_{\mathcal{C}_1}; H)$. In particular, $\nu(\mathcal{D}(L^\dagger)) = 1$.

Proof. By assumption, $L^\dagger \in \text{HS}(H_{\mathcal{K}}; H)$ and $H_{\mathcal{C}_1}$ is (strictly) contained in $H_{\mathcal{K}}$, so we immediately have $L^\dagger \in \mathcal{L}(H_{\mathcal{C}_1}; H)$. As $\langle h, \mathcal{C}_1 h \rangle \lesssim \langle h, \mathcal{K} h \rangle$ for all $h \in H$ [55, Lem. 6.15], it follows that $\sum_j \lambda_j^2 \langle L^\dagger \varphi_i, \phi_j \rangle^2 \lesssim \kappa_i^2 \sum_j \langle L^\dagger \varphi_i, \phi_j \rangle^2$ for all $i \in \mathbb{N}$, where $\{\phi_j\}, \{\varphi_j\}$ are the orthonormal eigenbases of $\mathcal{C}_1, \mathcal{K}$ corresponding to eigenvalues $\{\lambda_j^2\}, \{\kappa_j^2\}$, respectively. The facts that L^\dagger is self-adjoint and $\{\lambda_j \phi_j\}, \{\kappa_j \varphi_j\}$ form orthonormal bases of $H_{\mathcal{C}_1}, H_{\mathcal{K}}$, respectively, yield

$$\|L^\dagger\|_{\text{HS}(H_{\mathcal{C}_1}; H)}^2 = \sum_{i,j} \langle \varphi_i, L^\dagger(\lambda_j \phi_j) \rangle^2 \lesssim \sum_{i,j} \langle L^\dagger(\kappa_i \varphi_i), \phi_j \rangle^2 = \|L^\dagger\|_{\text{HS}(H_{\mathcal{K}}; H)}^2 < \infty.$$

Therefore, by (2.7), $L^\dagger x \in H$ ν -a.s. for $x \sim \nu$, that is, $\mathcal{D}(L^\dagger)$ has full measure under ν . \blacksquare

Our use of the weighted Hilbert–Schmidt space and its norm (2.7) is closely related to the notion of ν' -measurable linear operators [40], and similar considerations were made in [35, Prop. 3.2] for unbounded linear functionals. Since \mathcal{C}_0 is compact, the norm in (2.7) is weak in the sense that $\text{HS}(H_{\mathcal{C}_0}; H) \supset \mathcal{L}(H) \supset \text{HS}(H)$ [26, Sec. 2.2]. We further note that, if \mathcal{K} commutes with self-adjoint operator T , then the \mathcal{K} -weighted Hilbert–Schmidt norm of T corresponds to a weighted ℓ^2 norm on the eigenvalues of T , and we use this fact to perform explicit calculations in Subsection 2.4. Aligned with the notion of test error from machine learning, the $L_{\nu'}^2$ Bochner norm leads to the following definition.

Definition 2.5 (Test Error: Bayesian). *The test error of the posterior sample estimator is*

$$(2.8) \quad \mathbb{E}^{X \sim \nu^{\otimes N}, E \sim \pi_0^{\otimes N}} \mathbb{E}^{L \sim \mu_N^{Y,X}} \|L - L^\dagger\|_{L_{\nu'}^2(H;H)}^2.$$

The two outer expectations are with respect to the data, and the inner expectation is with respect to the Bayesian posterior. The definition of test error for the posterior mean is similar.

Definition 2.6 (Test Error: Point). *The test error of the posterior mean estimator is*

$$(2.9) \quad \mathbb{E}^{X \sim \nu^{\otimes N}, E \sim \pi_0^{\otimes N}} \|\bar{L}^{(N)} - L^\dagger\|_{L_{\nu'}^2(H;H)}^2.$$

We say that (2.8) or (2.9) tests *in-distribution* if $\nu' = \nu$ and *out-of-distribution* otherwise; these quantities are also referred to as *prediction error* [16]. In Section 3, we study the asymptotic performance of the estimators $L^{(N)} \sim \mu_N^{Y,X}$ and $\bar{L}^{(N)}$ as $N \rightarrow \infty$, using the notion of posterior contraction with respect to the $L_{\nu'}^2(H;H)$ error, leading to analysis of (2.8) and (2.9).

2.3. Statistical Learning. We now adopt a statistical learning theory perspective for the operator regression problem. Let \mathcal{P} denote the joint probability measure implied by data model (2.1) and Assumption 2.1. The given data is then $(x_n, y_n) \sim \mathcal{P}$ i.i.d., $n \in \{1, \dots, N\}$. Since regression is the focus, it is natural to work with the square loss on H so that $\mathbb{E}^{(x,y) \sim \mathcal{P}} \frac{1}{2} \|y - Lx\|^2$ and $\frac{1}{N} \sum_{n=1}^N \frac{1}{2} \|y_n - Lx_n\|^2$ define the expected risk and empirical risk, respectively, with minimizers of the latter viewed as point estimators. These two standard definitions need careful interpretation, however, since under the infinite-dimensional model $y = Lx + \eta$ with \mathcal{C}_2 *not* trace-class on H , the cylinder measure $\pi_0 = \mathcal{N}(0, \mathcal{C}_2)$ is not a proper Gaussian measure on H so that $\|y\| = \|\eta\| = \infty$ a.s.; as a consequence, both risks are infinite a.s.. No generality is lost by directly working in this setting, as pre-whitening yields the transformed equation $\mathcal{C}_2^{-1/2}y = \mathcal{C}_2^{-1/2}Lx + \xi$, where $\xi \sim \mathcal{N}(0, \text{Id})$ is Gaussian white noise (and $\text{Id} \in \mathcal{L}(H)$ is not trace-class on H); the learning can then proceed as usual with data $(x, \mathcal{C}_2^{-1/2}y)$ and unknown (in general, unbounded) operator $\mathcal{C}_2^{-1/2}L$. In fact, we assume more.

Assumption 2.7. $\mathcal{C}_2 \in \mathcal{L}(\mathcal{D}(L^\dagger))$ and L^\dagger commutes with \mathcal{C}_2 on $\mathcal{D}(L^\dagger)$.

Under Assumption 2.7, $\mathcal{C}_2^{-1/2}L^\dagger x_n = L^\dagger \mathcal{C}_2^{-1/2}x_n$ so that pre-whitening may be applied to the data pairs themselves: $(X, Y) \mapsto \{(\mathcal{C}_2^{-1/2}x_n, \mathcal{C}_2^{-1/2}y_n)\}_{n=1}^N$, while L^\dagger remains unchanged. Henceforth, it suffices to only consider the white noise case, and we impose this explicitly.

Assumption 2.8. *The Gaussian noise covariance operator is white: $\mathcal{C}_2 = \gamma^2 \text{Id}$, $\gamma > 0$.*

The Bayesian formula (2.6) provides inspiration for an alternative definition of risk in infinite dimensions, namely, to redefine it as the negative log likelihood of $\mu_N^{Y,X}$ [8, 45].

Definition 2.9 (Expected Risk). *The expected risk (or population risk) is*

$$(2.10) \quad \mathcal{R}_\infty(L) := \mathbb{E}^{(x,y) \sim \mathcal{P}} \left[\frac{1}{2} \|Lx\|^2 - \langle y, Lx \rangle \right].$$

Definition 2.10 (Empirical Risk). *The empirical risk is*

$$(2.11) \quad \mathcal{R}_N(L) := \frac{1}{N} \sum_{n=1}^N \left[\frac{1}{2} \|Lx_n\|^2 - \langle y_n, Lx_n \rangle \right],$$

and the regularized empirical risk is

$$(2.12) \quad \mathcal{R}_{N,R}(L) := \mathcal{R}_N(L) + \frac{1}{2N} \|W^{-1/2}L\|_{\text{HS}(H_{\mathcal{K}};H)}^2,$$

where $W \in \mathcal{L}(\text{HS}(H_{\mathcal{K}};H))$ is a linear weighting operator and \mathcal{K} is as in [Assumption 2.1](#).

Equations (2.10) and (2.11) are indeed finite: the infinite “constant” $\|y\|^2$ is subtracted out and the inner product contributions are finite a.s. even though $y \notin H$ a.s., as may be deduced from a limiting procedure [55, Thm. 6.14]. The second term in (2.12) acts like a Tikhonov penalty on the least squares functional, but generalized to operators.

Learning an input-output model for \mathcal{P} given finite data requires a suitable hypothesis class of linear operators on H , henceforth denoted by \mathcal{L} , and leads to the optimization problems

$$(2.13) \quad \inf_{L \in \mathcal{L}} \mathcal{R}_N(L), \quad \inf_{L \in \mathcal{L}} \mathcal{R}_{N,R}(L).$$

For now, we take \mathcal{L} to be the separable Hilbert space $\text{HS}(H_{\mathcal{K}};H)$ from [Assumption 2.1](#). The first problem in (2.13) is known as empirical risk minimization [41, 56]; the second is its regularized counterpart, and its minimizer, denoted by $\hat{L}^{(N,R)}$, is the focus of our discussion. These problems are meant to approximate the best linear estimator given infinite data, denoted by \hat{L} , which is a solution to $\inf\{\mathcal{R}_{\infty}(L) : L \text{ linear and measurable}\}$; this agrees with solutions to $\inf_{L \in \mathcal{L}} \mathcal{R}_{\infty}(L)$ since $L^{\dagger} \in \mathcal{L}$. Indeed, it is well known [40, Thm. 2] that

$$(2.14) \quad \hat{L} = \left(\mathbb{E}[y \otimes x] (\mathbb{E}[x \otimes x])^{-1/2} \right) (\mathbb{E}[x \otimes x])^{-1/2}.$$

Relatedly, formal computations reveal that $\hat{L}^{(N,R)}$ is essentially a regularized empirical approximation to the right hand side of (2.14) and additionally may be identified as the posterior mean $\bar{L}^{(N)}$ from (2.4) whenever $W = \mathcal{C}_3$ (the prior covariance). In this case the quadratic penalty in (2.12) is the Cameron–Martin norm of the prior μ_0 , and the *maximum a posteriori* estimator (which equals the posterior mean in our linear Gaussian setting) must minimize the Onsager–Machlup functional (2.12) [21]; henceforth we write $\hat{L}^{(N,R)} = \bar{L}^{(N)}$. To quantify the performance of $\bar{L}^{(N)}$, we employ the following well known error functionals.

Definition 2.11 (Excess Risk). *The excess risk of the posterior mean is*

$$(2.15) \quad \mathcal{R}_{\infty}(\bar{L}^{(N)}) - \mathcal{R}_{\infty}(\hat{L}).$$

Definition 2.12 (Generalization Gap). *The generalization gap of the posterior mean is*

$$(2.16) \quad \mathcal{R}_{\infty}(\bar{L}^{(N)}) - \mathcal{R}_N(\bar{L}^{(N)}).$$

The excess risk is always nonnegative and provides a notion of consistency for the estimator $\bar{L}^{(N)}$, while the generalization gap can take any sign and is defined with respect to the unregularized empirical risk. Subtracting the generalization gap from the excess risk gives $\mathcal{R}_N(\bar{L}^{(N)}) - \mathcal{R}_{\infty}(\hat{L})$ which, if known, enables estimation of the expected risk of the best linear estimator because $\mathcal{R}_N(\bar{L}^{(N)})$ is computable. We now elaborate on [Definitions 2.11](#) and [2.12](#).

2.3.1. Excess Risk. Under [Assumption 2.1](#), it holds that

$$(2.17) \quad \mathcal{R}_\infty(L) = \frac{1}{2} \mathbb{E}^{x \sim \nu} \|L^\dagger x - Lx\|^2 - \frac{1}{2} \mathbb{E}^{x \sim \nu} \|L^\dagger x\|^2.$$

This implies $\inf_L \mathcal{R}_\infty(L) = \mathcal{R}_\infty(L^\dagger) = -\frac{1}{2} \|L^\dagger\|_{L^2_\nu(H;H)}^2 =: b$ since $\inf_L \mathcal{R}_\infty(L) \geq b$ and $\inf_L \mathcal{R}_\infty(L) \leq \mathcal{R}_\infty(L^\dagger) = b$. Equations (2.1) and (2.14) also imply that $\hat{L} = L^\dagger = \mathbb{E}[y|x = \cdot]$ is the so-called *regression function* [17]. Hence, the excess risk has an explicit formula.

Proposition 2.13. *The excess risk of the posterior mean estimator is $\frac{1}{2} \|L^\dagger - \bar{L}^{(N)}\|_{L^2_\nu(H;H)}^2$.*

Our main results in [Section 3](#) control the excess risk either in expectation or with high probability over the input training samples X . We denote these quantities by

$$(2.18a) \quad \mathcal{E}_N^\mathbb{E} := \mathbb{E} \|L^\dagger - \bar{L}^{(N)}\|_{L^2_\nu(H;H)}^2,$$

$$(2.18b) \quad \mathcal{E}_N^\mathbb{P} := \mathbb{E}^{E \sim \pi_0^{\otimes N}} \|L^\dagger - \bar{L}^{(N)}\|_{L^2_\nu(H;H)}^2,$$

respectively, neglecting the factor 1/2 from [Proposition 2.13](#) to ease notation.

2.3.2. Generalization Gap. Similarly, the generalization gap (2.16) may be written in terms of L^\dagger instead of y (see (A.8) in the proof of [Theorem 3.8](#)). In [Section 3](#), we focus on obtaining explicit bounds for the *expected generalization gap*, defined by

$$(2.19) \quad \mathcal{G}_N^\mathbb{E} := \mathbb{E} |\mathcal{R}_\infty(\bar{L}^{(N)}) - \mathcal{R}_N(\bar{L}^{(N)})|.$$

2.4. Theorem Setting. We now describe the mathematical setting under which our theorems are developed; this setting is not as general as that outlined in the previous two subsections, but it is more amenable to analysis. To this end, let $\{\varphi_j\}_{j \in \mathbb{N}}$ be an orthonormal basis of H . We specialize our hypothesis class \mathcal{L} of linear operators from [Subsection 2.3](#) to the set

$$(2.20) \quad \mathcal{L}_s := \left\{ L \in \text{HS}(H_{\mathcal{C}(s)}; H) \mid L = \sum_{j=1}^{\infty} \ell_j \varphi_j \otimes \varphi_j, \ell = \{\ell_j\}_{j \in \mathbb{N}} \in \mathbb{R}^\infty \right\},$$

where $\mathcal{K} := \mathcal{C}(s)$ defines the input Hilbert space for L^\dagger and has eigenvectors $\{\varphi_j\}_{j \in \mathbb{N}}$ and corresponding eigenvalues $\{c_j^2\}_{j \in \mathbb{N}}$ with property $c_j^2 \asymp j^{2s}$ for some $s < -1/2$; thus $\mathcal{C}(s)$ is trace-class. The series representation of $L \in \mathcal{L}_s$ converges in $\text{HS}(H_{\mathcal{C}(s)}; H)$, which is equivalent to convergence of eigenvalue sequence ℓ in \mathcal{H}^s : $\|L\mathcal{C}(s)^{1/2}\|_{\text{HS}} \asymp \|\ell\|_{\mathcal{H}^s}$. Although taking $\mathcal{K} = \mathcal{C}(s)$ with $s \geq -1/2$ violates the imposed trace-class assumption, the norm $\|\cdot\|_{\mathcal{H}^s}$ on ℓ is still well-defined, and hence we abuse notation and allow any $s \in \mathbb{R}$ to determine \mathcal{L}_s .

In (2.20), we view any $L \in \mathcal{L}_s$ as a densely defined operator $L: \mathcal{D}(L) \subseteq H \rightarrow H$ with

$$(2.21) \quad \mathcal{D}(L) := \left\{ h \in H \mid \|Lh\|^2 = \sum_{j=1}^{\infty} \ell_j^2 \langle \varphi_j, h \rangle^2 < \infty \right\}.$$

It is straightforward to check that any such L equipped with domain (2.21) is self-adjoint on H . In words, \mathcal{L}_s consists of densely defined, self-adjoint linear operators diagonal in the basis $\{\varphi_j\}$ with controlled growth or decay of their eigenvalues $\ell = \{\ell_j\} \in \mathcal{H}^s$, $s \in \mathbb{R}$. In particular, \mathcal{L} contains unbounded operators, for example, any operator with $\ell \in \mathcal{H}^s$ such that $|\ell_j| \rightarrow \infty$. However, self-adjoint linear operators with nonempty *continuous spectra* fall outside the scope of our framework and require more advanced techniques to handle (see, e.g., [19]).

2.4.1. Main Assumptions. Recall that $\nu = \mathcal{N}(0, \mathcal{C}_1)$ is the input training data probability distribution, $\pi_0 = \mathcal{N}(0, \gamma^2 \text{Id})$ is the observational noise distribution, and $\mu_0 = \mathcal{N}(0, \mathcal{C}_3)$ is the prior. As in [Definition 2.5](#), ν' is an arbitrary test data distribution which we restrict to be the Gaussian measure $\mathcal{N}(0, \mathcal{C}_0)$ for simplicity. We impose the following main assumptions:

Assumption 2.14. (i) [Assumption 2.1](#), [2.7](#), and [2.8](#) hold; (ii) the truth is an element of the hypothesis class, that is, $L^\dagger := \sum_{j=1}^\infty \ell_j^\dagger \varphi_j \otimes \varphi_j \in \mathcal{L}_s$; (iii) \mathcal{C}_1 and \mathcal{C}_0 are simultaneously diagonalizable in the same orthonormal basis of eigenvectors, denoted by $\{\phi_k\}_{k \in \mathbb{N}}$, with eigenvalue sequences $\{\lambda_k^2\}_{k \in \mathbb{N}}$ and $\{\lambda_{0k}^2\}_{k \in \mathbb{N}}$, respectively; (iv) as a mapping on operators, \mathcal{C}_3 is “doubly diagonal” in the same basis as L^\dagger , in the sense that $\mathcal{C}_3 \varphi_i \otimes \varphi_j = c_j^2 \sigma_j^2 \delta_{ij} \varphi_i \otimes \varphi_j$; (v) there exist $\alpha_0, \alpha_1 > 1/2$ and $\alpha_3, \beta \in \mathbb{R}$ such that $\sigma_j^2 \asymp j^{-2\alpha_3}$, $\{\lambda_k^2\}_{k \in \mathbb{N}}$ and $\{\lambda_{0k}^2\}_{k \in \mathbb{N}}$ are such that $w_j^2 := \sum_{k=1}^\infty \lambda_k^2 \langle \varphi_j, \phi_k \rangle^2 \asymp j^{-2\alpha_1}$ and $\vartheta_j^2 := \sum_{k=1}^\infty \lambda_{0k}^2 \langle \varphi_j, \phi_k \rangle^2 \asymp j^{-2\alpha_0}$, and $\ell^\dagger := \{\ell_j^\dagger\}_{j \in \mathbb{N}} \in \mathcal{H}^{s'}$ for any $s' < \beta$; and finally, (vi) $\beta > s > -(\alpha_0 \wedge \alpha_1)$ and (vii) $(\alpha_0 \wedge \alpha_1) + \alpha_3 > 1/2$.

Remark 2.15 (Interpretation of Assumption 2.14). Since $L^\dagger \in \mathcal{L}_s$ by (ii), the class \mathcal{L}_s does not introduce any approximation error, and hence we need only control the estimation error due to finite data; this is the *well-specified* setting. Condition (iii) is provided for convenience, while (iv) guarantees that samples $L \sim \mu_0$ from the prior are actually diagonal in the basis $\{\varphi_j\}$, and KL expansion shows that $\ell_j \sim \mathcal{N}(0, \sigma_j^2)$. The algebraic growth or decay of the sequences in (v) closely align our setting with that in [\[35\]](#), although natural extensions to exponential rates is also possible, as in [\[7\]](#). Additionally, the parameter β is introduced as the “largest” s' such that $\ell^\dagger \in \mathcal{H}^{s'}$ (and such β is generally unknown *a priori*), but this lacks precision at the boundary $s' = \beta$ for certain sequences, as noted in [\[35\]](#). One remedy would be to consider slowly varying functions as in that paper, but we do not pursue this. Condition (vi) restates part of [Assumption 2.1](#) in diagonal form, that is, ensures $\nu'(\mathcal{D}(L^\dagger)) = \nu(\mathcal{D}(L^\dagger)) = 1$, while (vii) ensures that $\mathcal{C}_3, \mathcal{C}^{(N)}$ are trace-class on both $\text{HS}(H_{\mathcal{C}_0}; H)$ and $\text{HS}(H_{\mathcal{C}_1}; H)$ a.s. \diamond

Remark 2.16. The prior and test measure “smoothness”, α_3 and α_0 , are viewed as tunable while α_1 and β are fixed by the data. Furthermore, by the Feldman–Hájek theorem [\[22, Thm. 37\]](#), test measure ν' is singular with respect to training measure ν whenever $\alpha_0 \neq \alpha_1$. \diamond

2.4.2. Diagonalization. Under [Assumption 2.14](#), we may work entirely in coordinates with respect to the orthonormal bases $\{\varphi_j\}_{j \in \mathbb{N}}$ and $\{\phi_k\}_{k \in \mathbb{N}}$, leading to a prior and posterior on sequences. Many papers have studied statistical inverse problems in a similar sequence model form and under related commuting assumptions [\[5, 7, 18, 35\]](#). However, our Gaussian diagonal setting is different as not only do $K_X^* K_X, \mathcal{C}_3$ commute, but so do L^\dagger, \mathcal{C}_2 (\mathcal{C}_2 whitened in [Assumption 2.8](#)) and $\mathcal{C}_1, \mathcal{C}_0$; there is no analog of this additional layer of diagonalizability in the previous works due to the unique linear operator structure of L^\dagger here. Also, the complete theory in [\[35\]](#) is not easy to apply in our setting since it requires an eigenbasis for $K_X^* K_X$.

We write the prior μ_0 in sequence form, identifying $L \sim \mu_0$ with its eigenvalues $\ell \sim \mu_{0, \text{seq}} := \bigotimes_{j=1}^\infty \mathcal{N}(0, \sigma_j^2)$. The subscript “seq” is used here, and in what follows below, to denote the sequence space representation of a measure that can be diagonalized. If L^\dagger and \mathcal{C}_1 commute (and hence L^\dagger with \mathcal{C}_0 also), then without loss of generality $\varphi_j = \phi_j$ for all $j \in \mathbb{N}$ so that the action of L^\dagger on $x \sim \nu$ in coordinates is given by $\langle \varphi_j, L^\dagger x \rangle = \ell_j^\dagger \langle \varphi_j, x \rangle \sim \ell_j^\dagger \mathcal{N}(0, \lambda_j^2)$.

However, we allow for the more general case in which L^\dagger and \mathcal{C}_1 do not commute, so that

$$(2.22) \quad y_{jn} = \sum_{k=1}^{\infty} \mathbb{L}_{jk}^\dagger x_{kn} + \gamma \xi_{jn}, \quad \mathbb{L}_{jk}^\dagger = \ell_j^\dagger \langle \varphi_j, \phi_k \rangle, \quad n \in \{1, \dots, N\}, \quad j \in \mathbb{N},$$

is equivalent to the given data (2.2), where $y_{jn} := \langle \varphi_j, y_n \rangle$, $\{x_{jn} := \langle \phi_j, x_n \rangle\}_{n=1}^N$ is an i.i.d. family with $x_{j1} \sim \mathcal{N}(0, \lambda_j^2)$, and $\xi_{jn} \sim \mathcal{N}(0, 1)$ i.i.d.. We further make the definition

$$(2.23) \quad g_{jn} := \langle \varphi_j, x_n \rangle = \sum_{k=1}^{\infty} \langle \varphi_j, \phi_k \rangle x_{kn} \sim \mathcal{N}(0, w_j^2),$$

recalling $\{w_j\}_{j \in \mathbb{N}}$ from Assumption 2.14; hence $y_{jn} = g_{jn} \ell_j^\dagger + \gamma \xi_{jn}$. The primary difficulty that arises when L^\dagger and \mathcal{C}_1 do not commute (i.e., not simultaneously diagonalizable) is that for fixed n , $\{g_{jn}\}_{j \in \mathbb{N}}$ is not an independent family due to the correlation introduced by all the $\{x_{kn}\}_{k \in \mathbb{N}}$; to deal with this, our proofs use some independence-agnostic results. Nonetheless, $\{g_{jn}\}_{n=1}^N$ is still an i.i.d. family for fixed $j \in \mathbb{N}$ and thus $\mathbb{E}[g_{jn} g_{jn'}] = w_j^2 \delta_{nn'}$ for $n, n' \leq N$.

Remark 2.17. We note that, in the absence of the additive noise η , determination of L^\dagger is trivial: it may be recovered from a *single* input-output pair, say $(x_1, L^\dagger x_1)$. This special structure of the operator learning problem is arguably unrealistic in the context of supervised learning, where noiseless data is often assumed in universal approximation theorems, and is a direct consequence of the diagonalization from Assumption 2.14. However, our numerics will demonstrate the relevance of our theory in settings that go beyond this assumption. \diamond

The diagonal structure of \mathcal{L}_s and Assumption 2.14 also provide a useful characterization of the Bochner norm (2.7): for $L \in \mathcal{L}_s$ and trace-class $\mathcal{C} \in \mathcal{L}(H)$ with eigenpairs $\{(\theta_k^2, \psi_k)\}$,

$$(2.24) \quad \mathbb{E}^{x \sim \mathcal{N}(0, \mathcal{C})} \|Lx\|^2 = \sum_{j=1}^{\infty} \sum_{k=1}^{\infty} \theta_k^2 \mathbb{L}_{jk}^2 = \sum_{j=1}^{\infty} \hat{\theta}_j^2 \ell_j^2, \quad \hat{\theta}_j^2 := \sum_{k=1}^{\infty} \theta_k^2 \langle \varphi_j, \psi_k \rangle^2,$$

where $\mathbb{L}_{jk} := \langle \varphi_j, L\psi_k \rangle = \ell_j \langle \varphi_j, \psi_k \rangle$ is a matrix coordinate representation of L .

2.4.3. Posterior Characterization. In diagonal coordinates, we have the Bayesian inverse problem of finding sequence $\{\ell_j | Y, X\}_{j \in \mathbb{N}}$, with $\ell_j \sim \mathcal{N}(0, \sigma_j^2)$ independently for each j , from

$$(2.25) \quad y_{jn} = g_{jn} \ell_j + \gamma \xi_{jn}, \quad n \in \{1, \dots, N\}, \quad j \in \mathbb{N}.$$

Here the $\{\xi_{jn}\}$ are i.i.d. unit Gaussian and the $\{g_{jn}\}$ are given by (2.23). This is an infinite collection of random, decoupled, scalar inverse problems and is equivalent to the full infinite-dimensional problem (2.3) under our assumptions. Now, for two sequences $\{a_{jn}\}$ and $\{b_{jn}\}$, we henceforth use the averaging notation $\overline{a_j b_j}^{(N)} := \frac{1}{N} \sum_{n=1}^N a_{jn} b_{jn}$. Although the abstract formulae (2.5) hold true, (2.25) is easier to manipulate and may be solved explicitly. Indeed, by completing the square [55, Ex. 6.23] we obtain the following solution.

Proposition 2.18. The posterior $\mu_{N, \text{seq}}^{Y, X}$ on eigenvalue sequence $\ell | Y, X = \{\ell_j | Y, X\}_{j \in \mathbb{N}}$ is

$$(2.26) \quad \mu_{N, \text{seq}}^{Y, X} = \bigotimes_{j=1}^{\infty} \mathcal{N}(\bar{\ell}_j^{(N)}, (c_j^{(N)})^2), \quad \bar{\ell}_j^{(N)} = \frac{N \gamma^{-2} \sigma_j^2 \overline{y_j g_j}^{(N)}}{1 + N \gamma^{-2} \sigma_j^2 \overline{g_j g_j}^{(N)}}, \quad (c_j^{(N)})^2 = \frac{\sigma_j^2}{1 + N \gamma^{-2} \sigma_j^2 \overline{g_j g_j}^{(N)}}.$$

Defining the linear bijection $B: \ell \mapsto \sum_{j=1}^{\infty} \ell_j \varphi_j \otimes \varphi_j$, it follows that the actual posterior $\mu_N^{Y,X}$ on $L|Y, X$ is the pushforward of $\mu_{N,\text{seq}}^{Y,X}$ under B , that is, $\mu_N^{Y,X} = B_{\#} \mu_{N,\text{seq}}^{Y,X}$. Hence, $\mu_{N,\text{seq}}^{Y,X}$ is a more convenient solution representation. We are now ready to prove its convergence.

3. Convergence Theory. One common approach to studying the consistency of solutions to Bayesian inverse problems is to analyze the rate of contraction of the posterior measure to a Dirac measure centered on the ground truth as $N \rightarrow \infty$. To make this idea concrete in our setting where the posterior is $\mu_N^{Y,X}$ and the ground truth is L^\dagger , we follow [4, 7, 34, 35] and consider finding a sequence $\varepsilon_N \rightarrow 0$ such that for any sequence $M_N \rightarrow \infty$,

$$(3.1) \quad \mathbb{E}^{Y,X} \mu_N^{Y,X} (\{L: \|L - L^\dagger\|_{L_{\nu'}^2(H;H)} \geq M_N \varepsilon_N\}) \rightarrow 0 \quad \text{as } N \rightarrow \infty,$$

where $\mathbb{E}^{Y,X} = \mathbb{E}_{\{x_n, y_n\}_{n=1}^N \sim \mathcal{P}^{\otimes N}} = \mathbb{E}^{E \sim \pi_0^{\otimes N}} \mathbb{E}^{X \sim \nu^{\otimes N}}$, and we say that ε_N is a *contraction rate* of posterior $\mu_N^{Y,X}$ with respect to the $L_{\nu'}^2(H;H)$ Bochner norm. By Chebyshev's inequality,

$$(3.2) \quad \frac{1}{M_N^2 \varepsilon_N^2} \mathbb{E}^{Y,X} \mathbb{E}^{L \sim \mu_N^{Y,X}} \|L - L^\dagger\|_{L_{\nu'}^2(H;H)}^2$$

is an upper bound for the left hand side of (3.1), so that the limit in (3.1) holds true if the *squared posterior contraction* (SPC), which equals $\mathbb{E} \|L - L^\dagger\|_{L_{\nu'}^2(H;H)}^2$, is $O(\varepsilon_N^2)$ as $N \rightarrow \infty$.

In Subsection 3.1, we provide asymptotic convergence rates of both the posterior sample and mean test errors (2.8)–(2.9), as well as related high probability versions (Theorems 3.1 and 3.3); the former implies a posterior contraction rate via (3.1)–(3.2). Corresponding lower bounds are given in Subsection 3.2 (Theorems 3.5 and 3.6). Finally, both upper and lower bounds are established, in expectation, for the excess risk (2.15) and generalization gap (2.16) in Subsection 3.3 (Theorems 3.7 to 3.9). We defer all the proofs to Appendix A.

3.1. Upper Bounds.

Theorem 3.1. *If L^\dagger , $\{\lambda_j\}$, $\{\lambda_{0j}\}$, and $\{\sigma_j\}$ are as in Assumption 2.14 and $L^{(N)} \sim \mu_N^{Y,X}$, then for any $s \in (-(\alpha_0 \wedge \alpha_1), \beta)$, as $N \rightarrow \infty$,*

$$(3.3) \quad \mathbb{E} \|L^\dagger - L^{(N)}\|_{L_{\nu'}^2(H;H)}^2 = o(N^{-\left(\frac{\alpha_0+s}{\alpha_1+\alpha_3}\right)}) + \begin{cases} O(N^{-(1-\frac{\alpha_1+1/2-\alpha_0}{\alpha_1+\alpha_3})}), & \text{if } \alpha_0 < \alpha_1 + 1/2 \\ O(N^{-1} \log N), & \text{if } \alpha_0 = \alpha_1 + 1/2 \\ O(N^{-1}), & \text{if } \alpha_0 > \alpha_1 + 1/2 \end{cases}$$

whenever $\frac{\alpha_0+s}{\alpha_1+\alpha_3} < 2$, and $\mathbb{E} \|L^\dagger - L^{(N)}\|_{L_{\nu'}^2(H;H)}^2$ is of the order of the second term on the right hand side of (3.3) whenever $\frac{\alpha_0+s}{\alpha_1+\alpha_3} \geq 2$. The same results hold under substitution of the posterior mean $\bar{L}^{(N)} \leftarrow L^{(N)}$ into the left hand side of (3.3).

The next corollary follows since the SPC is precisely the posterior sample test error (2.8).

Corollary 3.2. *Any sequence $\{\varepsilon_N\}_{N \in \mathbb{N}}$ such that ε_N^2 is of the order of the right hand side of (3.3) as $N \rightarrow \infty$ is a contraction rate of $\mu_N^{Y,X}$ with respect to the $L_{\nu'}^2(H;H)$ Bochner norm.*

We now produce closely related estimates which hold with high probability over the input data $X \sim \nu^{\otimes N}$. To state the results cleanly we define $N_{\delta-} := (1 - \delta)N$ for any $\delta \in (0, 1)$.

Theorem 3.3 (High Probability Upper Bound). *Let L^\dagger , $\{\lambda_j\}$, $\{\lambda_{0j}\}$, and $\{\sigma_j\}$ be as in [Assumption 2.14](#) and $L^{(N)} \sim \mu_N^{Y,X}$. Fix $s \in (-(\alpha_0 \wedge \alpha_1), \beta)$ and $\delta \in (0, 1)$. Then there exists a constant $C > 0$ (depending only on δ , α_1 , and α_3) and $c \in (0, 1/8)$ such that as $N \rightarrow \infty$, with probability at least $1 - C \exp(-cN\delta^2)$ over input data $X \sim \nu^{\otimes N}$,*

$$(3.4) \quad \mathbb{E}^{\pi_0^{\otimes N}} \mathbb{E}^{\mu_N^{Y,X}} \|L^\dagger - L^{(N)}\|_{L_{\nu'}^2(H;H)}^2 = o(N_{\delta^-}^{-\frac{\alpha_0+s}{\alpha_1+\alpha_3}}) + \begin{cases} O\left(\left(\frac{1+\delta}{1-\delta}\right) N_{\delta^-}^{-\left(1-\frac{\alpha_1+1/2-\alpha_0}{\alpha_1+\alpha_3}\right)}\right), & \alpha_0 < \alpha_1 + \frac{1}{2} \\ O(N_{\delta^-}^{-1} \log N_{\delta^-}), & \alpha_0 = \alpha_1 + \frac{1}{2} \\ O(N_{\delta^-}^{-1}), & \alpha_0 > \alpha_1 + \frac{1}{2} \end{cases}$$

whenever $\frac{\alpha_0+s}{\alpha_1+\alpha_3} < 2$, and $\mathbb{E}^{\pi_0^{\otimes N}} \mathbb{E}^{\mu_N^{Y,X}} \|L^\dagger - L^{(N)}\|_{L_{\nu'}^2(H;H)}^2$ is of the order of the second term on the right hand side of (3.4) whenever $\frac{\alpha_0+s}{\alpha_1+\alpha_3} \geq 2$. The same results hold under substitution of the posterior mean $\bar{L}^{(N)} \leftarrow L^{(N)}$ into the left hand side of (3.4).

Remark 3.4. While the notion of averaging only over the noise realizations $E \sim \pi_0^{\otimes N}$ may seem unnatural, in practice such a situation could arise when the random data generating process is observed simultaneously with different measurement devices: so, there is only a single realization of $X \sim \nu^{\otimes N}$, but the N random measurement errors $E \sim \pi_0^{\otimes N}$ vary. \diamond

3.2. Lower Bounds.

Theorem 3.5. *If L^\dagger , $\{\lambda_j\}$, $\{\lambda_{0j}\}$, and $\{\sigma_j\}$ are as in [Assumption 2.14](#) and $L^{(N)} \sim \mu_N^{Y,X}$, then for any sequence $\tau_N \rightarrow 0$, as $N \rightarrow \infty$,*

$$(3.5) \quad \mathbb{E} \|L^\dagger - \bar{L}^{(N)}\|_{L_{\nu'}^2(H;H)}^2 = \begin{cases} \Omega(\tau_N N^{-\left(1-\frac{\alpha_1+1/2-\alpha_0}{\alpha_1+\alpha_3}\right)}), & \text{if } \alpha_0 < \alpha_1 + 1/2 \\ \Omega(\tau_N N^{-1} \log N), & \text{if } \alpha_0 = \alpha_1 + 1/2 \\ \Omega(\tau_N N^{-1}), & \text{if } \alpha_0 > \alpha_1 + 1/2, \end{cases}$$

and $\mathbb{E} \|L^\dagger - L^{(N)}\|_{L_{\nu'}^2(H;H)}^2$ is of the order of the right hand side of (3.5), but without τ_N .

The lower bounds with high probability over the input training samples are analogous.

Theorem 3.6 (High Probability Lower Bound). *If L^\dagger , $\{\lambda_j\}$, $\{\lambda_{0j}\}$, and $\{\sigma_j\}$ are as in [Assumption 2.14](#) and $L^{(N)} \sim \mu_N^{Y,X}$, then for any $\delta \in (0, 1)$, there exists a constant $C > 0$ (depending only on δ , α_1 , and α_3) and $c \in (0, 1/8)$ such that as $N \rightarrow \infty$, with probability at least $1 - C \exp(-cN\delta^2)$ over input data $X \sim \nu^{\otimes N}$,*

$$(3.6) \quad \mathbb{E}^{\pi_0^{\otimes N}} \mathbb{E}^{\mu_N^{Y,X}} \|L^\dagger - L^{(N)}\|_{L_{\nu'}^2(H;H)}^2 = \begin{cases} \Omega(N^{-\left(1-\frac{\alpha_1+1/2-\alpha_0}{\alpha_1+\alpha_3}\right)}), & \text{if } \alpha_0 < \alpha_1 + 1/2 \\ \Omega(N^{-1} \log N), & \text{if } \alpha_0 = \alpha_1 + 1/2 \\ \Omega(N^{-1}), & \text{if } \alpha_0 > \alpha_1 + 1/2. \end{cases}$$

The same result holds under substitution of the posterior mean $\bar{L}^{(N)} \leftarrow L^{(N)}$ into the left hand side of the display (3.6), but multiplied by $1 - \delta$ on the right hand side.

3.3. Bounds for Statistical Learning Functionals. In the previous two subsections, we bounded the SPC from above and below. It follows from [Subsection 2.3.1](#) that corresponding bounds for the excess risks [\(2.18a\)](#) and [\(2.18b\)](#) may be obtained by specializing to the case $\alpha_0 = \alpha_1$ and using only the posterior mean estimator. Concretely, we have the following theorem, proved as a consequence of [Theorems 3.1](#) and [3.5](#).

Theorem 3.7 (Upper and Lower Bounds on Excess Risk). *If L^\dagger , $\{\lambda_j\}$, and $\{\sigma_j\}$ are as in [Assumption 2.14](#), then for any $s \in (-\alpha_1, \beta)$, as $N \rightarrow \infty$,*

$$(3.7) \quad \mathcal{E}_N^{\mathbb{E}} = O\left(N^{-\left(\frac{\alpha_1+\alpha_3-1/2}{\alpha_1+\alpha_3}\right)}\right) + \begin{cases} o\left(N^{-\left(\frac{\alpha_1+s}{\alpha_1+\alpha_3}\right)}\right), & \text{if } \frac{\alpha_1+s}{\alpha_1+\alpha_3} < 2 \\ 0, & \text{if } \frac{\alpha_1+s}{\alpha_1+\alpha_3} \geq 2, \end{cases}$$

and for any sequence $\tau_N \rightarrow 0$, as $N \rightarrow \infty$,

$$(3.8) \quad \mathcal{E}_N^{\mathbb{E}} = \Omega\left(\tau_N N^{-\left(\frac{\alpha_1+\alpha_3-1/2}{\alpha_1+\alpha_3}\right)}\right).$$

A similar result may be established for $\mathcal{E}_N^{\mathbb{P}}$, defined in [\(2.18b\)](#), by using [Theorems 3.3](#) and [3.6](#); we omit the details for brevity. It remains to estimate the generalization gap [\(2.19\)](#), and we now establish rates of convergence for it in $L^1(\Omega, \mathcal{F}, \mathbb{P})$.

Theorem 3.8 (Upper Bound on Generalization Gap). *If L^\dagger , $\{\lambda_j\}$, and $\{\sigma_j\}$ are as in [Assumption 2.14](#), then as $N \rightarrow \infty$,*

$$(3.9) \quad \mathcal{G}_N^{\mathbb{E}} = O\left(N^{-\left(\frac{1}{2} \wedge \frac{\alpha_1+\alpha_3-1/2}{\alpha_1+\alpha_3}\right)}\right).$$

Thus the generalization gap decays at least as fast as the typical Monte Carlo rate $N^{-1/2}$ if $\alpha_1 + \alpha_3 \geq 1$; otherwise it decays at a slower rate which is arbitrarily slow as $\alpha_1 + \alpha_3$ approaches $1/2$ from above. The lower bound, shown next, only includes the latter contribution.

Theorem 3.9 (Lower Bound on Generalization Gap). *If L^\dagger , $\{\lambda_j\}$, and $\{\sigma_j\}$ are as in [Assumption 2.14](#), then for any sequence $\tau_N \rightarrow 0$ satisfying $N^{-1/2} = o(\tau_N)$, as $N \rightarrow \infty$,*

$$(3.10) \quad \mathcal{G}_N^{\mathbb{E}} = \Omega\left(\tau_N N^{-\left(\frac{\alpha_1+\alpha_3-1/2}{\alpha_1+\alpha_3}\right)}\right)$$

whenever $\frac{\alpha_1+\beta}{\alpha_1+\alpha_3} > 2$. When $\frac{\alpha_1+\beta}{\alpha_1+\alpha_3} \leq 2$, [\(3.10\)](#) still holds provided $\alpha_3 < 1 + \alpha_1 + 2\beta$ and $\max\{N^{-\frac{1}{2}}, N^{-(\frac{1}{2}-\varepsilon)-(\frac{1}{2}+\beta-\alpha_3)/(\alpha_1+\alpha_3)}\} = o(\tau_N)$ as $N \rightarrow \infty$ for any $\varepsilon > 0$ sufficiently small.

4. Discussion. Since $\mu_N^{Y,X}$ is Gaussian (with a conjugate prior μ_0), the SPC admits a standard decomposition into three terms (see [\(A.1\)–\(A.3\)](#)): the squared estimation bias, estimation variance, and posterior spread (i.e., the trace of $\mathcal{C}^{(N)}$) [[5](#), Sec. 1.1]. Inspection of the proof of [Theorem 3.1](#) shows that the first term on the right of [\(3.3\)](#) is the contribution from the squared estimation bias, while the second is from both the estimation variance and posterior spread (both being of equal order); interpretations are similar for the remaining theorems. Furthermore, we may conclude that our operator learning inverse problem under [Assumption 2.14](#) is *moderately ill-posed* [[5](#), Sec. 4] since ε_N follows a power law.

Now note that the first term in (3.3) increases with α_3 while the second decreases (in the first case $\alpha_0 < \alpha_1 + 1/2$). This suggests the choice of prior smoothness $\alpha_3 = s + 1/2$ (i.e., exactly matching the regularity of the truth ℓ^\dagger) to balance the contributions from both terms. The rate then becomes $N^{-(2\alpha_0+2s)/(1+2\alpha_1+2s)}$ when $\alpha_0 < \alpha_1 + 1/2$ (which for $\alpha_0 = \alpha_1$ is minimax optimal in the sense described in [34, 35], [18, Table 1]) or N^{-1} (up to log terms) when $\alpha_0 \geq \alpha_1 + 1/2$. In the first case, it becomes clear that as s decreases (meaning L^\dagger becomes “less compact” and possibly even unbounded), the rates degrade. One could further analyze the minimax optimality by scaling the prior sequence $\{\sigma_j^2\}$ by an N -dependent sequence (precisely as was done in [4, 35] as well as in deterministic regularization theory, where this is standard practice). Furthermore, our convergence rates are not sharp in the over-smoothed prior regime $\alpha_3 > s + 1/2$ (but are sharp when $\alpha_3 \leq s + 1/2$, see Theorem 3.5), in which the term $N^{-(\alpha_0+s)/(\alpha_1+\alpha_3)}$ dominates the upper bound in (3.3) (the rate decaying to zero as $\alpha_3 \rightarrow \infty$) and is not captured by the lower bound (3.5). Nevertheless, our numerical results in the next section (Figure 2a) strongly suggest that the upper bound is in fact sharp. The rates should also be compared to those obtained from the *direct estimation problem* (i.e., white noise model), see [34, Sec. 3.2] and [6, Sec. 5.1–5.3]. The similarity in the rates arises from our use of weak Bochner (*prediction*) norms for recovery and the fact that K_X is not compact.

In our infinite-dimensional setting, it is also interesting to see from Theorem 3.7 that so-called “fast rates” for the excess risk (i.e., faster than $N^{-1/2}$, see [41]) may be attained by the posterior operator estimator in certain parameter regimes. The usual statistical learning theory techniques based on bounding suprema of empirical processes typically yield “slow” $N^{-1/2}$ rates or worse (see, e.g., [56]). Our results are sharper here since we exploit explicit calculations under our diagonalization assumptions. The training data smoothness α_1 also plays a crucial role in the rates, namely, the in-distribution test error (3.7) tends to the best possible parametric (Monte Carlo) rate N^{-1} as $\alpha_1 \rightarrow \infty$. Further, as α_0 in ν' increases, the rate in the $L_{\nu'}^2(H; H)$ norm always improves, and the three smoothness cases in (3.3) are similar to those in functional linear regression [16]. Finally, Figure 1 visualizes our convergence rates in various settings and suggests the increased difficulty of learning unbounded operators.

We conclude this section with a discussion on the direct learning of inverse maps arising from ill-posed inverse problems. This is currently a popular research area, catalyzed by the success of deep neural networks deployed in these and related problems [8, 9, 10, 15, 23, 25]. If one has access to data of the form $y = Lx + \eta$ for an unknown compact forward operator L and of the form $y' = L^{-1}x' + \eta'$ for the corresponding unknown unbounded inverse operator L^{-1} , then our theory suggests (in the diagonal setting) that the learning of L in the first case enjoys better sample complexity than that of learning L^{-1} in the second (Figure 1c); although less common, data of the latter form could arise from, e.g., noisy differentiation of time series in system identification for time-dependent PDEs. However, our theory does not account for “errors-in-covariates” that distinguishes true inverse map learning, where L^{-1} must be estimated only from noisy forward map samples $y = Lx + \eta$, because the statistical structure does not fit into Assumption 2.14. This may be handled finite-dimensionally with total least squares [28], and the infinite-dimensional setting was considered in [12] but with non-Bayesian methods. We leave inverse operator learning in this challenging setting to future research.

5. Numerical Experiments. In this section, we instantiate our operator learning framework numerically, both according to the theory (Subsection 5.1) and beyond (Subsection 5.2). For clarity, we only implement the posterior mean $\bar{L}^{(N)}$ (via $\bar{\ell}^{(N)}$ in (2.26)). Although realistically the noise variance γ^2 must be estimated from the data, here we simply assume it is known and equal to the value prescribed when generating datasets. Moreover, our conceptually infinite-dimensional formulation in the previous sections must be discretized carefully when implemented on a computer to ensure that its effect does not obscure the theoretical infinite-dimensional behavior [3, Sec. 1.2]. We choose spectral truncation as the discretization method [3, 6] as this naturally finite-dimensionalizes the infinite sequence setting from Subsection 2.4; indeed, for $v = \{v_j\}_{j \in \mathbb{N}} \in \mathbb{R}^\infty$ (the coordinates of a vector or operator in the orthonormal basis $\{\varphi_j\}_{j \in \mathbb{N}}$), its spectral truncation is $v^{(J)} := \{v_j\}_{j \leq J} \in \mathbb{R}^J$ for some $J \in \mathbb{N}$. Lastly, we use the *relative* expected squared $L^2_{\nu'}$ Bochner norm as a numerical error metric:

$$(5.1) \quad \frac{\mathbb{E} \|L^\dagger - \bar{L}^{(N)}\|_{L^2_{\nu'}(H;H)}^2}{\|L^\dagger\|_{L^2_{\nu'}(H;H)}^2} = \frac{\mathbb{E}^{Y,X} \mathbb{E}^{x' \sim \nu'} \|L^\dagger x' - \bar{L}^{(N)} x'\|^2}{\mathbb{E}^{x' \sim \nu'} \|L^\dagger x'\|^2} = \frac{\mathbb{E} \sum_{j=1}^\infty \vartheta_j^2 |\ell_j^\dagger - \bar{\ell}_j^{(N)}|^2}{\sum_{j=1}^\infty \vartheta_j^2 |\ell_j^\dagger|^2}.$$

Finally, similar results to those that follow were obtained in the setting of high probability, rather than expectation (5.1), over $X \sim \nu^{\otimes N}$ as suggested by the theory in Section 3.

5.1. Within the Theory. We now perform a set of carefully designed numerical experiments to confirm the theoretical results of the paper. Our running example in this section involves target operators given by powers of the negative Laplacian in one spatial dimension. Specifically, define $A: \mathcal{D}(A) \subset H \rightarrow H$ by $h \mapsto Ah := -\Delta h$ with domain $\mathcal{D}(A) := H_0^1(I; \mathbb{R}) \cap H^2(I; \mathbb{R})$, where $I := (0, 1)$, $H := L^2(I; \mathbb{R})$, and Δ is the Laplacian. We then consider $L^\dagger = A, \text{Id}, A^{-1}$ corresponding to unbounded, bounded, and compact self-adjoint operators on H , respectively. It is well known that A is diagonalized in the orthonormal basis $\{\varphi_j\}_{j \in \mathbb{N}}$ of H given by $z \mapsto \varphi_j(z) = \sqrt{2} \sin(j\pi z)$, with eigenvalues $\{(j\pi)^2\}_{j \in \mathbb{N}}$. Defining \mathcal{L}_s in (2.20) with respect to $\{\varphi_j\}_j$, it holds that $L^\dagger = A, \text{Id}, A^{-1}$ have eigenvalues $\ell^\dagger = \{(j\pi)^2\}_j, \{1\}_j, \{(j\pi)^{-2}\}_j \in \mathcal{H}^s$ for any $s < \beta$, where $\beta = -5/2, -1/2, 3/2$, respectively.

Turning to the data model, for ν and ν' we choose Matérn-like covariance operators

$$(5.2) \quad \mathcal{C}_i = \tau_i^{2\alpha_i - 1} (A + \tau_i^2 \text{Id})^{-\alpha_i}, \quad i \in \{0, 1\},$$

where τ_i is an inverse length scale and α_i is the smoothness parameter. In this setting, L^\dagger and \mathcal{C}_i are simultaneously diagonalizable in $\{\phi_j = \varphi_j\}$, the eigenvalues of \mathcal{C}_i being $\tau_i^{2\alpha_i - 1} ((j\pi)^2 + \tau_i^2)^{-\alpha_i} \asymp j^{-2\alpha_i}$, $j \in \mathbb{N}$, so that $\lambda_j^2 = w_j^2$ and $\lambda_{0j}^2 = \vartheta_j^2$ satisfy Assumption 2.14 (v). Additionally, we directly define the prior covariance \mathcal{C}_3 in sequence space according to Assumption 2.14 (iv), choosing $\sigma_j^2 := \tau_3^{2\alpha_3 - 1} ((j\pi)^2 + \tau_3^2)^{-\alpha_3} \asymp j^{-2\alpha_3}$. The numerical values of $\alpha_0, \alpha_1, \alpha_3$ are always chosen to satisfy Assumption 2.14 (vi)-(vii).

In what follows we work entirely in Fourier (i.e., coordinate) space, and for each of the three target operators, we generate synthetic datasets according to (2.22) and compute the errors (5.1) and associated experimental convergence rates, both for in-distribution ($\nu' = \nu$) and out-of-distribution ($\nu' \neq \nu$) cases. We fix $\gamma = 10^{-1}, 10^{-3}, 10^{-5}$ for $L^\dagger = A, \text{Id}, A^{-1}$, respectively, in $\pi_0 = \mathcal{N}(0, \gamma^2 \text{Id})$ unless otherwise noted. In particular, each change in sample size $N \in \mathbb{N}$ corresponds to a fresh random dataset $\{(x_n, y_n)\}_{n=1}^N \sim \mathcal{P}^{\otimes N}$ used to construct

$\bar{L}^{(N)}$, and this process is repeated 250, 500, or 1000 times for $L^\dagger = A, \text{Id}, A^{-1}$, respectively, to approximate the outer expectation in (5.1) by sample averages. All estimates of convergence rates are produced by linear least square fits to the logarithm of the experimental error data.

5.1.1. In-Distribution. Here we take $\alpha_0 = \alpha_1 = 4.5$, $\tau_1 = 15$ and define the prior smoothness to depend on L^\dagger , that is, $\alpha_3 = \alpha_3(L^\dagger) = 1/2 + \beta(L^\dagger) + z$, where $z = -0.75, 0, 0.75$ is a fixed shift to replicate rough, matching, or smooth priors, respectively, and $\tau_3 = 1$ is fixed. We numerically study the in-distribution error ((5.1) with $\nu' = \nu$) for the posterior mean estimator, which is equivalent to the expected excess risk (2.18a) up to scaling. Function samples are discretized by keeping up to $J = 2^{16} = 65536$ Fourier modes, and the sample size is $N \in \{2^4, 2^5, \dots, 2^{14}\}$. Table 1 suggests that our theoretical predictions (Theorem 3.7) for the convergence rates are correct and that the asymptotic upper bound (3.7) is in fact sharp.

Table 1

(Matching test measure) Theoretical (upper bounds) v.s. experimental (in parentheses) rates of convergence r in $O(N^{-r})$ of the expected relative squared $L_\nu^2(H; H)$ in-distribution error (i.e., the scaled excess risk $\mathcal{E}_N^\mathbb{E}$).

L^\dagger	[Operator Class]	Rough Prior	Matching Prior	Smooth Prior
A	[Unbounded]	0.714 (0.714)	0.800 (0.809)	0.615 (0.616)
Id	[Bounded]	0.867 (0.865)	0.889 (0.889)	0.762 (0.762)
A^{-1}	[Compact]	0.913 (0.913)	0.923 (0.920)	0.828 (0.830)

Moving on to study the rates of convergence of $\mathcal{E}_N^\mathbb{E}$ and $\mathcal{G}_N^\mathbb{E}$ for the unbounded operator $L^\dagger = A$, we work in the same setting considered previously except now we vary the prior regularity shift z and discretize all functions (and hence the series in (5.1)) with N -dependent spectral truncation; that is, for each N , we only use Fourier modes in the set $\{j \in \mathbb{N} : j \leq cN^{1/u}\}$ for $c > 0$ a tunable constant and $u := 2(\alpha_1 + \alpha_3)$. This approach more easily yields the asymptotic convergence rate of $\mathcal{G}_N^\mathbb{E}$ as predicted by Theorem 3.8 and is justified by careful examination of the proofs in Appendix A; indeed, all convergence rates in Section 3 are maintained with this particular choice of N -dependent spectral truncation since contributions from the tail set $\{j \in \mathbb{N} : j > cN^{1/u}\}$ are of equal order or negligible, asymptotically.

Figure 2 shows the results using N up to 2^{21} samples and c chosen such that $c(2^{21})^{1/u} \approx 2^{14}$ (the maximum spectral truncation level). The finite-dimensionality of the numerics clearly manifests via the influence of the noise variance γ^2 ; for $\mathcal{E}_N^\mathbb{E}$, the rough prior regions ($z < 0$) are relatively insensitive to γ and closely match upper bound (3.7), but in the over-smoothed region $z > 0$ for large z , the rates begin to diverge from the theory (with larger γ yielding larger discrepancy). This is likely due to large γ producing large constants in the error estimates which mask the correct asymptotic rates for the finite N regime considered here.

Similarly, for $\mathcal{G}_N^\mathbb{E}$, the noise level controls the size of constants pre-multiplying the competing terms in the error estimates. For small γ , Figure 2e shows that terms $O(N^{-1/2})$ have the largest constants for finite N (missing the asymptotically smaller contributions for small z) while Figure 2h suggests that large γ ensures the terms $o(N^{-1/2})$ have larger constants for small z (tightly matching the theoretical lower bound curve); the transition to terms $O(N^{-1/2})$ occurs numerically between $-0.8 \leq z \leq -0.5$ (when theoretically it should occur exactly at $z = -1.5$), and this transition point moves to the right when γ is increased.

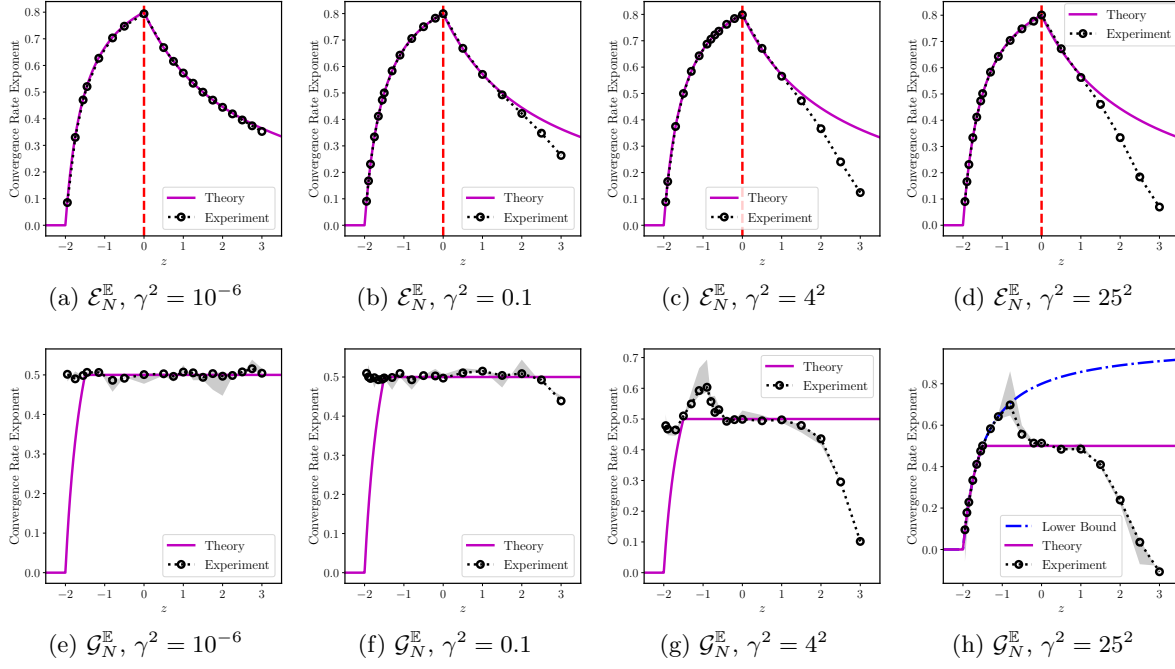


Figure 2. The numerical influence of data noise level for $L^\dagger = A$: **Figures 2a to 2d** show convergence rate exponents for \mathcal{E}_N^E v.s. z for various noise levels γ , with z being the prior smoothness shift parameter, while **Figures 2e to 2h** display rates for \mathcal{G}_N^E v.s. z . Throughout, the solid magenta “Theory” labels denote the theoretical upper bound rates, vertical red dashed lines denote matching priors ($z = 0$), and the shaded regions denote one standard deviation from the mean rate computed from 250 repetitions of the experiment.

5.1.2. Out-of-Distribution. We now vary α_0 , the regularity of test data measure ν' ; we do so in such a way so as to be on either side of the **Theorem 3.1** boundary $\alpha_0 = \alpha_1 + 1/2$. It is clear from the theory (**Theorem 3.1**) that as α_0 increases, that is, the estimator is applied to smoother functions than it was trained on, the convergence rates improve (linearly).

Table 2

(Rougher testing measure) Theoretical (upper bounds) v.s. experimental (in parentheses) rates of convergence r in $O(N^{-r})$ of the out-of-distribution error (5.1): Here, ν' is such that $\alpha_0 = 4 < \alpha_1 = 4.5$.

L^\dagger	[Operator Class]	Rough Prior	Matching Prior	Smooth Prior
A	[Unbounded]	0.429 (0.428)	0.600 (0.607)	0.462 (0.462)
Id	[Bounded]	0.733 (0.734)	0.778 (0.788)	0.667 (0.667)
A^{-1}	[Compact]	0.826 (0.837)	0.846 (0.861)	0.759 (0.764)

Analogously to **Table 1** and the setup there (in particular, keeping $J = 2^{16}$ fixed and using the same trained estimators $\bar{\ell}^{(N)}$ for fixed N), **Tables 2** and **3** show excellent numerical agreement with the asymptotic upper bounds predicted by **Theorem 3.1** in both out-of-distribution regimes. Focusing in on the matching prior case ($z = 0$), in **Figure 3**, we plot the decay of the error (5.1) as N is increased over several orders of magnitude. The solid magenta lines

Table 3

(Smoother testing measure) Same as Table 2 except that here, ν' is such that $\alpha_0 = 5.25 > \alpha_1 = 4.5$.

L^\dagger [Operator Class]	Rough Prior	Matching Prior	Smooth Prior
A [Unbounded]	1.000 (0.992)	1.000 (0.996)	0.846 (0.849)
Id [Bounded]	1.000 (0.986)	1.000 (0.979)	0.905 (0.905)
A^{-1} [Compact]	1.000 (0.981)	1.000 (0.975)	0.931 (0.926)

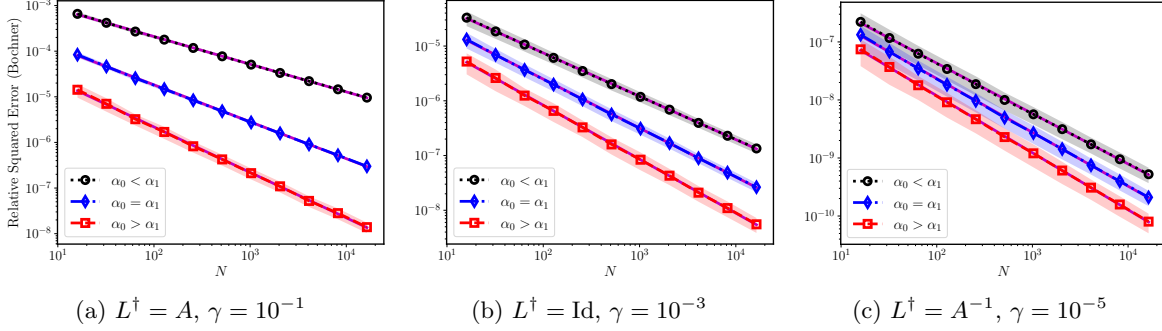


Figure 3. The relative error (5.1) v.s. training sample size N : Here the prior smoothness matches that of the truth, and we use testing measures ν' that are either equal to ($\alpha_0 = \alpha_1$), rougher than ($\alpha_0 < \alpha_1$), or smoother than ($\alpha_0 > \alpha_1$) the training measure ν . For fixed L^\dagger , Figures 3a to 3c show that the same posterior mean estimator achieves smaller error as α_0 increases, that is, when testing against smoother input functions.

represent least square fits and the shaded regions denote one standard deviation from the mean with respect to training set resampling. The excellent numerical fits indicate that in this particular setup, the asymptotic regime described by the theory is rapidly achieved.

5.2. Beyond the Theory. In this subsection, we consider the case where $L^\dagger \in \text{HS}(H_{\mathcal{K}}; H)$ (with \mathcal{K} as in Assumption 2.1) is *not necessarily diagonalized* in orthonormal basis $\{\varphi_j\}$ (here $\{\varphi_j\}$ is assumed fixed, arbitrary, and known) so that the infinite matrix coordinate sequence $\{\mathbb{L}_{jk}^\dagger\}_{j,k \geq 1}$ from (2.22) must be estimated from data instead of the eigenvalues ℓ^\dagger , where now $\mathbb{L}_{jk}^\dagger = \langle \varphi_j, L^\dagger \phi_k \rangle$. By Proposition 2.4, $L^\dagger \in \text{HS}(H_{C_1}; H)$ so that the expansion $L^\dagger = \sum_{i,j} (\lambda_j \mathbb{L}_{ij}^\dagger) \varphi_i \otimes_{H_{C_1}} (\lambda_j \phi_j) = \sum_{i,j} \mathbb{L}_{ij}^\dagger \varphi_i \otimes \phi_j$ always exists and is unique; yet, we have no theory for posterior estimators of L^\dagger . To this end, we briefly derive the posterior mean.

Since the data in $\{\varphi_j\}$ -coordinates are given as the first equality in (2.22), the inverse problem for $\mathbb{L}|Y, X$ decouples along rows of $\mathbb{L} = \{\mathbb{L}_{jk}\}$, denoted by $\mathbb{L}_{j\cdot}$, $j \in \mathbb{N}$. By choosing a Gaussian prior $\mathbb{L}_{j\cdot} \sim \mathcal{N}(0, \Sigma_j)$, we induce a prior on \mathbb{L} . For simplicity, we take $\Sigma_j = \text{diag}(\{\sigma_{jk}^2\}_{k \in \mathbb{N}})$ so that $(\mathbb{L}_{j\cdot})_k = \mathbb{L}_{jk} \sim \mathcal{N}(0, \sigma_{jk}^2)$. After writing down the Onsager–Machlup functional and deriving the normal equations, we obtain for $j, k, \ell \in \mathbb{N}$ the posterior mean

$$(5.3) \quad \bar{\mathbb{L}}_{j\cdot}^{(N)} = (A + \frac{\gamma^2}{N} \Sigma_j^{-1})^{-1} \mathbf{b}_j, \quad A_{\ell k} := \frac{1}{N} \sum_{n=1}^N x_{\ell n} x_{kn}, \quad (\mathbf{b}_j)_\ell := \frac{1}{N} \sum_{n=1}^N y_{jn} x_{\ell n}.$$

For the computation, we use the same covariance (5.2) diagonalized in Fourier sine input

basis $\{\phi_j\}$, but now use Volterra cosine output basis $\{\varphi_j\}$, $z \mapsto \varphi_j(z) := \sqrt{2} \cos((j - \frac{1}{2})\pi z)$. Defining *divergence form elliptic operator* $A_a: \mathcal{D}(A_a) \subset H \rightarrow H$ by $h \mapsto A_a h := -\nabla \cdot (a \nabla h)$, with domain $\mathcal{D}(A_a) = \mathcal{D}(A)$ from before and where $z \mapsto a(z) := \exp(-3z)$ is smooth, we learn (via $\bar{\mathbf{L}}^{(N)}$) unbounded, bounded, and compact self-adjoint operators $L^\dagger = A_a, \text{Id}, A_a^{-1}$, respectively. We pick prior sequences $\sigma_{jk}^2 = \sigma_{jk}^2(L^\dagger)$ for each of the three L^\dagger given by

$$(5.4) \quad \sigma_{jk}^2(L^\dagger) := \begin{cases} (jk)^{-(z-2)} \left(\frac{1+(k/j)^2}{1+(j-k)^2} \right)^2, & \text{if } L^\dagger = A_a \\ (jk)^{-z} \left(\frac{k+k/j}{1+j+(j-k)^2} \right)^2, & \text{if } L^\dagger = \text{Id} \\ (jk)^{-(z+2)} \left(\frac{1+j/k}{1+(j-k)^2} \right)^2, & \text{if } L^\dagger = A_a^{-1}. \end{cases}$$

These priors make \mathbf{L} match the exact asymptotic behavior (as $j \rightarrow \infty$, $k \rightarrow \infty$, and $j = k \rightarrow \infty$) of L^\dagger when $z = 0$. The remaining experimental setup is the same as in [Subsection 5.1](#), except with $J = 2^{12}$, N up to 2^{14} , and only 100 Monte Carlo repetitions. Although A_a is not diagonal in $\varphi_j \neq \phi_j$ (hence each L^\dagger is dense) and the posterior mean estimator is now a doubly indexed sequence, our results in [Figure 4](#) support the same conclusions previously asserted.

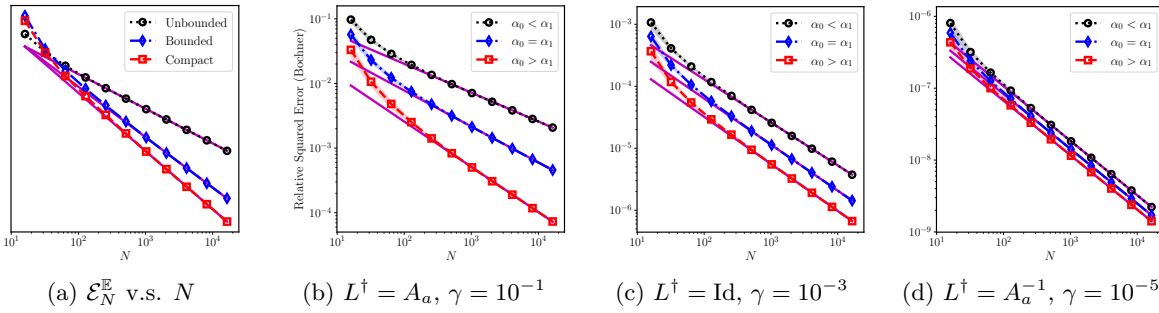


Figure 4. Beyond the theory: Same as [Figures 1a](#) and [3](#) except with non-diagonal elliptic operator A_a .

6. Conclusion. We have theoretically analyzed the nonparametric learning problem of regressing a densely defined linear operator on an infinite-dimensional Hilbert space from infinite-dimensional training data. Though formulated as a linear Bayesian inverse problem, the framework was also interpreted naturally from a statistical learning theory perspective, albeit modified to fit the infinite-dimensional white noise setting considered here. Our theory and numerics both confirmed the intuition that unbounded linear operators are more difficult to learn from noisy data relative to better behaved objects such as bounded or compact operators, in a sense made precise by convergence rates of the test error/squared posterior contraction and excess risk in the large data limit.

We leave open various directions for future work. As we assumed that the target operator was diagonal in a known Hilbert basis, it would be natural to generalize our convergence rate theorems to the non-diagonal target operator setting. We also assumed full knowledge of the input-output bases, but in practice, one or both may be unknown and would have to be estimated from data, for example, by using (functional) principle components analysis (PCA) [[11](#), [20](#), [30](#)]. How the convergence rates would be affected by this additional PCA

approximation, and whether the rates improve in the small noise setting $\gamma = \gamma(N) \rightarrow 0$ as $N \rightarrow \infty$, are also interesting research questions. Finally, it would be of interest to study analogous problems for *nonlinear operators* between Hilbert spaces. Pointwise evaluation remains a linear operation on operators and hence the underlying inverse problem remains linear. However, simplifications such as diagonalization, which rendered the analysis in this paper tractable, are no longer available, and hence the methods of proof will be substantially different; in particular, the architecture used to parametrize unknown operators (e.g., neural operators, random features, and other such methods) will be central to the analysis.

Appendix A. Proofs of Main Results.

Proof of Theorem 3.1. We use (2.24) and (2.26) to deduce that the squared Bochner error, averaged over the posterior and over realizations of the white noise E in the data, but with design points X fixed, satisfies $\mathbb{E}^{\pi_0^{\otimes N}} \mathbb{E}^{\mu_N^{Y,X}} \|L^\dagger - L^{(N)}\|_{L_{\nu'}^2(H;H)}^2 = \mathcal{I}_1 + \mathcal{I}_2 + \mathcal{I}_3$, where

$$(A.1) \quad \mathcal{I}_1 = \sum_{j=1}^{\infty} \frac{\vartheta_j^2 (\ell_j^\dagger)^2}{(1 + N\gamma^{-2} \sigma_j^2 \overline{g_j g_j^{(N)}})^2},$$

$$(A.2) \quad \mathcal{I}_2 = \sum_{j=1}^{\infty} \frac{N \vartheta_j^2 \gamma^{-2} \sigma_j^4 \overline{g_j g_j^{(N)}}}{(1 + N\gamma^{-2} \sigma_j^2 \overline{g_j g_j^{(N)}})^2},$$

$$(A.3) \quad \mathcal{I}_3 = \sum_{j=1}^{\infty} \frac{\vartheta_j^2 \sigma_j^2}{1 + N\gamma^{-2} \sigma_j^2 \overline{g_j g_j^{(N)}}}.$$

Recall that ϑ_j decays as $j^{-\alpha_0}$, reflecting the measure ν' with respect to which the Bochner norm is defined, and σ_j decays as $j^{-\alpha_3}$, reflecting prior assumptions on the operator L . The dependence of the three sums on design X is through the $\{g_{jn}\}$; these centered Gaussian random variables have standard deviation w_j decaying as $j^{-\alpha_1}$. The parameter β enters through the assumed regularity of the truth ℓ^\dagger . Moreover, use of Assumption 2.14 with Lemma B.3 guarantees that (A.1)–(A.3) converge \mathbb{P} -a.s. with respect to the random design X .

In the remainder of the proof we let $u = 2(\alpha_1 + \alpha_3)$ (which is positive by Assumption 2.14) and write $\overline{g_j g_j^{(N)}} = w_j^2 Z_j^{(N)}$, where $N Z_j^{(N)} \sim \chi^2(N)$, for ease of presentation. For intuition it is useful to note that Lemma B.7 quantifies precisely the sense in which $Z_j^{(N)}$ is of order 1 for $N \gg 1$. We split each of the three series (A.1)–(A.3) that compose the Bochner squared error into sums over the two disjoint index sets $\{j \in \mathbb{N} : j \leq N^{1/u}\}$ and $\{j \in \mathbb{N} : j > N^{1/u}\}$ for any $N \in \mathbb{N}$. We denote such sums by \mathcal{I}_i^{\leq} and $\mathcal{I}_i^{>}$, respectively, for each $i \in \{1, 2, 3\}$.

Beginning with $\mathbb{E}\mathcal{I}_2$, we estimate $\mathbb{E}\mathcal{I}_2^{\leq}$ by

$$\mathbb{E} \sum_{j \leq N^{1/u}} \frac{N \vartheta_j^2 \gamma^{-2} \sigma_j^4 \overline{g_j g_j^{(N)}}}{(1 + N\gamma^{-2} \sigma_j^2 \overline{g_j g_j^{(N)}})^2} \leq \sum_{j \leq N^{1/u}} \frac{\vartheta_j^2 \gamma^2 w_j^{-2} \mathbb{E}[(Z_j^{(N)})^{-1}]}{N} \asymp \sum_{j \leq N^{1/u}} \frac{N j^{-2(\alpha_0 - \alpha_1)}}{(N - 2)N},$$

where we used Lemma B.7 to evaluate the negative moment. Since

$$\sum_{j \leq N^{1/u}} \frac{j^{-2(\alpha_0 - \alpha_1)}}{N} = \sum_{j \leq N^{1/u}} \frac{j^{-2(\alpha_0 + \alpha_3)}}{N j^{-u}} \asymp \sum_{j \leq N^{1/u}} \frac{j^{-2(\alpha_0 + \alpha_3)}}{1 + N j^{-u}}$$

(using $Nj^{-u} \asymp 1 + Nj^{-u}$ whenever $j \leq N^{1/u}$) and $N/(N-2) = \Theta(1)$ as $N \rightarrow \infty$, by (B.3b) in Lemma B.2 (applied with $t = 2(\alpha_0 + \alpha_3)$, $v = 1$, and condition $t > 1$ satisfied by $(\alpha_0 \wedge \alpha_1) + \alpha_3 > 1/2$ from Assumption 2.14) we deduce that $\mathbb{E}\mathcal{I}_2^{\leq}$ is of order the right most term in the display (3.3). Upper bounding the denominator of the tail series by one yields

$$\mathbb{E}\mathcal{I}_2^> \leq \sum_{j>N^{1/u}} N\vartheta_j^2 \gamma^{-2} \sigma_j^4 \mathbb{E}[w_j^2 Z_j^{(N)}] \asymp N \sum_{j>N^{1/u}} j^{-2(\alpha_0 + \alpha_1 + 2\alpha_3)} = O(N^{-(1 - \frac{\alpha_1 + 1/2 - \alpha_0}{\alpha_1 + \alpha_3})})$$

as $N \rightarrow \infty$ by Lemma B.7 and (B.3a) in Lemma B.2 (applied with $t = 2(\alpha_0 + \alpha_1 + 2\alpha_3)$, and condition $t > 1$ satisfied by $(\alpha_0 \wedge \alpha_1) + \alpha_3 > 1/2$ from Assumption 2.14), which is always of the same order or negligible compared to the upper bound on $\mathbb{E}\mathcal{I}_2^{\leq}$.

Proceeding with $\mathbb{E}\mathcal{I}_3$, again by Lemma B.7 it holds that

$$\mathbb{E}\mathcal{I}_3^{\leq} = \sum_{j \leq N^{1/u}} \frac{\vartheta_j^2 \sigma_j^2}{1 + N\gamma^{-2} \sigma_j^2 g_j g_j^{(N)}} \leq \sum_{j \leq N^{1/u}} \frac{\vartheta_j^2 \sigma_j^2 \mathbb{E}[(Z_j^{(N)})^{-1}]}{N\gamma^{-2} \sigma_j^2 w_j^2} \asymp \sum_{j \leq N^{1/u}} \frac{Nj^{-2(\alpha_0 - \alpha_1)}}{(N-2)N}$$

which is precisely the same as the asymptotic contribution from $\mathbb{E}\mathcal{I}_2$. The tail is bounded by

$$\mathbb{E}\mathcal{I}_3^> \leq \sum_{j>N^{1/u}} \vartheta_j^2 \sigma_j^2 \asymp \sum_{j>N^{1/u}} j^{-2(\alpha_0 + \alpha_3)} = O(N^{-(1 - \frac{\alpha_1 + 1/2 - \alpha_0}{\alpha_1 + \alpha_3})})$$

as $N \rightarrow \infty$ by (B.3a) in Lemma B.2 (applied with $t = 2(\alpha_0 + \alpha_3)$, and condition $t > 1$ satisfied by $(\alpha_0 \wedge \alpha_1) + \alpha_3 > 1/2$ from Assumption 2.14), just as it was for $\mathbb{E}\mathcal{I}_2^>$.

Finally, observe that $\mathbb{E}\mathcal{I}_1^{\leq}$ admits the upper bound

$$\mathbb{E} \sum_{j \leq N^{\frac{1}{u}}} \frac{\vartheta_j^2 (\ell_j^\dagger)^2}{(1 + N\gamma^{-2} \sigma_j^2 g_j g_j^{(N)})^2} \leq \sum_{j \leq N^{\frac{1}{u}}} \frac{\vartheta_j^2 (\ell_j^\dagger)^2 \mathbb{E}[(Z_j^{(N)})^{-2}]}{(N\gamma^{-2} \sigma_j^2 w_j^2)^2} \asymp \sum_{j \leq N^{\frac{1}{u}}} \frac{j^{-2\alpha_0} (\ell_j^\dagger)^2 \mathbb{E}[(Z_j^{(N)})^{-2}]}{(Nj^{-u})^2}.$$

By Lemma B.7, the expectation in the rightmost sum equals $N^2(N-2)^{-1}(N-4)^{-1} = \Theta(1)$ as $N \rightarrow \infty$. This and the fact that $Nj^{-u} \asymp 1 + Nj^{-u}$ whenever $j \leq N^{1/u}$ implies

$$(A.4) \quad \mathbb{E}\mathcal{I}_1^{\leq} = O\left(\sum_{j \leq N^{1/u}} \frac{j^{-2\alpha_0} (\ell_j^\dagger)^2}{(1 + Nj^{-u})^2}\right) = \begin{cases} o(N^{-(\frac{\alpha_0 + s}{\alpha_1 + \alpha_3})}), & \text{if } \frac{\alpha_0 + s}{\alpha_1 + \alpha_3} < 2 \\ O(N^{-2}), & \text{if } \frac{\alpha_0 + s}{\alpha_1 + \alpha_3} \geq 2 \end{cases}$$

as $N \rightarrow \infty$, where we applied (B.2) in Lemma B.1 (with $t = 2\alpha_0$, $q = s$, $v = 2$, and condition $t \geq -2q$ satisfied by the theorem assumption $s > -(\alpha_0 \wedge \alpha_1)$) to get the last bound. By upper bounding the denominator of the tail series by one, we obtain

$$(A.5) \quad \mathbb{E}\mathcal{I}_1^> \leq \sum_{j>N^{1/u}} \vartheta_j^2 (\ell_j^\dagger)^2 \asymp \sum_{j>N^{1/u}} j^{-2\alpha_0} (\ell_j^\dagger)^2$$

which by (B.1) in Lemma B.1 (applied with $t = 2\alpha_0$, $q = s$) has the same asymptotic upper bound as $\mathbb{E}\mathcal{I}_1^{\leq}$ in the case $\frac{\alpha_0 + s}{\alpha_1 + \alpha_3} < 2$ and is of strictly smaller order otherwise.

All together, we deduce $\mathbb{E}\mathcal{I}_2$ and $\mathbb{E}\mathcal{I}_3$ have the same upper bound, and if $\frac{\alpha_0 + s}{\alpha_1 + \alpha_3} \geq 2$, then the $O(N^{-2})$ contribution from $\mathbb{E}\mathcal{I}_1$ is negligible relative to this bound (yielding the second assertion). Otherwise, we obtain (3.3) as asserted. Since the posterior mean estimator only corresponds to the two error terms \mathcal{I}_1 and \mathcal{I}_2 , the final assertion of the theorem is proved. ■

Proof of Theorem 3.3. Fix $\delta \in (0, 1)$ and define $N_{\delta^-} := (1 - \delta)N$. Following the same conventions as in the proof of Theorem 3.1, except now summing over $\{j \in \mathbb{N} : j \leq N_{\delta^-}^{1/u}\}$ and $\{j \in \mathbb{N} : j > N_{\delta^-}^{1/u}\}$, we begin by estimating \mathcal{I}_1 : using Lemma B.6 (lower tail only),

$$\mathcal{I}_1^{\leq} = \sum_{j \leq N_{\delta^-}^{1/u}} \frac{\vartheta_j^2(\ell_j^\dagger)^2}{(1 + N\gamma^{-2}\sigma_j^2 w_j^2 Z_j^{(N)})^2} \leq \sum_{j \leq N_{\delta^-}^{1/u}} \frac{\vartheta_j^2(\ell_j^\dagger)^2}{(1 + N_{\delta^-}\gamma^{-2}\sigma_j^2 w_j^2)^2} \asymp \sum_{j \leq N_{\delta^-}^{1/u}} \frac{j^{-2\alpha_0}(\ell_j^\dagger)^2}{(1 + N_{\delta^-}j^{-u})^2}$$

with probability at least $1 - N_{\delta^-}^{1/u} \exp(-N\delta^2/8)$. The remaining (almost sure) bounds for \mathcal{I}_1 (including that for the tail $\mathcal{I}_1^>$) are the same as those found in the proof of Theorem 3.1.

Again, similar to the proof of Theorem 3.1, \mathcal{I}_2^{\leq} is upper bounded by

$$\sum_{j \leq N_{\delta^-}^{1/u}} \frac{\vartheta_j^2 \gamma^2}{N w_j^2 Z_j^{(N)}} \leq \sum_{j \leq N_{\delta^-}^{1/u}} \frac{\vartheta_j^2 \gamma^2}{N(1 - \delta) w_j^2} \asymp \sum_{j \leq N_{\delta^-}^{1/u}} \frac{j^{-2(\alpha_0 + \alpha_3)}}{1 + N_{\delta^-} j^{-u}},$$

where the first inequality follows from the same application of Lemma B.6 above; the right hand side is of order the second term in (3.4) (except without the factor $(1 + \delta)/(1 - \delta)$ in the first case). To control the tail with high probability, we appeal to Lemma B.5. First note that since $\alpha_0 + \alpha_1 + 2\alpha_3 > 1$ by Assumption 2.14, the positive sequence $\{j^{-2(\alpha_0 + \alpha_1 + 2\alpha_3)}\}$ is summable (and so is any shift of the sequence). Bounding the denominator by one, we obtain

$$\mathcal{I}_2^> \lesssim N \sum_{j > N_{\delta^-}^{1/u}} j^{-2(\alpha_0 + \alpha_1 + 2\alpha_3)} Z_j^{(N)} \leq N(1 + \delta) \sum_{j > N_{\delta^-}^{1/u}} j^{-2(\alpha_0 + \alpha_1 + 2\alpha_3)} = O\left(\left(\frac{1 + \delta}{1 - \delta}\right) N_{\delta^-}^{-\left(1 - \frac{\alpha_1 + 1/2 - \alpha_0}{\alpha_1 + \alpha_3}\right)}\right)$$

where the second inequality is from Lemma B.5 (upper tail only) and holds with probability at least $1 - \exp(-N\delta^2/8)$, and the last bound is Lemma B.2 similarly to preceding arguments.

The posterior covariance contribution \mathcal{I}_3^{\leq} has the same order as \mathcal{I}_2^{\leq} (by the same concentration inequality in Lemma B.6) and its tail $\mathcal{I}_3^>$ bounded a.s. as in the proof of Theorem 3.1.

Putting together the above bounds, noticing that $a + b(1 + \delta)/(1 - \delta) \lesssim (1 + \delta)/(1 - \delta)$ for any constants $a, b > 0$, and combining the probabilities with the union bound and the fact that there exists $C_{\delta, u} > 0$ and $0 < c < c'$ such that $\sup_{n \geq 1} n^{1/u} \exp(-(c' - c)n\delta^2) < C_{\delta, u}$, the theorem is proved. The final assertion follows immediately by dropping \mathcal{I}_3 from the analysis. ■

Proof of Theorem 3.5. Once again, recall that $\overline{g_j g_j^{(N)}} = w_j^2 Z_j^{(N)}$, where $N Z_j^{(N)} \sim \chi^2(N)$ and denote $u = 2(\alpha_1 + \alpha_3)$. The posterior sample test error satisfies $\mathbb{E}\mathcal{I}_1 + \mathbb{E}\mathcal{I}_2 + \mathbb{E}\mathcal{I}_3 \geq \mathbb{E}\mathcal{I}_2 + \mathbb{E}\mathcal{I}_3$; we do not lose any sharpness of our estimate by bounding $\mathbb{E}\mathcal{I}_1$ below by zero, since the lower bound contribution of $\Omega(N^{-2})$ from $\mathbb{E}\mathcal{I}_1$ (here obtained via Jensen's inequality) using the second case in (B.2) of Lemma B.1 is negligible when valid. For $\mathbb{E}\mathcal{I}_3$, since $r \mapsto (1 + ar)^{-1}$ is convex on $[0, \infty)$ for all $a \geq 0$, Jensen's inequality and Lemma B.7 yields

$$(A.6) \quad \mathbb{E} \sum_{j=1}^{\infty} \frac{\vartheta_j^2 \sigma_j^2}{1 + N\gamma^{-2}\sigma_j^2 \overline{g_j g_j^{(N)}}} \geq \sum_{j=1}^{\infty} \frac{\vartheta_j^2 \sigma_j^2}{1 + N\gamma^{-2}\sigma_j^2 \mathbb{E}[w_j^2 Z_j^{(N)}]} \asymp \sum_{j=1}^{\infty} \frac{j^{-2(\alpha_0 + \alpha_3)}}{1 + Nj^{-2(\alpha_1 + \alpha_3)}},$$

which has exact order the right hand side of (3.5) (without τ_N) by (B.3b) in Lemma B.2 (applied with $t = 2(\alpha_0 + \alpha_3)$, $v = 1$, and condition $t > 1$ satisfied by Assumption 2.14). Next,

$$\mathbb{E}\mathcal{I}_2 \geq \mathbb{E}\mathcal{I}_2^{\leq} = \mathbb{E} \sum_{j \leq N^{1/u}} \frac{N\vartheta_j^2 \gamma^{-2} \sigma_j^4 \overline{g_j g_j^{(N)}}}{(1 + N\gamma^{-2} \sigma_j^2 \overline{g_j g_j^{(N)}})^2} = \mathbb{E} \sum_{j \leq N^{1/u}} \frac{N\vartheta_j^2 \gamma^{-2} \sigma_j^4 w_j^2 Z_j^{(N)}}{(1 + N\gamma^{-2} \sigma_j^2 w_j^2 Z_j^{(N)})^2}.$$

For any $\tau_N \rightarrow 0$, define the sequence of events $\{A_j^{(N)}\}_{j, N \geq 1}$ by

$$(A.7) \quad A_j^{(N)} := \{\omega \in \Omega : Z_j^{(N)}(\omega) \geq \tau_N\}.$$

By the law of total expectation,

$$\begin{aligned} \mathbb{E}\mathcal{I}_2^{\leq} &= N \sum_{j \leq N^{1/u}} \vartheta_j^2 \gamma^{-2} \sigma_j^4 w_j^2 \mathbb{E} \left[\frac{Z_j^{(N)}}{(1 + N\gamma^{-2} \sigma_j^2 w_j^2 Z_j^{(N)})^2} \middle| A_j^{(N)} \right] \mathbb{P}(A_j^{(N)}) \\ &\quad + N \sum_{j \leq N^{1/u}} \vartheta_j^2 \gamma^{-2} \sigma_j^4 w_j^2 \mathbb{E} \left[\frac{Z_j^{(N)}}{(1 + N\gamma^{-2} \sigma_j^2 w_j^2 Z_j^{(N)})^2} \middle| (A_j^{(N)})^c \right] \mathbb{P}(A_j^{(N)})^c. \end{aligned}$$

The second term in the above display is nonnegative, so we obtain

$$\begin{aligned} \mathbb{E}\mathcal{I}_2^{\leq} &\geq \tau_N N \sum_{j \leq N^{1/u}} \vartheta_j^2 \gamma^{-2} \sigma_j^4 w_j^2 \mathbb{E} \left[\frac{1}{(1 + N\gamma^{-2} \sigma_j^2 w_j^2 Z_j^{(N)})^2} \middle| A_j^{(N)} \right] \mathbb{P}(A_j^{(N)}) \\ &\geq \tau_N N \sum_{j \leq N^{1/u}} \frac{\vartheta_j^2 \gamma^{-2} \sigma_j^4 w_j^2 \mathbb{P}(A_j^{(N)})}{(1 + N\gamma^{-2} \sigma_j^2 w_j^2 \mathbb{E}[Z_j^{(N)} | A_j^{(N)}])^2} \\ &= \tau_N N \sum_{j \leq N^{1/u}} \frac{\vartheta_j^2 \gamma^{-2} \sigma_j^4 w_j^2 \mathbb{P}(A_j^{(N)})^3}{(\mathbb{P}(A_j^{(N)}) + N\gamma^{-2} \sigma_j^2 w_j^2 \mathbb{E}[\mathbb{1}_{A_j^{(N)}} Z_j^{(N)}])^2}, \end{aligned}$$

where we applied conditional Jensen's inequality to yield the second line since $r \mapsto (1 + ar)^{-2}$ is convex on $[0, \infty)$ for any $a \geq 0$. Using $\mathbb{E}[\mathbb{1}_A Z_j^{(N)}] \leq \mathbb{E}[Z_j^{(N)}] = 1$, $\mathbb{P}(A) \leq 1$ for any $A \in \mathcal{F}$ and applying $1 + Nj^{-u} \asymp Nj^{-u}$ for $j \leq N^{1/u}$ twice yields

$$\mathbb{E}\mathcal{I}_2^{\leq} \gtrsim \tau_N N \sum_{j \leq N^{1/u}} \frac{j^{-2(\alpha_0 + \alpha_1 + 2\alpha_3)} \mathbb{P}(A_j^{(N)})^3}{(1 + Nj^{-2(\alpha_1 + \alpha_3)})^2} \asymp \tau_N \sum_{j \leq N^{1/u}} \frac{j^{-2(\alpha_0 + \alpha_3)} \mathbb{P}(A_j^{(N)})^3}{1 + Nj^{-2(\alpha_1 + \alpha_3)}}.$$

Next, by Markov's inequality and Lemma B.7 it holds that

$$\sup_{j \geq 1} \mathbb{P}(A_j^{(N)})^c = \sup_{j \geq 1} \mathbb{P}\{(Z_j^{(N)})^{-1} > \tau_N^{-1}\} \leq \sup_{j \geq 1} \tau_N \mathbb{E}[(Z_j^{(N)})^{-1}] = \tau_N N(N-2)^{-1} \rightarrow 0$$

as $N \rightarrow \infty$, which implies $\inf_{j \geq 1} \mathbb{P}(A_j^{(N)}) \rightarrow 1$. Using this, a short calculation gives

$$\mathbb{E}\mathcal{I}_2 \gtrsim \tau_N \sum_{j \leq N^{1/u}} \frac{j^{-2(\alpha_0 + \alpha_3)} \mathbb{P}(A_j^{(N)})^3}{1 + Nj^{-2(\alpha_1 + \alpha_3)}} = \Theta \left(\tau_N \sum_{j \leq N^{1/u}} \frac{j^{-2(\alpha_0 + \alpha_3)}}{1 + Nj^{-2(\alpha_1 + \alpha_3)}} \right)$$

as $N \rightarrow \infty$. Since $\tau_N \rightarrow 0$, the above display is negligible relative to lower bound (A.6) from $\mathbb{E}\mathcal{I}_3$. Thus the posterior sample estimator enjoys the asserted rate, while $\mathbb{E}\|L^\dagger - \bar{L}^{(N)}\|_{L^2_\nu(H;H)}^2 = \mathbb{E}\mathcal{I}_1 + \mathbb{E}\mathcal{I}_2$ implies that the posterior mean only admits lower bound (3.5). ■

Proof of Theorem 3.6. As $1 < 1 + \delta < 2$ is bounded above and below, we need not track the dependence on $1 + \delta$ in what follows. The proof proceeds as in the proof of Theorem 3.3 by applying Lemma B.6, but this time splitting the series at the critical index $j = N^{1/u}$ (since $N_{\delta+} := (1 + \delta)N \asymp N$). We first lower bound the error by $\mathcal{I}_2^\leq + \mathcal{I}_3^\leq$ as all terms are nonnegative. The remaining calculations follow directly from Lemma B.2 and are omitted. For the error associated with the posterior mean, the only contribution comes from \mathcal{I}_2^\leq , whose lower bound via Lemma B.6 has the (potentially very small) pre-factor $1 - \delta$ as asserted, and this factor is ignored in the display (3.6) since the similar bound from \mathcal{I}_3^\leq has no $1 - \delta$ dependence. ■

Proof of Theorem 3.8. The generalization gap involves expectation over \mathcal{P} due to the expected risk, and its modulus averaged over the data, $\mathcal{G}_N^\mathbb{E}$, involves expectation with respect to $(x_n, y_n) \sim \mathcal{P}$ i.i.d.. Note that \mathcal{P} is defined by (2.1) with $x \sim \mathcal{N}(0, \mathcal{C}_1)$, white noise $\eta \sim \mathcal{N}(0, \gamma^2 \text{Id})$ and $L = L^\dagger$; the properties of \mathcal{C}_1 and L^\dagger are given in Assumption 2.14. Again we recall that $\overline{g_j g_j}^{(N)} = w_j^2 Z_j^{(N)}$, where $N Z_j^{(N)} \sim \chi^2(N)$ and that $Z_j^{(N)}$ is of order 1 with respect to N by Lemma B.7. Revisiting (2.19) in Subsection 2.3.2 and (2.24), in coordinates with respect to the output eigenbasis $\{\varphi_j\}_{j \in \mathbb{N}}$ of L^\dagger we have $\mathcal{G}_N^\mathbb{E} = \mathbb{E}|\mathcal{J}_1 + \mathcal{J}_2 + \mathcal{J}_3|$, where (A.8)

$$\mathcal{J}_1 = \frac{1}{2} \sum_{j=1}^{\infty} (w_j^2 - \overline{g_j g_j}^{(N)}) |\bar{\ell}_j^{(N)} - \ell_j^\dagger|^2, \quad \mathcal{J}_2 = \frac{1}{2} \sum_{j=1}^{\infty} (\overline{g_j g_j}^{(N)} - w_j^2) (\ell_j^\dagger)^2, \quad \mathcal{J}_3 = \sum_{j=1}^{\infty} \gamma \overline{g_j \xi_j}^{(N)} \bar{\ell}_j^{(N)}.$$

Using the explicit form (2.26) of posterior mean sequence $\{\bar{\ell}_j^{(N)}\}$, we find that $\mathcal{G}_N^\mathbb{E}$ equals (A.9)

$$\mathbb{E} \left| \frac{1}{2} \sum_{j=1}^{\infty} (w_j^2 - \overline{g_j g_j}^{(N)}) \frac{(\ell_j^\dagger)^2 + N \gamma^{-2} \sigma_j^4 \overline{g_j g_j}^{(N)}}{(1 + N \gamma^{-2} \sigma_j^2 \overline{g_j g_j}^{(N)})^2} + \mathcal{J}_2 + \sum_{j=1}^{\infty} \frac{\gamma^2 (\overline{g_j \xi_j}^{(N)})^2 + \ell_j^\dagger \gamma \overline{g_j g_j}^{(N)} \overline{g_j \xi_j}^{(N)}}{N^{-1} \gamma^2 \sigma_j^{-2} + \overline{g_j g_j}^{(N)}} \right|.$$

Although at this point (A.8) and (A.9) involve purely formal random series, the computations in the remainder of the proof along with Lemma B.3 imply their \mathbb{P} -a.s. convergence.

Again following the conventions in the proof of Theorem 3.1, by the triangle inequality, we upper bound $\mathcal{G}_N^\mathbb{E}$ in (A.9) by five terms G_i , $i \in \{1, \dots, 5\}$ (with G_1, G_2 corresponding to $\mathbb{E}|\mathcal{J}_1|$, G_3 to $\mathbb{E}|\mathcal{J}_2|$, and G_4, G_5 to $\mathbb{E}|\mathcal{J}_3|$).

By triangle and Jensen's inequality, we have the $G_3 = \mathbb{E}|\mathcal{J}_2|$ upper bound

$$(A.10) \quad \frac{1}{2} \sum_{j=1}^{\infty} (\ell_j^\dagger)^2 \mathbb{E} |\overline{g_j g_j}^{(N)} - w_j^2| \leq \frac{1}{2} \sum_{j=1}^{\infty} (\ell_j^\dagger)^2 \sqrt{\text{Var}[\overline{g_j g_j}^{(N)}]} = \frac{\sqrt{2}}{2} \|L^\dagger\|_{L^2_\nu(H;H)}^2 N^{-1/2}.$$

using the fact that the variance factor evaluates to $(2w_j^4/N)^{1/2}$.

For the first term G_1 , we upper bound G_1^\leq by

$$\frac{1}{2} \mathbb{E} \sum_{j \leq N^{1/u}} \frac{|w_j^2 - \overline{g_j g_j}^{(N)}| (\ell_j^\dagger)^2}{(1 + N \gamma^{-2} \sigma_j^2 \overline{g_j g_j}^{(N)})^2} \leq \frac{1}{2} \mathbb{E} |(Z_1^{(N)})^{-2} - (Z_1^{(N)})^{-1}| \sum_{j \leq N^{1/u}} \frac{w_j^2 (\ell_j^\dagger)^2}{(N \gamma^{-2} \sigma_j^2 w_j^2)^2}$$

since the $\{Z_j^{(N)}\}_{j \geq 1}$ are identically distributed. Then

$$(\mathbb{E}|(Z_1^{(N)})^{-2} - (Z_1^{(N)})^{-1}|)^2 \leq \mathbb{E}|(Z_1^{(N)})^{-2} - (Z_1^{(N)})^{-1}|^2 = N^{-2}(2N + 48) = \Theta(N^{-1})$$

as $N \rightarrow \infty$ by Jensen's inequality and Lemma B.7 applied three times. Similar to (A.4) in the proof of Theorem 3.1, for any $-\alpha_1 < s < \beta$, by Lemma B.1 (with $t = 2\alpha_1$, $q = s$, $v = 2$),

$$(A.11) \quad G_1^{\leq} = O\left(N^{-\frac{1}{2}} \sum_{j \leq N^{1/u}} \frac{j^{-2\alpha_1}(\ell_j^\dagger)^2}{(1 + Nj^{-u})^2}\right) = \begin{cases} o(N^{-(\frac{1}{2} + \frac{\alpha_1+s}{\alpha_1+\alpha_3})}), & \text{if } \frac{\alpha_1+s}{\alpha_1+\alpha_3} < 2 \\ O(N^{-5/2}), & \text{if } \frac{\alpha_1+s}{\alpha_1+\alpha_3} \geq 2. \end{cases}$$

In particular, $G_1^{\leq} = o(N^{-1/2})$ in the first case since $\alpha_1 + s > 0$. Using the Jensen's inequality variance bound as for G_3 and similar to (A.5) in the proof of Theorem 3.1, we find that the tail sum $G_1^>$ is never bigger than G_1^{\leq} in (A.11) (and of matching order in the first case).

For the second term G_2 associated with \mathcal{J}_1 , we upper bound G_2^{\leq} by

$$\frac{1}{2} \mathbb{E} \sum_{j \leq N^{1/u}} \frac{|w_j^2 - \overline{g_j g_j^{(N)}}| N \gamma^{-2} \sigma_j^4 \overline{g_j g_j^{(N)}}}{(1 + N \gamma^{-2} \sigma_j^2 \overline{g_j g_j^{(N)}})^2} \leq \frac{1}{2} \mathbb{E}|(Z_1^{(N)})^{-1} - 1| \sum_{j \leq N^{1/u}} \frac{\gamma^2}{N},$$

and application of Jensen's inequality and Lemma B.7 twice yields

$$(A.12) \quad G_2^{\leq} \leq \frac{\gamma^2}{2} \sqrt{\mathbb{E}|(Z_1^{(N)})^{-1} - 1|^2} N^{-(1-\frac{1}{u})} = \frac{\gamma^2}{2} \sqrt{\frac{2N+8}{(N-2)(N-4)}} N^{-(1-\frac{1}{u})} = O(N^{-(\frac{3}{2} - \frac{1/2}{\alpha_1+\alpha_3})})$$

as $N \rightarrow \infty$. To check that the tail series is never bigger than this sum, we estimate

$$G_2^> \leq \frac{1}{2} \sqrt{\mathbb{E}|Z_j^{(N)} - (Z_j^{(N)})^2|^2} \sum_{j > N^{1/u}} N \gamma^{-2} \sigma_j^4 w_j^4 = O\left(N^{-1/2} \sum_{j > N^{1/u}} N j^{-2(2\alpha_1+2\alpha_3)}\right)$$

as $N \rightarrow \infty$ by Lemma B.7 three times (since the square of the square root expectation factor evaluates to $N^{-3}(2N^2 + 28N + 48) = \Theta(N^{-1})$ as $N \rightarrow \infty$). By (B.3a) in Lemma B.2 (applied with $t = 2(2\alpha_1 + 2\alpha_3)$), the rightmost series above is of order the upper bound (A.12) on G_2^{\leq} .

Moving on to the contributions G_4 and G_5 from $\mathbb{E}|\mathcal{J}_3|$, we first average out the random variables $\{\xi_{jn}\} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$ associated with the noise structure of our statistical model. First,

$$G_4 = \mathbb{E} \sum_{j=1}^{\infty} \frac{\gamma^2 (\overline{g_j \xi_j^{(N)}})^2}{N^{-1} \gamma^2 \sigma_j^{-2} + \overline{g_j g_j^{(N)}}} = \mathbb{E}^\nu \sum_{j=1}^{\infty} \frac{N \sigma_j^2 \mathbb{E}^{\pi_0}[(\overline{g_j \xi_j^{(N)}})^2]}{1 + N \gamma^{-2} \sigma_j^2 \overline{g_j g_j^{(N)}}} = \mathbb{E}^\nu \sum_{j=1}^{\infty} \frac{\sigma_j^2 \overline{g_j g_j^{(N)}}}{1 + N \gamma^{-2} \sigma_j^2 \overline{g_j g_j^{(N)}}}.$$

Since the map $r \mapsto r(1+ar)^{-1}$ is concave on $[0, \infty)$ for all $a \geq 0$, Jensen's inequality yields

$$(A.13) \quad G_4 \leq \sum_{j=1}^{\infty} \frac{\sigma_j^2 w_j^2}{1 + N \gamma^{-2} \sigma_j^2 w_j^2} \asymp \sum_{j=1}^{\infty} \frac{j^{-u}}{1 + N j^{-u}} = O(N^{-(\frac{\alpha_1+\alpha_3-1/2}{\alpha_1+\alpha_3})})$$

as $N \rightarrow \infty$, where the last bound follows from (B.3b) in Lemma B.2 (applied with $t = 2(\alpha_1 + \alpha_3)$ and $v = 1$, which satisfies the first case).

Similarly, for the last series we have, by Jensen's inequality applied to the entire series,

$$G_5 = \mathbb{E} \left| \sum_{j=1}^{\infty} \frac{\ell_j^\dagger \gamma \overline{g_j g_j^{(N)}} \overline{g_j \xi_j^{(N)}}}{N^{-1} \gamma^2 \sigma_j^{-2} + \overline{g_j g_j^{(N)}}} \right| \leq \left(\mathbb{E}^\nu \mathbb{E}^{\pi_0} \left| \sum_{j=1}^{\infty} \frac{N \gamma^{-1} \sigma_j^2 \ell_j^\dagger \overline{g_j g_j^{(N)}} \overline{g_j \xi_j^{(N)}}}{1 + N \gamma^{-2} \sigma_j^2 \overline{g_j g_j^{(N)}}} \right|^2 \right)^{1/2}.$$

Bringing the expectation over π_0 into the now doubly indexed series, we compute the factor

$$\mathbb{E}^{\pi_0} [\overline{g_j \xi_j^{(N)}} \overline{g_{j'} \xi_{j'}^{(N)}}] = \frac{1}{N^2} \sum_{n, n' \leq N} g_{jn} g_{j'n'} \mathbb{E}^{\pi_0} [\xi_{jn} \xi_{j'n'}] = \frac{\delta_{jj'}}{N} \left(\frac{1}{N} \sum_{n=1}^N g_{jn} g_{j'n} \right)$$

for any $j, j' \geq 1$, since the $\{\xi_{jn}\}$ are i.i.d. standard Gaussian. Thus,

(A.14)

$$G_5 \leq \left(\mathbb{E} \sum_{j=1}^{\infty} \frac{N (\ell_j^\dagger)^2 \gamma^{-2} \sigma_j^4 (\overline{g_j g_j^{(N)}})^3}{(1 + N \gamma^{-2} \sigma_j^2 \overline{g_j g_j^{(N)}})^2} \right)^{1/2} \leq \left(\sum_{j=1}^{\infty} \frac{(\ell_j^\dagger)^2 \gamma^2 w_j^2}{N} \right)^{1/2} = \gamma \|L^\dagger\|_{L_v^2(H; H)} N^{-1/2}.$$

Comparing (A.10)–(A.14) and deducing that $G_2 = o(N^{-1/2})$ since $\alpha_1 + \alpha_3 > 1/2$ by [Assumption 2.14](#), we conclude that $\mathcal{G}_N^\mathbb{E} = O(N^{-1/2} + G_4)$ as desired. \blacksquare

Proof of Theorem 3.9. By the triangle inequality, $\mathcal{G}_N^\mathbb{E} \geq \mathbb{E}|\mathcal{J}_3 + \mathcal{J}_2| - \mathbb{E}|\mathcal{J}_1| \geq |\mathbb{E}\mathcal{J}_3 + \mathbb{E}\mathcal{J}_2| - \mathbb{E}|\mathcal{J}_1| = |\mathbb{E}\mathcal{J}_3| - \mathbb{E}|\mathcal{J}_1|$. We first develop a lower bound on $|\mathbb{E}\mathcal{J}_3|$. For any $\tau_N \rightarrow 0$, we define the sequence of events $\{A_j^{(N)}\}_{j, N \geq 1}$ from (A.7), and following the same argument in the proof of [Theorem 3.5](#) (details omitted), we obtain that $|\mathbb{E}\mathcal{J}_3|$ is lower bounded by

$$(A.15) \quad \tau_N \sum_{j=1}^{\infty} \frac{j^{-2(\alpha_1 + \alpha_3)} \mathbb{P}(A_j^{(N)})^2}{1 + N j^{-2(\alpha_1 + \alpha_3)}} = \Theta \left(\tau_N \sum_{j=1}^{\infty} \frac{j^{-2(\alpha_1 + \alpha_3)}}{1 + N j^{-2(\alpha_1 + \alpha_3)}} \right) = \Theta \left(\tau_N N^{-\left(\frac{\alpha_1 + \alpha_3 - 1/2}{\alpha_1 + \alpha_3}\right)} \right)$$

as $N \rightarrow \infty$, the last bound by routine application of [Lemma B.2](#) (with $t = 2(\alpha_1 + \alpha_3) = u$, $v = 1$, and condition $t > 1$ satisfied by [Assumption 2.14](#)). To conclude the proof, we claim that the upper bound on $\mathbb{E}|\mathcal{J}_1|$ from (A.11) and (A.12) in [Theorem 3.8](#) is asymptotically negligible relative to the above display (A.15) under the conditions assumed. Indeed, choosing $\tau_N \gg N^{-1/2}$ (i.e., $N^{-1/2} = o(\tau_N)$) ensures that (A.15) asymptotically dominates (A.12). The other $\mathbb{E}|\mathcal{J}_1|$ contribution (A.11) depends on regularity $s \in (-\alpha_1, \beta)$, and the tightest bound is obtained for maximal s , that is, $s = \beta - \varepsilon$ for any $\varepsilon > 0$ arbitrarily small. Consider the case $\frac{\alpha_1 + s}{\alpha_1 + \alpha_3} \geq 2$, which holds whenever $\frac{\alpha_1 + \beta}{\alpha_1 + \alpha_3} > 2$ for ε sufficiently small. Then (A.11) is $O(N^{-5/2})$, and this is certainly negligible relative to (A.15) since we have $\tau_N \gg N^{-1/2} \gg N^{-1/2} N^{-1 - \frac{1/2}{\alpha_1 + \alpha_3}} = N^{-5/2} N^{\frac{\alpha_1 + \alpha_3 - 1/2}{\alpha_1 + \alpha_3}}$. To prove the final assertion in the case $\frac{\alpha_1 + \beta}{\alpha_1 + \alpha_3} \leq 2$, we see that (A.15) dominates (A.11) if $\tau_N \gg N^{\frac{\alpha_1 + \alpha_3 - 1/2}{\alpha_1 + \alpha_3}} N^{-\left(\frac{1}{2} + \frac{\alpha_1 + \beta - \varepsilon}{\alpha_1 + \alpha_3}\right)} = N^{-\left(\frac{1}{2} - \varepsilon' + \frac{1/2 + \beta - \alpha_3}{\alpha_1 + \alpha_3}\right)}$ for any sufficiently small $\varepsilon' > 0$. This statement is non-vacuous if the exponent in parentheses on the right is positive, and to this end it suffices that $\alpha_3 < 1 + \alpha_1 + 2\beta$ as asserted. \blacksquare

Appendix B. Supporting Lemmas.

Lemma B.1 ([35], Lemma 8.1, pg. 2653). *Let $q \in \mathbb{R}$, $t \geq -2q$, $u > 0$, and $v \geq 0$. Then for any $a \in \mathcal{H}^q(\mathbb{N}; \mathbb{R})$ and $N \in \mathbb{N}$,*

$$(B.1) \quad \sum_{j > N^{1/u}} \frac{j^{-t} a_j^2}{(1 + N j^{-u})^v} \asymp \sum_{j > N^{1/u}} j^{-t} a_j^2 \leq \left(\sum_{j > N^{1/u}} j^{2q} a_j^2 \right) N^{-\left(\frac{t+2q}{u}\right)}.$$

Furthermore, as $N \rightarrow \infty$,

$$(B.2) \quad \sum_{j \leq N^{1/u}} \frac{j^{-t} a_j^2}{(1 + Nj^{-u})^v} = \begin{cases} o(N^{-(\frac{t+2q}{u})}), & \text{if } (t+2q)/u < v \\ \Theta\left(\sum_{j=1}^{\infty} \frac{j^{-t} a_j^2}{(1 + Nj^{-u})^v}\right) = \Theta(N^{-v}), & \text{if } (t+2q)/u \geq v. \end{cases}$$

Lemma B.2 ([35], Lemma 8.2, pg. 2654). *Let $t > 1$, $u > 0$, and $v \geq 0$. Then as $N \rightarrow \infty$,*

$$(B.3a) \quad \sum_{j > N^{1/u}} \frac{j^{-t}}{(1 + Nj^{-u})^v} \asymp \sum_{j > N^{1/u}} j^{-t} = \Theta(N^{-(\frac{t-1}{u})}),$$

$$(B.3b) \quad \sum_{j=1}^{\infty} \frac{j^{-t}}{(1 + Nj^{-u})^v} = \Theta\left(\sum_{j \leq N^{1/u}} \frac{j^{-t}}{(1 + Nj^{-u})^v}\right) = \begin{cases} \Theta(N^{-(\frac{t-1}{u})}), & \text{if } (t-1)/u < v \\ \Theta(N^{-v} \log N), & \text{if } (t-1)/u = v \\ \Theta(N^{-v}), & \text{if } (t-1)/u > v. \end{cases}$$

Proof. The assertions follow from [35, pg. 2654–2655]: choose the slowly varying function used there to be identically constant, $q = -1/2$, and use the fact that $\sum_{j=1}^J 1/j = \Theta(\log J)$. ■

Lemma B.3. *Let $\{X_k\}_{k \geq 1}$ be a sequence of (possibly dependent) real random variables. If $\sum_{k=1}^{\infty} \mathbb{E}|X_k| < \infty$, then $\sum_{k=1}^n X_k \xrightarrow{a.s.} \sum_{k=1}^{\infty} X_k$ as $n \rightarrow \infty$.*

Proof. An application of monotone convergence shows that $\sum_k |X_k|$ converges a.s. ■

In what follows, let $\text{SE}(v^2, a)$ denote the set of real-valued sub-exponential random variables with parameters $(v^2, a) \in \mathbb{R}_{\geq 0}^2$ (also $v \geq 0$), so that $X \in \text{SE}(v^2, a)$ satisfies the moment generating function (MGF) bound $\mathbb{E} \exp(\theta(X - \mathbb{E}X)) \leq \exp(v^2 \theta^2 / 2)$ for all $|\theta| < 1/a$.

Lemma B.4 (Closure of Sub-Exponential Random Variables Under Addition). *For $n \in \mathbb{N}$, if $\{X_k\}_{k=1}^n$ are (possibly dependent) real-valued random variables with $X_k \in \text{SE}(v_k^2, a_k)$, then*

$$(B.4) \quad \sum_{k=1}^n (X_k - \mathbb{E}X_k) \in \text{SE}\left(\left(\sum_{k=1}^n v_k\right)^2, \left(\sum_{k=1}^n v_k\right) \max_{i \leq n} \frac{a_i}{v_i}\right).$$

Proof. We estimate

$$\begin{aligned} \mathbb{E} \exp\left(\theta \sum_{k=1}^n (X_k - \mathbb{E}X_k)\right) &= \mathbb{E} \prod_{k=1}^n \exp(\theta(X_k - \mathbb{E}X_k)) \leq \prod_{k=1}^n (\mathbb{E} \exp(\theta(X_k - \mathbb{E}X_k) p_k))^{1/p_k} \\ &\leq \prod_{k=1}^n (\exp(v_k^2 \theta^2 p_k^2 / 2))^{1/p_k} = \exp\left(\left(\sum_{k=1}^n v_k\right)^2 \frac{\theta^2}{2}\right), \end{aligned}$$

where we used the generalized Hölder's inequality to yield the first inequality with $\sum_{i=1}^n 1/p_i = 1$, $p_i := v_i^{-1} \sum_{k=1}^n v_k$, and the sub-exponential MGF bound for all $k \leq n$ to yield the second, which holds for all $|\theta| < \min_{i \leq n} (p_i a_i)^{-1} = (\max_{i \leq n} p_i a_i)^{-1}$ as required. ■

The following lemma is useful for controlling tail sums in our high probability bounds.

Lemma B.5. *Let $\{Z_j\}_{j \geq 1}$ be a (possibly dependent) identically distributed family of $\chi^2(n)$ random variables for some $n \in \mathbb{N}$, and let $u \in \ell^1(\mathbb{N}; \mathbb{R})$ be a nonnegative sequence. Fix $\delta \in (0, 1)$. Then with probability at least $1 - 2\exp(-n\delta^2/8)$,*

$$(B.5) \quad (1 - \delta) \sum_{j=1}^{\infty} u_j \leq \sum_{j=1}^{\infty} u_j Z_j / n \leq (1 + \delta) \sum_{j=1}^{\infty} u_j.$$

Proof. Define $Y_J := \sum_{j \leq J} u_j Z_j / n$ for any $J \in \mathbb{N}$. Since $Z_j / n \in \text{SE}(4/n, 4/n)$ [57, Sec. 2.1.3], it follows from Lemma B.4 that $Y_J \in \text{SE}(\frac{4}{n} \|\{u_j\}_{j \leq J}\|_1^2, \frac{4}{n} \|\{u_j\}_{j \leq J}\|_1)$. Noting that

$$\sum_{j=1}^{\infty} \mathbb{E}[u_j Z_j / n] = \sum_{j=1}^{\infty} \mathbb{E}[u_j Z_j / n] = \sum_{j=1}^{\infty} (u_j / n)(n) = \|u\|_{\ell^1} < \infty$$

holds by supposition, we deduce $Y_J \rightarrow Y_{\infty}$ as $J \rightarrow \infty$ \mathbb{P} -a.s. by monotone convergence (Lemma B.3). By Fatou's lemma applied to the MGF bound for Y_J , we conclude $Y_{\infty} \in \text{SE}(\frac{4}{n} \|u\|_{\ell^1}^2, \frac{4}{n} \|u\|_{\ell^1})$. Therefore, using $\mathbb{E}Y_{\infty} = \|u\|_{\ell^1}$ and the sub-exponential tail bound [57, Prop. 2.9] gives $\mathbb{P}\{|Y_{\infty} - \mathbb{E}Y_{\infty}| \leq \mathbb{E}Y_{\infty} \delta\} \geq 1 - 2\exp(-n\delta^2/8)$ for any $\delta \in (0, 1)$ as desired. ■

The next lemma, while not as tight as the previous one, is elementary and easier to apply to the rational functions of random variables that arise in the proofs of results in Section 3.

Lemma B.6. *Let $\{Z_j^{(n)}\}_{j \leq J}$ be a (possibly dependent) identically distributed family of random variables such that $nZ_j^{(n)} \sim \chi^2(n)$ for some $n, J \in \mathbb{N}$, and fix $\delta \in (0, 1)$. Then with probability at least $1 - 2J\exp(-n\delta^2/8)$, $(1 - \delta) \leq Z_j^{(n)} \leq (1 + \delta)$ for all $j \leq J$.*

Proof. The result follows immediately from [57, Ex. 2.11, pg. 29] and the union bound. ■

Lemma B.7 (Chi-Square Moments). *Let $W \sim \chi^2(n)$ be a chi-square random variable with $n \in \mathbb{N}$ degrees of freedom. Then for any $p > -n/2$,*

$$(B.6) \quad \mathbb{E}[W^p] = 2^p \frac{\Gamma(p + n/2)}{\Gamma(n/2)},$$

where Γ is Euler's complete gamma function. In particular,

$$(B.7) \quad \mathbb{E}[W^{-1}] = \frac{1}{(n-2)}, \quad \mathbb{E}[W^{-2}] = \frac{1}{(n-2)(n-4)}, \quad \mathbb{E}[W^{-3}] = \frac{1}{(n-2)(n-4)(n-6)}.$$

Proof. A direct calculation with the PDF of $\chi^2(n)$ yields (B.6) in closed form. ■

Acknowledgments. The authors thank Kamyar Azizzadenesheli and Joel A. Tropp for helpful discussions about statistical learning. The computations presented in this paper were conducted on the Resnick High Performance Computing Center, a facility supported by the Resnick Sustainability Institute at the California Institute of Technology.

REFERENCES

- [1] J. ABERNETHY, F. BACH, T. EVGENIOU, AND J.-P. VERT, *A new approach to collaborative filtering: Operator estimation with spectral regularization.*, Journal of Machine Learning Research, 10 (2009).
- [2] B. ADCOCK, S. BRUGIAPAGLIA, N. DEXTER, AND S. MORAGA, *Deep neural networks are effective at learning high-dimensional Hilbert-valued functions from limited data*, arXiv preprint arXiv:2012.06081, (2020).

- [3] S. AGAPIOU, J. M. BARDSLEY, O. PAPASPILIOPOULOS, AND A. M. STUART, *Analysis of the Gibbs sampler for hierarchical inverse problems*, SIAM/ASA Journal on Uncertainty Quantification, 2 (2014), pp. 511–544.
- [4] S. AGAPIOU, S. LARSSON, AND A. M. STUART, *Posterior contraction rates for the Bayesian approach to linear ill-posed inverse problems*, Stochastic Processes and their Applications, 123 (2013), pp. 3828–3860.
- [5] S. AGAPIOU AND P. MATHÉ, *Posterior contraction in Bayesian inverse problems under Gaussian priors*, in New trends in parameter identification for mathematical models, Springer, 2018, pp. 1–29.
- [6] S. AGAPIOU AND P. MATHÉ, *Designing truncated priors for direct and inverse Bayesian problems*, arXiv preprint arXiv:2105.10254, (2021).
- [7] S. AGAPIOU, A. M. STUART, AND Y.-X. ZHANG, *Bayesian posterior contraction rates for linear severely ill-posed inverse problems*, Journal of Inverse and Ill-posed Problems, 22 (2014), pp. 297–321.
- [8] G. S. ALBERTI, E. DE VITO, M. LASSAS, L. RATTI, AND M. SANTACESARIA, *Learning the optimal regularizer for inverse problems*, arXiv preprint arXiv:2106.06513, (2021).
- [9] S. ARRIDGE, P. MAASS, O. ÖKTEM, AND C.-B. SCHÖNLIEB, *Solving inverse problems using data-driven models*, Acta Numerica, 28 (2019), pp. 1–174.
- [10] A. ASPRI, Y. KOROLEV, AND O. SCHERZER, *Data driven regularization by projection*, Inverse Problems, 36 (2020), p. 125009.
- [11] K. BHATTACHARYA, B. HOSSEINI, N. B. KOVACHKI, AND A. M. STUART, *Model reduction and neural networks for parametric PDEs*, The SMAI Journal of Computational Mathematics, 7 (2021), pp. 121–157.
- [12] I. R. BLEYER AND R. RAMLAU, *A double regularization approach for inverse problems with noisy data and inexact operator*, Inverse Problems, 29 (2013), p. 025004.
- [13] N. BOULLÉ AND A. TOWNSEND, *Learning elliptic partial differential equations with randomized linear algebra*, arXiv preprint arXiv:2102.00491, (2021).
- [14] S. L. BRUNTON, J. L. PROCTOR, AND J. N. KUTZ, *Discovering governing equations from data by sparse identification of nonlinear dynamical systems*, Proceedings of the national academy of sciences, 113 (2016), pp. 3932–3937.
- [15] T. A. BUBBA, M. GALINIER, M. LASSAS, M. PRATO, L. RATTI, AND S. SILTANEN, *Deep neural networks for inverse problems with pseudodifferential operators: An application to limited-angle tomography*, SIAM Journal on Imaging Sciences, 14 (2021), pp. 470–505, <https://doi.org/10.1137/20M1343075>.
- [16] T. T. CAI AND P. HALL, *Prediction in functional linear regression*, The Annals of Statistics, 34 (2006), pp. 2159–2179.
- [17] A. CAPONNETTO AND E. DE VITO, *Optimal rates for the regularized least-squares algorithm*, Foundations of Computational Mathematics, 7 (2007), pp. 331–368.
- [18] L. CAVALIER, *Nonparametric statistical inverse problems*, Inverse Problems, 24 (2008), p. 034004.
- [19] M. J. COLBROOK, A. HORNING, AND A. TOWNSEND, *Computing spectral measures of self-adjoint operators*, SIAM Review, 63 (2021), pp. 489–524.
- [20] C. CRAMBES AND A. MAS, *Asymptotics of prediction in functional linear regression with functional outputs*, Bernoulli, 19 (2013), pp. 2627–2651.
- [21] M. DASHTI, K. J. LAW, A. M. STUART, AND J. VOSS, *MAP estimators and their consistency in Bayesian nonparametric inverse problems*, Inverse Problems, 29 (2013), p. 095017.
- [22] M. DASHTI AND A. M. STUART, *The Bayesian Approach to Inverse Problems*, Springer International Publishing, Cham, 2017, pp. 311–428, https://doi.org/10.1007/978-3-319-12385-1_7.
- [23] M. V. DE HOOP, M. LASSAS, AND C. A. WONG, *Deep learning architectures for nonlinear operator functions and nonlinear inverse problems*, arXiv preprint arXiv:1912.11090, (2019).
- [24] E. DE VITO, L. ROSASCO, A. CAPONNETTO, U. D. GIOVANNINI, AND F. ODONE, *Learning from examples as an inverse problem*, Journal of Machine Learning Research, 6 (2005), pp. 883–904.
- [25] Y. FAN AND L. YING, *Solving electrical impedance tomography with deep learning*, Journal of Computational Physics, 404 (2020), p. 109119.
- [26] L. GAWARECKI AND V. MANDREKAR, *Stochastic differential equations in infinite dimensions: with applications to stochastic partial differential equations*, Springer Science & Business Media, 2010.
- [27] D. GIANNAKIS, *Data-driven spectral decomposition and forecasting of ergodic dynamical systems*, Applied and Computational Harmonic Analysis, 47 (2019), pp. 338–396.

- [28] G. H. GOLUB AND C. F. VAN LOAN, *An analysis of the total least squares problem*, SIAM journal on numerical analysis, 17 (1980), pp. 883–893.
- [29] S. GRÜNEWÄLDER, G. LEVER, L. BALDASSARRE, S. PATTERSON, A. GRETTON, AND M. PONTIL, *Conditional mean embeddings as regressors*, arXiv preprint arXiv:1205.4656, (2012).
- [30] S. HÖRMANN AND Ł. KIDZIŃSKI, *A note on estimation in Hilbertian linear models*, Scandinavian journal of statistics, 42 (2015), pp. 43–62.
- [31] I. KLEBANOV, I. SCHUSTER, AND T. J. SULLIVAN, *A rigorous theory of conditional mean embeddings*, SIAM Journal on Mathematics of Data Science, 2 (2020), pp. 583–606.
- [32] S. KLUS, F. NÜSKE, S. PEITZ, J.-H. NIEMANN, C. CLEMENTI, AND C. SCHÜTTE, *Data-driven approximation of the Koopman generator: Model reduction, system identification, and control*, Physica D: Nonlinear Phenomena, 406 (2020), p. 132416.
- [33] S. KLUS, I. SCHUSTER, AND K. MUANDET, *Eigendecompositions of transfer operators in reproducing kernel Hilbert spaces*, Journal of Nonlinear Science, 30 (2020), pp. 283–315.
- [34] B. KNAPIK, J.-B. SALOMOND, ET AL., *A general approach to posterior contraction in nonparametric inverse problems*, Bernoulli, 24 (2018), pp. 2091–2121.
- [35] B. T. KNAPIK, A. W. VAN DER VAART, J. H. VAN ZANTEN, ET AL., *Bayesian inverse problems with Gaussian priors*, The Annals of Statistics, 39 (2011), pp. 2626–2657.
- [36] Y. KOROLEV, *Two-layer neural networks with values in a Banach space*, arXiv preprint arXiv:2105.02095, (2021).
- [37] M. S. LEHTINEN, L. PAIVARINTA, AND E. SOMERSALO, *Linear inverse problems for generalised random variables*, Inverse Problems, 5 (1989), p. 599.
- [38] Z. LI, N. B. KOVACHKI, K. AZIZZADENESHELI, B. LIU, K. BHATTACHARYA, A. M. STUART, AND A. ANANDKUMAR, *Neural operator: Graph kernel network for partial differential equations*, ICLR 2020 Workshop on Integration of Deep Neural Models and Differential Equations, (2020), <https://arxiv.org/abs/2003.03485>.
- [39] L. LU, P. JIN, AND G. E. KARNIADAKIS, *Deeponet: Learning nonlinear operators for identifying differential equations based on the universal approximation theorem of operators*, arXiv preprint arXiv:1910.03193, (2019).
- [40] A. MANDELBAUM, *Linear estimators and measurable linear transformations on a Hilbert space*, Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete, 65 (1984), pp. 385–397.
- [41] T. MATHIEU AND S. MINSKER, *Excess risk bounds in robust empirical risk minimization*, Information and Inference: A Journal of the IMA, (2021).
- [42] C. A. MICCHELLI AND M. PONTIL, *On learning vector-valued functions*, Neural computation, 17 (2005), pp. 177–204.
- [43] M. MOLLENHAUER AND P. KOLTAI, *Nonparametric approximation of conditional expectation operators*, arXiv preprint arXiv:2012.12917, (2020).
- [44] N. H. NELSEN AND A. M. STUART, *The random feature model for input-output maps between Banach spaces*, arXiv preprint arXiv:2005.10224, (2020).
- [45] R. NICKL, S. VAN DE GEER, AND S. WANG, *Convergence rates for penalized least squares estimators in PDE constrained regression problems*, SIAM/ASA Journal on Uncertainty Quantification, 8 (2020), pp. 374–413.
- [46] T. O’LEARY-ROSEBERRY, U. VILLA, P. CHEN, AND O. GHATTAS, *Derivative-informed projected neural networks for high-dimensional parametric maps governed by PDEs*, arXiv preprint arXiv:2011.15110, (2020).
- [47] R. G. PATEL, N. A. TRASK, M. A. WOOD, AND E. C. CYR, *A physics-informed operator regression framework for extracting data-driven continuum models*, arXiv preprint arXiv:2009.11992, (2020).
- [48] T. PORTONE AND R. D. MOSER, *Bayesian inference of an uncertain generalized diffusion operator*, arXiv preprint arXiv:2105.01807, (2021).
- [49] J. O. RAMSAY AND B. W. SILVERMAN, *Functional data analysis*, Springer Series in Statistics, Springer, New York, second ed., 2005.
- [50] A. RASTOGI, G. BLANCHARD, P. MATHÉ, ET AL., *Convergence analysis of Tikhonov regularization for non-linear statistical inverse problems*, Electronic Journal of Statistics, 14 (2020), pp. 2798–2841.
- [51] M. REIMHERR, *Functional regression with repeated eigenvalues*, Statistics & Probability Letters, 107 (2015), pp. 62–70.

- [52] C. SCHWAB AND J. ZECH, *Deep learning in high dimension: Neural network expression rates for generalized polynomial chaos expansions in UQ*, Analysis and Applications, 17 (2019), pp. 19–55.
- [53] L. SONG, J. HUANG, A. SMOLA, AND K. FUKUMIZU, *Hilbert space embeddings of conditional distributions with applications to dynamical systems*, in Proceedings of the 26th Annual International Conference on Machine Learning, 2009, pp. 961–968.
- [54] I. STEINWART, *Convergence types and rates in generic Karhunen-Loeve expansions with applications to sample path properties*, Potential Analysis, 51 (2019), pp. 361–395.
- [55] A. M. STUART, *Inverse problems: a Bayesian perspective*, Acta numerica, 19 (2010), pp. 451–559.
- [56] P. TABAGHI, M. DE HOOP, AND I. DOKMANIĆ, *Learning Schatten–von Neumann operators*, arXiv preprint arXiv:1901.10076, (2019).
- [57] M. J. WAINWRIGHT, *High-dimensional statistics: A non-asymptotic viewpoint*, vol. 48, Cambridge University Press, 2019.
- [58] D. WANG, Z. ZHAO, Y. YU, AND R. WILLETT, *Functional linear regression with mixed predictors*, arXiv preprint arXiv:2012.00460, (2020).