

## COMPUTING ACTIVE SUBSPACES WITH MONTE CARLO

PAUL G. CONSTANTINE\* AND DAVID F. GLEICH†

**Abstract.** Active subspaces can effectively reduce the dimension of high-dimensional parameter studies enabling otherwise infeasible experiments with expensive simulations. The key components of active subspace methods are the eigenvectors of a symmetric, positive semidefinite matrix whose elements are the average products of partial derivatives of the simulation’s input/output map. We study a Monte Carlo method for approximating the eigenpairs of this matrix. We offer both theoretical results based on recent non-asymptotic random matrix theory and a practical approach based on the bootstrap. We extend the analysis to the case when the gradients are approximated, for example, with finite differences. Our goal is to provide guidance for two questions that arise in active subspaces: (i) How many gradient samples does one need to accurately approximate the eigenvalues and subspaces? (ii) What can be said about the accuracy of the estimated subspace, both theoretically and practically? We test the approach on both simple quadratic functions where the active subspace is known and a parameterized PDE with 100 variables characterizing the coefficients of the differential operator.

**Key words.** active subspaces, dimension reduction

**AMS subject classifications.**

**1. Introduction.** Engineering models typically contain several input parameters that must be specified to produce a set of model outputs that contains one or more quantities of interest. The engineer’s goal is to characterize the behavior of the quantities of interest as functions of the model’s inputs. However, parameter studies—such as optimization and uncertainty quantification—are challenging when the number of inputs is large and the model involves an expensive computer simulation. In such cases, the engineer may analyze the output’s sensitivity with respect to inputs to identify a subset of inputs whose variation changes the outputs the most [26]. In the best case, she can then limit parameter studies to key parameters and thus reduce the *dimension* of the parameter study. This approach is appropriate when varying important parameters changes the outputs much more than varying the unimportant parameters. However, a model’s output may depend on all the parameters through certain linear combinations, which generalizes seeking key parameters to seeking key directions in the parameter space. The *active subspace* identifies important directions in the model’s input space with respect to a particular quantity of interest; perturbing the inputs along these important directions changes the quantity of interest more, on average, than perturbing the inputs in orthogonal directions [7]. For parameter studies whose work depends exponentially on the number of parameters—e.g., integration or response surface construction—the active subspace-enabled dimension reduction can permit otherwise infeasible studies.

The active subspace is defined by a set of eigenvectors corresponding to large eigenvalues of the average outer product of the gradient with itself. These eigenpairs are properties of the map between model inputs and outputs, like Fourier coefficients or the Lipschitz constant. To determine if a function admits a low-dimensional active subspace—and thus reduce the dimension of the parameter studies—one must estimate these eigenpairs. This estimation is problematic because the elements of the matrix defining the eigenpairs are themselves high-dimensional integrals. Most de-

---

\*Department of Applied Mathematics and Statistics, Colorado School of Mines, Golden, CO 80401 (paul.constantine@mines.edu).

†Department of Computer Science, Purdue University, West Lafayette, IN (dgleich@purdue.edu)

terministic numerical integration rules are impractical beyond a handful of variables, especially if the integrand is costly to evaluate. We therefore focus on a Monte Carlo approach to approximate the eigenpairs, where we take advantage of recent theoretical results that bound the number of samples needed to approximate the spectrum of sums of random matrices. Monte Carlo is attractive because it makes few restrictions on the function defining the quantity of interest. Under additional assumptions, one may be able to outperform Monte Carlo with specialized integration rules for integrands that depend on many variables, e.g., with sparse grids [5] or quasi-Monte Carlo [6].

In what follows, we analyze a Monte Carlo method for estimating the eigenpairs that uses independent samples of the function's gradient. After formally defining the active subspace in Section 2, we employ results from Tropp [30] and Gittens and Tropp [14] to bound the probability that the estimated eigenvalues deviate from the true eigenvalues, which yields lower bounds on the number of samples needed for accurate estimation. We extend these results to the case where samples are approximate gradients (e.g., finite difference approximations). In Section 4 we discuss a practical bootstrap approach to study the variability in the estimated eigenvalues, and we demonstrate these procedures numerically in Section 5.

*Notation.* We use bold lower case letters to denote vectors and bold upper case letters to denote matrices. Finite sample estimates are denoted with hats, e.g.,  $\hat{\mathbf{C}} \approx \mathbf{C}$ . The functional  $\lambda_k(\cdot)$  denotes the  $k$ th eigenvalue of its argument, ordered from algebraically largest to smallest; all matrices are symmetric, so the ordering is meaningful. A  $\lambda$  on its own is an eigenvalue. Norms of vectors and matrices are 2-norms; the matrix 2-norm is the operator-induced norm given by the largest singular value. The partial ordering operator  $\preceq$  is defined as follows:  $\mathbf{A} \preceq \mathbf{B}$  means that  $\mathbf{B} - \mathbf{A}$  is positive semidefinite.

**2. Active subspaces.** We represent the map from simulation inputs to the scalar-valued quantity of interest by a function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , where  $\mathcal{X} \subseteq \mathbb{R}^m$ , with  $m > 1$ , represents the set of simulation inputs, which we assume is centered at the origin and scaled so that each component of  $\mathbf{x} \in \mathcal{X}$  has the same range. Let  $\mathbb{R}^m$  be equipped with a weight function  $\rho : \mathbb{R}^m \rightarrow \mathbb{R}_+$  that is bounded, strictly positive on the domain  $\mathcal{X}$ , and zero outside of  $\mathcal{X}$ . We also assume that  $\rho$  is both separable and normalized to integrate to 1. In the context of uncertainty quantification, this weight function represents a given probability density function on the inputs; examples in this context include Gaussian, uniform, or data-conditioned Bayesian posterior density functions. We assume  $f$  is differentiable and absolutely continuous, and we denote the gradient  $\nabla_{\mathbf{x}} f(\mathbf{x}) = [\partial f / \partial x_1, \dots, \partial f / \partial x_m]^T$  oriented as a column vector.

We are interested in the following matrix, denoted  $\mathbf{C}$  and defined as

$$\mathbf{C} = \int (\nabla_{\mathbf{x}} f)(\nabla_{\mathbf{x}} f)^T \rho d\mathbf{x}. \quad (2.1)$$

Samarov studied this matrix as one of several *average derivative functionals* in the context of regression, where  $f$  is the regression function [27]. The matrix  $\mathbf{C}$  is symmetric and positive semi-definite, so it has a real eigenvalue decomposition

$$\mathbf{C} = \mathbf{W}\mathbf{\Lambda}\mathbf{W}^T, \quad \mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_m), \quad \lambda_1 \geq \dots \geq \lambda_m \geq 0, \quad (2.2)$$

where  $\mathbf{W}$  is the orthogonal matrix of eigenvectors. Partition the eigenpairs,

$$\mathbf{W} = [\mathbf{W}_1 \quad \mathbf{W}_2], \quad \mathbf{\Lambda} = \begin{bmatrix} \Lambda_1 & \\ & \Lambda_2 \end{bmatrix}, \quad (2.3)$$

where  $\mathbf{W}_1$  contains the first  $n < m$  eigenvectors, and  $\Lambda_1$  contains the  $n$  largest eigenvalues. The eigenvectors define new coordinates

$$\mathbf{y} = \mathbf{W}_1^T \mathbf{x}, \quad \mathbf{z} = \mathbf{W}_2^T \mathbf{x}. \quad (2.4)$$

We call the column space of  $\mathbf{W}_1$  the *active subspace* and the corresponding  $\mathbf{y}$  the *active variables*. Similarly,  $\mathbf{W}_2$  defines the *inactive subspace* with corresponding *inactive variables*  $\mathbf{z}$ . The following two lemmas justify this characterization; these are proved in [7].

LEMMA 2.1. *The mean-squared directional derivative of  $f$  with respect to the eigenvector  $\mathbf{w}_i$  is equal to the corresponding eigenvalue,*

$$\int ((\nabla_{\mathbf{x}} f)^T \mathbf{w}_i)^2 \rho d\mathbf{x} = \mathbf{w}_i^T \mathbf{C} \mathbf{w}_i = \lambda_i. \quad (2.5)$$

LEMMA 2.2. *The mean-squared gradients of  $f$  with respect to the coordinates  $\mathbf{y}$  and  $\mathbf{z}$  satisfy*

$$\begin{aligned} \int (\nabla_{\mathbf{y}} f)^T (\nabla_{\mathbf{y}} f) \rho d\mathbf{x} &= \text{trace}(\mathbf{W}_1^T \mathbf{C} \mathbf{W}_1) = \lambda_1 + \cdots + \lambda_n, \\ \int (\nabla_{\mathbf{z}} f)^T (\nabla_{\mathbf{z}} f) \rho d\mathbf{x} &= \text{trace}(\mathbf{W}_2^T \mathbf{C} \mathbf{W}_2) = \lambda_{n+1} + \cdots + \lambda_m. \end{aligned} \quad (2.6)$$

The eigenvalues  $\Lambda$  and eigenvectors  $\mathbf{W}$  are properties of  $f$ . If the  $m - n$  trailing eigenvalues  $\Lambda_2$  are exactly zero, then  $f$  is constant along the directions corresponding to  $\mathbf{W}_2$ . If  $\Lambda_2$  is not exactly zero but significantly smaller than  $\Lambda_1$ , then  $f$  changes less, on average, in response to small changes in  $\mathbf{z}$  than small changes in  $\mathbf{y}$ . If  $f$  admits such a property, we would like to discover and exploit it in parameter studies by focusing on the variables  $\mathbf{y}$ . In other words, we can reduce the dimension of the parameter studies from  $m$  to  $n < m$ .

Two special cases illustrate the active subspace. The first class of functions are index models that have the form  $f(\mathbf{x}) = h(\mathbf{A}^T \mathbf{x})$ , where  $\mathbf{A} \in \mathbb{R}^{m \times k}$  and  $h : \mathbb{R}^k \rightarrow \mathbb{R}$ . In this case,  $\mathbf{C}$  has rank at most  $k$ , and the active subspace is a subspace of the range of  $\mathbf{A}$ . If  $k = 1$ , then the one-dimensional active subspace can be discovered with a single evaluation of  $\nabla_{\mathbf{x}} f$  at any  $\mathbf{x} \in \mathcal{X}$  such that the derivative  $h'(\mathbf{A}^T \mathbf{x})$  is not zero. The second special case is a function of the form  $f(\mathbf{x}) = h(\mathbf{x}^T \mathbf{H} \mathbf{x})/2$ , where  $h : \mathbb{R} \rightarrow \mathbb{R}$ , and  $\mathbf{H}$  is a symmetric  $m \times m$  matrix. In this case

$$\mathbf{C} = \mathbf{H} \left( \int (h')^2 \mathbf{x} \mathbf{x}^T \rho d\mathbf{x} \right) \mathbf{H}^T, \quad (2.7)$$

where  $h' = h'(\mathbf{x}^T \mathbf{H} \mathbf{x})$  is the derivative of  $h$ . This implies that the null space of  $\mathbf{C}$  is the null space of  $\mathbf{H}$  provided that  $h'$  is non-degenerate. We study the example where  $h(t) = t$  in Section 5.

If we can estimate  $\Lambda$  and  $\mathbf{W}$  from (2.2), then we can approximate  $f(\mathbf{x})$  with a model of the form

$$f(\mathbf{x}) \approx g(\mathbf{W}_1^T \mathbf{x}), \quad (2.8)$$

where  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  is an appropriately constructed map. In [7], we derive error bounds for this approximation with a particular choice of  $g$ . We extend those error bounds to the case when  $\mathbf{W}_1$  is estimated with some error. The main goal of this paper is to study the error in  $\mathbf{W}_1$  when  $\mathbf{C}$  is estimated with Monte Carlo.

**2.1. Related literature.** The idea of studying the eigenpairs of the average outer product of the gradients arose in statistics as *average derivative functionals* [27, 18] for exploring structure in regression functions. In contrast to our  $f(\mathbf{x})$ , the regression function is generally unknown; to estimate the gradients of the unknown regression function, one can first fit a kernel-based model to a set of predictor/response pairs and then compute gradients from the approximation [31, 12]. In our case, the function is a map between the inputs and outputs of an engineering simulation; there is no random noise as in the regression problem. The set up in Russi’s Ph.D. thesis [25] is closer to ours. He applies the methods to physical simulations of chemical kinetics; this work is where we encountered the term *active subspace*. Recent work in approximation theory by Fornasier, Schnass, and Vybiral [11] attempts to discover the parameters of the active subspace solely through queries of the function; guarantees on reconstruction follow from compressed sensing results under the assumption that  $f$  is an index model.

If the matrix  $\mathbf{C}$  were given as an input matrix, or if we could easily compute matrix-vector products with  $\mathbf{C}$ , then we could employ recent procedures for randomized low-rank approximation to estimate the desired eigenpairs [17, 13]—assuming  $\mathbf{C}$  is well approximated by a low-rank matrix, which is often the case in practice. Unfortunately, we do not have easy access to the elements of  $\mathbf{C}$ ; estimating its eigenpairs requires estimating its elements. There may be fruitful relationships with low-rank approximation of quasimatrices and cmatrices [29] that are worth exploring.

**3. Computing active subspaces.** If drawing independent samples from the density  $\rho$  is cheap and simple, then a straightforward and easy-to-implement random sampling method to approximate the eigenvalues  $\Lambda$  and eigenvectors  $\mathbf{W}$  proceeds as follows.

1. Draw  $N$  samples  $\mathbf{x}_j$  independently from the measure  $\rho$ .
2. For each  $\mathbf{x}_j$ , compute  $\nabla_{\mathbf{x}} f_j = \nabla_{\mathbf{x}} f(\mathbf{x}_j)$ .
3. Approximate

$$\mathbf{C} \approx \hat{\mathbf{C}} = \frac{1}{N} \sum_{j=1}^N (\nabla_{\mathbf{x}} f_j)(\nabla_{\mathbf{x}} f_j)^T. \quad (3.1)$$

4. Compute the eigendecomposition  $\hat{\mathbf{C}} = \hat{\mathbf{W}} \hat{\Lambda} \hat{\mathbf{W}}^T$ .

The last step is equivalent to computing the full SVD of the matrix

$$\hat{\mathbf{B}} = \frac{1}{\sqrt{N}} [\nabla_{\mathbf{x}} f_1 \quad \cdots \quad \nabla_{\mathbf{x}} f_N] = \hat{\mathbf{W}} \hat{\Sigma} \hat{\mathbf{V}}^T, \quad (3.2)$$

where standard manipulations show that the  $\hat{\Lambda} = \hat{\Sigma} \hat{\Sigma}^T$ , and the left singular vectors are the desired eigenvectors. The SVD perspective was developed by Russi [25] as the method to discover the active subspace. This SVD perspective also calls to mind randomized methods for subsampling the columns of  $\hat{\mathbf{B}}$ , where  $N \gg m$  [13]. If it were possible to evaluate the importance of a column of  $\hat{\mathbf{B}}$  without explicitly computing  $\nabla_{\mathbf{x}} f(\mathbf{x})$ , then such ideas might prove useful.

For many simulations, the number  $m$  of input parameters is small enough (e.g., tens to thousands) that computing the full eigendecomposition (3.1) or singular value decomposition (3.2) is negligible compared to the cost of computing the gradient  $N$  times; we consider this to be our case of interest. We are therefore concerned with understanding the number of gradient samples needed so that the estimates  $\hat{\Lambda}$  and  $\hat{\mathbf{W}}$  are close to the true  $\Lambda$  and  $\mathbf{W}$ .

We apply recent work by Tropp [30] and Gittens and Tropp [14] on the spectrum of sums of random matrices to answer these questions. We were motivated to use these tools by Section 7 in Gittens and Tropp [14], which studies the spectrum of a finite sample estimate of a covariance matrix for a Gaussian random vector. In the present case, the gradient vector  $\nabla_{\mathbf{x}}f(\mathbf{x})$  is a deterministic function of  $\mathbf{x}$ . However, if  $\mathbf{x}_j$  are drawn independently at random according to the density  $\rho$ , then we can interpret  $\nabla_{\mathbf{x}}f(\mathbf{x}_j)$  as a random draw from an unknown density. This is a standard interpretation of Monte Carlo techniques for integration [24]. In principle, our analysis approach could apply to model reduction of high-dimensional systems that use Gramian matrices [1].

**THEOREM 3.1.** *Assume that  $\|\nabla_{\mathbf{x}}f\| \leq L$  for all  $\mathbf{x} \in \mathcal{X}$ . Then for  $0 < \varepsilon \leq 1$ ,*

$$\mathbb{P}\left\{\hat{\lambda}_k \geq (1 + \varepsilon)\lambda_k\right\} \leq (m - k + 1) \exp\left(\frac{-N\lambda_k\varepsilon^2}{4L^2}\right), \quad (3.3)$$

and

$$\mathbb{P}\left\{\hat{\lambda}_k \leq (1 - \varepsilon)\lambda_k\right\} \leq k \exp\left(\frac{-N\lambda_k^2\varepsilon^2}{4\lambda_1L^2}\right). \quad (3.4)$$

The key to establishing Theorem 3.1 is a matrix Bernstein inequality from Theorem 5.3 in Gittens and Tropp [14]. When we apply this concentration result, we set

$$\mathbf{X}_j = \nabla_{\mathbf{x}}f_j \nabla_{\mathbf{x}}f_j^T. \quad (3.5)$$

Thus, each  $\mathbf{X}_j$  is an independent random sample of a matrix from the same distribution. Under this notion of randomness,

$$\mathbb{E}[\mathbf{X}_j] = \int \nabla_{\mathbf{x}}f_j \nabla_{\mathbf{x}}f_j^T \rho d\mathbf{x} = \int \nabla_{\mathbf{x}}f \nabla_{\mathbf{x}}f^T \rho d\mathbf{x} = \mathbf{C}. \quad (3.6)$$

For completeness, we restate Theorem 5.3 from [14].

**THEOREM 3.2** (Eigenvalue Bernstein Inequality for Subexponential Matrices, Theorem 5.3 [14]). *Consider a finite sequence  $\{\mathbf{X}_j\}$  of independent, random, self-adjoint matrices with dimension  $n$ , all of which satisfy the subexponential moment growth condition*

$$\mathbb{E}[\mathbf{X}_j^m] \preceq \frac{m!}{2} B^{m-2} \Sigma_j^2 \quad \text{for } m = 2, 3, 4, \dots$$

where  $B$  is a positive constant and  $\Sigma_j^2$  are positive-semidefinite matrices. Given an integer  $k \leq n$ , set

$$\mu_k = \lambda_k \left( \sum_j \mathbb{E}[\mathbf{X}_j] \right).$$

Choose  $\mathbf{V}_+$  as an orthogonal matrix of size  $n \times n - k + 1$  that satisfies

$$\mu_k = \lambda_{\max} \left( \sum_j \mathbf{V}_+^T (\mathbb{E}[\mathbf{X}_j]) \mathbf{V}_+ \right),$$

and define

$$\sigma_k^2 = \lambda_{\max} \left( \sum_j \mathbf{V}_+^T \boldsymbol{\Sigma}_j^2 \mathbf{V}_+ \right).$$

Then, for any  $\tau \geq 0$ ,

$$\mathbb{P} \left\{ \lambda_k \left( \sum_j \mathbf{X}_j \right) \geq \mu_k + \tau \right\} \leq \begin{cases} (n-k+1) \exp(-\tau^2/(4\sigma_k^2)), & \tau \leq \sigma_k^2/B, \\ (n-k+1) \exp(-\tau/(4B)), & \tau \geq \sigma_k^2/B. \end{cases}$$

*Proof.* (Theorem 3.1.) We begin with the upper estimate (3.3). First note that

$$\mathbb{P} \left\{ \lambda_k(\hat{\mathbf{C}}) \geq \lambda_k(\mathbf{C}) + t \right\} = \mathbb{P} \left\{ \lambda_k \left( \sum_{j=1}^N \nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T \right) \geq N\lambda_k + Nt \right\}. \quad (3.7)$$

In this form we can apply Theorem 3.2. We check that the bound on the gradient's norm implies that the matrix  $\nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T$  satisfies the subexponential growth condition:

$$\begin{aligned} \int (\nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T)^p \rho \, d\mathbf{x} &= \int (\nabla_{\mathbf{x}} f^T \nabla_{\mathbf{x}} f)^{p-1} \nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T \rho \, d\mathbf{x} \\ &\leq (L^2)^{p-1} \int \nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T \rho \, d\mathbf{x} \\ &\leq \frac{p!}{2} (L^2)^{p-2} (L^2 \mathbf{C}). \end{aligned} \quad (3.8)$$

Next we set

$$\mu_k = \lambda_k \left( \sum_{j=1}^N \int \nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T \rho \, d\mathbf{x} \right) = N\lambda_k, \quad (3.9)$$

where we simplified using the identically distributed samples of  $\mathbf{x}_j$ . Choose  $\mathbf{W}_+ = \mathbf{W}(:, k:m)$  to be the last  $m-k+1$  eigenvectors of  $\mathbf{C}$ , and note that

$$\lambda_{\max} \left( \sum_{j=1}^N \mathbf{W}_+^T \left( \int \nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T \rho \, d\mathbf{x} \right) \mathbf{W}_+ \right) = N\lambda_{\max}(\mathbf{W}_+^T \mathbf{C} \mathbf{W}_+) = N\lambda_k = \mu_k, \quad (3.10)$$

as required by Theorem 3.2. Define

$$\sigma_k^2 = \lambda_{\max} \left( \sum_{j=1}^N \mathbf{W}_+^T (L^2 \mathbf{C}) \mathbf{W}_+ \right) = NL^2 \lambda_{\max}(\mathbf{W}_+^T \mathbf{C} \mathbf{W}_+) = NL^2 \lambda_k. \quad (3.11)$$

With these quantities, Theorem 3.2 states

$$\mathbb{P} \left\{ \lambda_k \left( \sum_{j=1}^N \nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T \right) \geq N\lambda_k + Nt \right\} \leq (m-k+1) \exp \left( \frac{-(Nt)^2}{4\sigma_k^2} \right) \quad (3.12)$$

when  $Nt \leq \sigma_k^2/L^2$ . Applying this theorem with  $t = \varepsilon\lambda_k$ ,  $\varepsilon \leq 1$ , and the computed  $\sigma_k^2 = NL^2\lambda_k$  yields the upper estimate (3.3).

For the lower estimate,

$$\begin{aligned}
& \mathbb{P} \left\{ \lambda_k(\hat{\mathbf{C}}) \leq \lambda_k(\mathbf{C}) - t \right\} \\
&= \mathbb{P} \left\{ -\lambda_k(\hat{\mathbf{C}}) \geq -\lambda_k(\mathbf{C}) + t \right\} \\
&= \mathbb{P} \left\{ -\lambda_k \left( \sum_{j=1}^N \nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T \right) \geq -N\lambda_k(\mathbf{C}) + Nt \right\} \\
&= \mathbb{P} \left\{ \lambda_{m-k+1} \left( \sum_{j=1}^N (-\nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T) \right) \geq N\lambda_{m-k+1}(-\mathbf{C}) + Nt \right\} \\
&= \mathbb{P} \left\{ \lambda_{k'} \left( \sum_{j=1}^N (-\nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T) \right) \geq N\lambda_{k'}(-\mathbf{C}) + Nt \right\},
\end{aligned} \tag{3.13}$$

for  $k' = m - k + 1$ . We can now apply Theorem 3.2 again. The subexponential growth condition is satisfied since

$$\int (-\nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T)^p \rho \, d\mathbf{x} \preceq \int (\nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T)^p \rho \, d\mathbf{x} \preceq \frac{p!}{2} (L^2)^{p-2} (L^2 \mathbf{C}). \tag{3.14}$$

Set

$$\mu_{k'} = \lambda_{k'} \left( \sum_{j=1}^N \int (-\nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T) \rho \, d\mathbf{x} \right) = N\lambda_{k'}(-\mathbf{C}). \tag{3.15}$$

Set  $\mathbf{W}_+ = \mathbf{W}(:, 1:k)$  to be the first  $k$  eigenvectors of  $\mathbf{C}$ , and note that

$$\begin{aligned}
\lambda_{\max} \left( \sum_{j=1}^N \mathbf{W}_+^T \left( \int (-\nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T) \rho \, d\mathbf{x} \right) \mathbf{W}_+ \right) &= N\lambda_{\max}(-\mathbf{W}_+^T \mathbf{C} \mathbf{W}_+) \\
&= N(-\lambda_k(\mathbf{C})) \\
&= N\lambda_{m-k+1}(-\mathbf{C}) \\
&= N\lambda_{k'}(-\mathbf{C}),
\end{aligned} \tag{3.16}$$

as required by Theorem 3.2. Set

$$\sigma_{k'}^2 = \lambda_{\max} \left( \sum_{j=1}^N \mathbf{W}_+^T (L^2 \mathbf{C}) \mathbf{W}_+ \right) = NL^2 \lambda_{\max}(\mathbf{W}_+^T \mathbf{C} \mathbf{W}_+) = NL^2 \lambda_1. \tag{3.17}$$

Theorem 3.2 states

$$\mathbb{P} \left\{ \lambda_{k'} \left( \sum_{j=1}^N (-\nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T) \right) \geq N\lambda_{k'}(-\mathbf{C}) + Nt \right\} \leq k \exp \left( \frac{-(Nt)^2}{4\sigma_{k'}^2} \right) \tag{3.18}$$

when  $Nt \leq \sigma_{k'}^2 / L^2$ . Plug in the computed quantities with  $t = -\varepsilon \lambda_{k'}(-\mathbf{C}) = \varepsilon \lambda_k(\mathbf{C})$  to achieve the lower estimate (3.4). Note that the condition  $\varepsilon \leq 1 \leq \lambda_1 / \lambda_{k'}$  allows us to apply Theorem 3.2.  $\square$

Next we use this result to derive a lower bound on the number of gradient samples needed for relative accuracy of  $\varepsilon$ . Recall the definition of *big omega* notation that  $a = \Omega(b)$  means  $a \geq cb$  for some positive constant  $c$ .

COROLLARY 3.3. *Let  $\kappa_k = \lambda_1/\lambda_k$ . Then for  $\varepsilon \in (0, 1]$ ,*

$$N = \Omega\left(\frac{L^2 \kappa_k^2}{\lambda_1 \varepsilon^2} \log(m)\right) \quad (3.19)$$

*implies  $|\hat{\lambda}_k - \lambda_k| \leq \varepsilon \lambda_k$  with high probability.*

*Proof.* Starting with the upper estimate from Theorem 3.1, if

$$N \geq \frac{4L^2}{\lambda_k \varepsilon^2} (\beta + 1) \log(m) \geq \frac{4L^2}{\lambda_k \varepsilon^2} (\beta \log(m) + \log(m - k + 1)), \quad (3.20)$$

then

$$\mathbb{P}\left\{\hat{\lambda}_k \geq (1 + \varepsilon)\lambda_k\right\} \leq m^{-\beta}. \quad (3.21)$$

Similarly for the lower estimate from Theorem 3.1, if

$$N \geq \frac{4L^2 \lambda_1}{\lambda_k^2 \varepsilon^2} (\beta + 1) \log(m) \geq \frac{4L^2 \lambda_1}{\lambda_k^2 \varepsilon^2} (\beta \log(m) + \log(k)), \quad (3.22)$$

then

$$\mathbb{P}\left\{\hat{\lambda}_k \leq (1 - \varepsilon)\lambda_k\right\} \leq m^{-\beta}. \quad (3.23)$$

Setting  $\kappa_k = \lambda_1/\lambda_k$  and taking

$$N \geq (\beta + 1) \frac{4L^2 \kappa_k^2}{\lambda_1 \varepsilon^2} \log(m) \quad (3.24)$$

satisfies both conditions.  $\square$

We can combine results from Golub and Van Loan [15, Chapter 8] with results from Tropp [30] to obtain an estimate of the distance between the subspace defined by the eigenvectors  $\mathbf{W}_1$  and the subspace defined by the eigenvectors  $\hat{\mathbf{W}}_1$ . This requires a different matrix Bernstein inequality in the form of Theorem 6.1 from Tropp [30], which we restate below. When we apply the theorem,  $\mathbf{X}_j = \nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T - \mathbf{C}$ , that is, the random matrix samples are the deviation of the  $j$ th sampled gradient outer product from the matrix  $\mathbf{C}$ .

THEOREM 3.4 (Matrix Bernstein: bounded case, Theorem 6.1 [30]). *Consider a finite sequence  $\{\mathbf{X}_j\}$  of independent, random, self-adjoint matrices with dimension  $n$ . Assume that*

$$\mathbb{E}[\mathbf{X}_j] = 0 \quad \text{and} \quad \lambda_{\max}(\mathbf{X}_j) \leq R \quad \text{almost surely.}$$

*Compute the norm of the total variance,*

$$\sigma^2 := \left\| \sum_j \mathbb{E}[\mathbf{X}_j^2] \right\|.$$



Then the following inequality holds for all  $\tau \geq 0$ :

$$\mathbb{P} \left\{ \lambda_{\max} \left( \sum_j \mathbf{X}_j \right) \geq \tau \right\} \leq \begin{cases} n \exp(-3\tau^2/(8\sigma^2)), & \tau \leq \sigma^2/R, \\ n \exp(-3\tau/(8R)), & \tau > \sigma^2/R. \end{cases}$$

**THEOREM 3.5.** Let  $\varepsilon > 0$ . Assume  $\|\nabla_{\mathbf{x}} f\| \leq L$  for all  $\mathbf{x} \in \mathcal{X}$ . Define the variance

$$\nu^2 = \left\| \int (\nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T - \mathbf{C})^2 \rho d\mathbf{x} \right\|, \quad (3.25)$$

and assume  $\nu^2 > 0$ . Then

$$\mathbb{P} \left\{ \|\hat{\mathbf{C}} - \mathbf{C}\| \geq \varepsilon \|\mathbf{C}\| \right\} \leq \begin{cases} 2m \exp\left(\frac{-3N\lambda_1^2\varepsilon^2}{8\nu^2}\right), & \text{if } \varepsilon \leq \nu^2/(\lambda_1 L^2), \\ 2m \exp\left(\frac{-3N\lambda_1\varepsilon}{8L^2}\right), & \text{if } \varepsilon > \nu^2/(\lambda_1 L^2). \end{cases} \quad (3.26)$$

*Proof.* Observe that

$$\begin{aligned} \mathbb{P} \left\{ \|\hat{\mathbf{C}} - \mathbf{C}\| \geq t \right\} &= \mathbb{P} \left\{ \lambda_{\max}(\hat{\mathbf{C}} - \mathbf{C}) \geq t \text{ or } \lambda_{\max}(\mathbf{C} - \hat{\mathbf{C}}) \geq t \right\} \\ &\leq \mathbb{P} \left\{ \lambda_{\max}(\hat{\mathbf{C}} - \mathbf{C}) \geq t \right\} + \mathbb{P} \left\{ \lambda_{\max}(\mathbf{C} - \hat{\mathbf{C}}) \geq t \right\} \\ &= \mathbb{P} \left\{ \lambda_{\max} \left( \sum_{j=1}^N (\nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T - \mathbf{C}) \right) \geq Nt \right\} \\ &\quad + \mathbb{P} \left\{ \lambda_{\max} \left( \sum_{j=1}^N (\mathbf{C} - \nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T) \right) \geq Nt \right\}. \end{aligned} \quad (3.27)$$

Note that both

$$\int (\nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T - \mathbf{C}) \rho d\mathbf{x} = \int (\mathbf{C} - \nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T) \rho d\mathbf{x} = \mathbf{0}. \quad (3.28)$$

Since  $\mathbf{C}$  is positive semidefinite and  $\|\nabla_{\mathbf{x}} f\| \leq L$ ,

$$\begin{aligned} \lambda_{\max}(\nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T - \mathbf{C}) &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^T (\nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T - \mathbf{C}) \mathbf{v} \\ &\leq \max_{\|\mathbf{v}\|=1} \mathbf{v}^T (\nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T) \mathbf{v} \leq L^2. \end{aligned} \quad (3.29)$$

This bound also holds for  $\lambda_{\max}(\mathbf{C} - \nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T)$ , since

$$\begin{aligned} \lambda_{\max}(\mathbf{C} - \nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T) &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^T (\mathbf{C} - \nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T) \mathbf{v} \\ &\leq \max_{\|\mathbf{v}\|=1} \mathbf{v}^T \mathbf{C} \mathbf{v} \leq \|\mathbf{C}\| \leq L^2. \end{aligned} \quad (3.30)$$

Thus, the upper-bound  $R$  in Theorem 3.4 is  $L^2$ . The variance parameter  $\sigma^2$  is

$$\sigma^2 = \left\| \left( \sum_{j=1}^N \int (\nabla_{\mathbf{x}} f_j \nabla_{\mathbf{x}} f_j^T - \mathbf{C})^2 \rho d\mathbf{x} \right) \right\| = N \left\| \int (\nabla_{\mathbf{x}} f \nabla_{\mathbf{x}} f^T - \mathbf{C})^2 \rho d\mathbf{x} \right\| = N \nu^2. \quad (3.31)$$

Assume  $\varepsilon \leq \nu^2/(\lambda_1 L^2)$ . Then  $N\lambda_1\varepsilon \leq N\nu^2/L^2$ , and we can apply the upper branch of Theorem 3.4 to the two terms at the end of (3.27) with  $\tau = \lambda_1\varepsilon = \|\mathbf{C}\|_\varepsilon$ , which produces the bound's upper branch in (3.26). Now assume  $\varepsilon > \nu^2/(\lambda_1 L^2)$ . Similarly,  $N\lambda_1\varepsilon \geq N\nu^2/L^2$ , and we can apply the lower branch of Theorem 3.4 with the same  $t = \lambda_1\varepsilon$  to (3.27) to produce the lower branch of (3.26).  $\square$

This result leads to a lower bound on the number of samples needed for a small relative error in  $\hat{\mathbf{C}}$  in the matrix 2-norm; compare the following corollary to Corollary 3.3.

COROLLARY 3.6. *Let  $\varepsilon > 0$ , and define*

$$\delta = \max\left(\frac{\nu^2}{\lambda_1\varepsilon}, L^2\right). \quad (3.32)$$

Then

$$N = \Omega\left(\frac{\delta}{\lambda_1\varepsilon} \log(2m)\right) \quad (3.33)$$

implies that  $\|\hat{\mathbf{C}} - \mathbf{C}\| \leq \varepsilon\|\mathbf{C}\|$  with high probability.

*Proof.* We consider the two cases of  $\varepsilon$  from Theorem 3.5. Assume  $\varepsilon \leq \nu^2/(\lambda_1 L^2)$ , so that  $N\lambda_1\varepsilon \leq N\nu^2/L^2$ , we follow the reasoning in the proof of Corollary 3.3 with the upper branch of the bound in Theorem 3.5 to get

$$N \geq \frac{8}{3}(\beta + 1)\frac{\nu^2}{\lambda_1^2\varepsilon^2} \log(2m). \quad (3.34)$$

Similarly, if  $\varepsilon > \nu^2/(\lambda_1 L^2)$ , then the lower branch from the bound in Theorem 3.5 produces

$$N \geq \frac{8}{3}(\beta + 1)\frac{L^2}{\lambda_1\varepsilon} \log(2m). \quad (3.35)$$

Using  $\delta$  from (3.32) chooses the larger lower bound between (3.34) and (3.35).  $\square$

We can combine Corollary 3.6 with [15, Corollary 8.1.11] to control the error in the estimated subspace defined by  $\hat{\mathbf{W}}_1$ . We quantify this error by the distance between the subspace defined by  $\mathbf{W}_1$  and the subspace defined by  $\hat{\mathbf{W}}_1$ . Recall the definition of the distance between subspaces [28],

$$\text{dist}(\text{ran}(\mathbf{W}_1), \text{ran}(\hat{\mathbf{W}}_1)) = \|\mathbf{W}_1\mathbf{W}_1^T - \hat{\mathbf{W}}_1\hat{\mathbf{W}}_1^T\| = \|\mathbf{W}_1^T\hat{\mathbf{W}}_2\|. \quad (3.36)$$

COROLLARY 3.7. *Let  $\varepsilon > 0$  be such that*

$$\varepsilon \leq \min(1, (\lambda_n - \lambda_{n+1})/(5\lambda_1)), \quad (3.37)$$

and choose  $N$  according to Corollary 3.6. Then

$$\text{dist}(\text{ran}(\mathbf{W}_1), \text{ran}(\hat{\mathbf{W}}_1)) \leq \frac{4\lambda_1\varepsilon}{\lambda_n - \lambda_{n+1}}, \quad (3.38)$$

with high probability.

*Proof.* Let  $\mathbf{E} = \hat{\mathbf{C}} - \mathbf{C}$ . For  $\varepsilon$  in (3.37) with  $N$  chosen according to Corollary 3.6, we have

$$\|\mathbf{E}\| \leq \varepsilon\|\mathbf{C}\| = \varepsilon\lambda_1 \leq (\lambda_n - \lambda_{n+1})/5, \quad (3.39)$$

with high probability. Under this condition on  $\|\mathbf{E}\|$ , [15, Corollary 8.1.11] states

$$\text{dist}(\text{ran}(\mathbf{W}_1), \text{ran}(\hat{\mathbf{W}}_1)) \leq \frac{4\|\mathbf{E}\|}{\lambda_n - \lambda_{n+1}} \leq \frac{4\lambda_1\varepsilon}{\lambda_n - \lambda_{n+1}}, \quad (3.40)$$

as required.  $\square$

Corollary 3.7 shows that control of the eigenvalues implies control of the subspace generated by the eigenvectors. However, the error in the estimated subspace is inversely proportional to the gap between the smallest eigenvalue associated with the active subspace and the largest eigenvalue associated with the inactive subspace. This implies, for example, if the gap between the second and third eigenvalues is larger than the gap between the first and second, then estimates of a two-dimensional active subspace are more accurate than estimates of a one-dimensional active subspace.

**3.1. Approximate gradients.** Many modern simulations have subroutines for estimating gradients with, e.g., adjoint methods [4, 3] or algorithmic differentiation [16]. However, legacy codes or simulations that couple multiple codes might not have such gradient capabilities. When there is no subroutine for gradients, finite difference approximations may suffice when  $m$  is not too large and  $f$  is neither too expensive nor too noisy. Recent work characterizes the gradient when function evaluations contain noise [22, 23].

Next, we extend the bounds on errors in the estimated eigenpairs to the case when the gradients are computed with some error. The gradient error model we analyze depends on a parameter that controls the amount of error. Let  $\mathbf{g}(\mathbf{x})$  denote the approximate gradient computed at  $\mathbf{x} \in \mathcal{X}$ . We assume that

$$\|\mathbf{g}(\mathbf{x}) - \nabla_{\mathbf{x}}f(\mathbf{x})\| \leq \sqrt{m}\gamma_h, \quad \mathbf{x} \in \mathcal{X}, \quad (3.41)$$

where

$$\lim_{h \rightarrow 0} \gamma_h = 0. \quad (3.42)$$

The parameter  $h$  may be a finite difference parameter, the grid spacing in a continuous adjoint computation, or the solver tolerance for a discrete adjoint computation.

Define the symmetric positive semidefinite matrix  $\mathbf{G}$  and its eigenvalue decomposition

$$\mathbf{G} = \int \mathbf{g} \mathbf{g}^T \rho d\mathbf{x} = \mathbf{U}\Theta\mathbf{U}^T, \quad \Theta = \text{diag}(\theta_1, \dots, \theta_m), \quad (3.43)$$

and define its random sample approximation

$$\hat{\mathbf{G}} = \frac{1}{N} \sum_{j=1}^N \mathbf{g}_j \mathbf{g}_j^T = \hat{\mathbf{U}}\hat{\Theta}\hat{\mathbf{U}}^T, \quad \hat{\Theta} = \text{diag}(\hat{\theta}_1, \dots, \hat{\theta}_m), \quad (3.44)$$

where  $\mathbf{g}_j = \mathbf{g}(\mathbf{x}_j)$  for  $\mathbf{x}_j$  drawn independently from  $\rho$ . With these quantities defined, we have the following lemma.

LEMMA 3.8. *Let  $\|\nabla_{\mathbf{x}}f\| \leq L$  for all  $\mathbf{x} \in \mathcal{X}$ . The norm of the difference between  $\hat{\mathbf{C}}$  and  $\hat{\mathbf{G}}$  is bounded by*

$$\|\hat{\mathbf{C}} - \hat{\mathbf{G}}\| \leq (\sqrt{m}\gamma_h + 2L)\sqrt{m}\gamma_h. \quad (3.45)$$

*Proof.* Let  $\mathbf{g} = \mathbf{g}(\mathbf{x})$  and  $\nabla_{\mathbf{x}}f = \nabla_{\mathbf{x}}f(\mathbf{x})$ . First observe

$$\|\mathbf{g} + \nabla_{\mathbf{x}}f\| = \|\mathbf{g} - \nabla_{\mathbf{x}}f + 2\nabla_{\mathbf{x}}f\| \leq \|\mathbf{g} - \nabla_{\mathbf{x}}f\| + 2\|\nabla_{\mathbf{x}}f\| \leq \sqrt{m}\gamma_h + 2L. \quad (3.46)$$

Next,

$$\begin{aligned} \|\mathbf{g}\mathbf{g}^T - \nabla_{\mathbf{x}}f\nabla_{\mathbf{x}}f^T\| &= \frac{1}{2}\|(\mathbf{g} + \nabla_{\mathbf{x}}f)(\mathbf{g} - \nabla_{\mathbf{x}}f)^T + (\mathbf{g} - \nabla_{\mathbf{x}}f)(\mathbf{g} + \nabla_{\mathbf{x}}f)^T\| \\ &\leq \|(\mathbf{g} + \nabla_{\mathbf{x}}f)(\mathbf{g} - \nabla_{\mathbf{x}}f)^T\| \\ &\leq (\sqrt{m}\gamma_h + 2L)\sqrt{m}\gamma_h. \end{aligned} \quad (3.47)$$

Then,

$$\begin{aligned} \|\hat{\mathbf{G}} - \hat{\mathbf{C}}\| &= \left\| \frac{1}{N} \sum_{j=1}^N \mathbf{g}_j \mathbf{g}_j^T - \frac{1}{N} \sum_{j=1}^N \nabla_{\mathbf{x}}f_j \nabla_{\mathbf{x}}f_j^T \right\| \\ &\leq \frac{1}{N} \sum_{j=1}^N \|\mathbf{g}_j \mathbf{g}_j^T - \nabla_{\mathbf{x}}f_j \nabla_{\mathbf{x}}f_j^T\| \\ &\leq \sqrt{m}\gamma_h(\sqrt{m}\gamma_h + 2L). \end{aligned} \quad (3.48)$$

□

We combine Lemma 3.8 with Corollary 3.3 to study the error in the eigenvalues of the random sample estimate with approximate gradients.

**THEOREM 3.9.** *For  $\varepsilon \in (0, 1]$ , if  $N$  is chosen as (3.19), then the difference between  $\lambda_k$  in (2.2) and the eigenvalue  $\hat{\theta}_k$  from (3.44) is bounded as*

$$|\lambda_k - \hat{\theta}_k| \leq \varepsilon\lambda_k + \sqrt{m}\gamma_h(\sqrt{m}\gamma_h + 2L), \quad (3.49)$$

with high probability.

*Proof.* Observe that

$$|\lambda_k - \hat{\theta}_k| \leq |\lambda_k - \hat{\lambda}_k| + |\hat{\lambda}_k - \hat{\theta}_k|. \quad (3.50)$$

Apply Corollary 3.3 to the first term. The second term follows from [15, Corollary 8.1.6] combined with Lemma 3.8, since

$$|\hat{\theta}_k - \hat{\lambda}_k| = |\lambda_k(\hat{\mathbf{G}}) - \lambda_k(\hat{\mathbf{C}})| \leq \|\hat{\mathbf{G}} - \hat{\mathbf{C}}\| \leq \sqrt{m}\gamma_h(\sqrt{m}\gamma_h + 2L). \quad (3.51)$$

□

The bias in the finite sample eigenvalue estimates using approximate gradients goes to zero at the same rate as the error in the approximate gradient. Next we attend to the error in the active subspace computed with Monte Carlo and approximate gradients.

**THEOREM 3.10.** *Choose  $\varepsilon > 0$  such that*

$$\varepsilon < \min\left(1, \frac{\lambda_n - \lambda_{n+1}}{5\lambda_1}, \frac{\lambda_n - \lambda_{n+1}}{\lambda_n + \lambda_{n+1}}\right) \quad (3.52)$$

*Choose  $N$  so that it satisfies both (3.33) and (3.19) with  $k = n + 1$ . Choose  $h$  small enough so that*

$$\sqrt{m}\gamma_h(\sqrt{m}\gamma_h + 2L) \leq \frac{(1 - \varepsilon)\lambda_n - (1 + \varepsilon)\lambda_{k+1}}{5}. \quad (3.53)$$

Then

$$\text{dist}(\text{ran}(\hat{\mathbf{U}}_1), \text{ran}(\mathbf{W}_1)) \leq \frac{4\sqrt{m}\gamma_h(\sqrt{m}\gamma_h + 2L)}{(1-\varepsilon)\lambda_n - (1+\varepsilon)\lambda_{n+1}} + \frac{4\lambda_1}{\lambda_n - \lambda_{n+1}}, \quad (3.54)$$

with high probability.

*Proof.* The conditions on  $N$  and  $\varepsilon$  imply  $|\hat{\lambda}_{n+1} - \lambda_{n+1}| \leq \varepsilon\lambda_{n+1}$  with high probability due to Corollary 3.3. Examining (3.19), we see that if  $N$  is large enough to estimate  $\lambda_{n+1}$ , then  $N$  is large enough to estimate  $\lambda_n$ , so  $|\hat{\lambda}_n - \lambda_n| \leq \varepsilon\lambda_n$  with high probability, too. Then

$$\begin{aligned} \lambda_n - \lambda_{n+1} &= |\lambda_n - \lambda_{n+1}| \\ &\leq |\lambda_n - \hat{\lambda}_n| + |\hat{\lambda}_{n+1} - \lambda_{n+1}| + (\hat{\lambda}_n - \hat{\lambda}_{n+1}) \\ &\leq \varepsilon\lambda_n + \varepsilon\lambda_{n+1} + (\hat{\lambda}_n - \hat{\lambda}_{n+1}), \end{aligned} \quad (3.55)$$

with high probability. Rearranging this inequality yields

$$\hat{\lambda}_n - \hat{\lambda}_{n+1} \geq (1-\varepsilon)\lambda_n - (1+\varepsilon)\lambda_{n+1}. \quad (3.56)$$

This relates the gap between the eigenvalue estimates to the gap between the true eigenvalues. The condition on  $\varepsilon$  ensures that

$$(1-\varepsilon)\lambda_n - (1+\varepsilon)\lambda_{n+1} > 0. \quad (3.57)$$

Next,

$$\text{dist}(\text{ran}(\hat{\mathbf{U}}_1), \text{ran}(\mathbf{W}_1)) \leq \text{dist}(\text{ran}(\hat{\mathbf{U}}_1), \text{ran}(\hat{\mathbf{W}}_1)) + \text{dist}(\text{ran}(\hat{\mathbf{W}}_1), \text{ran}(\mathbf{W}_1)). \quad (3.58)$$

The second term on the right is bounded in Corollary 3.7 under the assumptions on  $N$  and  $\varepsilon$ . The condition (3.53) on  $h$  and (3.56) imply

$$\sqrt{m}\gamma_h(\sqrt{m}\gamma_h + 2L) \leq \frac{\hat{\lambda}_n - \hat{\lambda}_{k+1}}{5}. \quad (3.59)$$

Then [15, Corollary 8.1.11] implies

$$\text{dist}(\text{ran}(\hat{\mathbf{U}}_1), \text{ran}(\hat{\mathbf{W}}_1)) \leq \frac{4}{\hat{\lambda}_n - \hat{\lambda}_{n+1}} \|\hat{\mathbf{G}} - \hat{\mathbf{C}}\|. \quad (3.60)$$

Combining this with (3.56) and the bound from Lemma 3.8 yields the result.  $\square$

In summary, the eigenvalues and the active subspace approximated with Monte Carlo and approximate gradients are well-behaved. The error bounds include a term that goes to zero like the error in the approximate gradient and a term that behaves like the finite sample approximation with exact gradients. Note that the error bound on the subspace estimate depends on both the gap between  $\lambda_n$  and  $\lambda_{n+1}$  and a smaller gap that depends on  $\varepsilon$ .

**4. Practical approach to computation.** The bounds we present in Section 3 provide a theoretical foundation for understanding the behavior of the Monte Carlo estimates. However, many of the quantities in the bounds may not be known a priori—such as the maximum norm of the gradient  $L$  and the true eigenvalues of the matrix  $\mathbf{C}$ . In this section we offer a practical recipe guided by the insights from the

theory. We caution that the following approach, which relies on a nonparametric bootstrap, can perform poorly for badly-behaved functions. For example, one could be unlucky and sample the gradient in regions that are not representative of the gradient over the entire domain; the bootstrap uses only the  $N$  samples used to compute the eigenpair estimates. Also, errors in the gradients could produce poor approximations of the eigenvalues and subspaces; we show an example of this in Section 5.1. Nevertheless, we have used the following approach on several problems in practice to reveal low-dimensional structure in complex functions of several variables coming from engineering simulations [7, 8, 21].

The first objective is to estimate the eigenvalues and a measure of the estimates' variability from the finite samples. Suppose one wishes to estimate the first  $k$  eigenpairs from the matrix  $\mathbf{C}$ . Practical considerations guide the choice of  $k$ . For example, if one wishes to build a response surface approximation of  $f$  on a low-dimensional domain, then five or six dimensions might be the most one can afford given the cost of computing  $f(\mathbf{x})$ . Hence  $k$  might be seven or eight to allow the possibility of finding a gap that indicates a sufficiently low-dimensional approximation. If a gap is not present in the first  $k$  eigenvalues, then  $f$  may not be amenable to dimension reduction via active subspaces for the desired purpose.

We recommend choosing the number  $N$  of independent gradient samples as

$$N = \alpha k \log(m), \quad (4.1)$$

where  $\alpha$  is a multiplier between 2 and 10. Taking at least  $k$  samples means that  $\mathbf{C}$  is a sum of  $k$  rank-one matrices and thus has a rank of at most  $k$ . This allows the possibility of estimating  $k$  non-zero eigenvalues. The  $\log(m)$  term follows from the bounds in Theorem 3.1. The  $\alpha$  between 2 and 10 is an ad hoc multiplicative factor that we have used on several problems. In principle,  $\alpha k$  is meant to model the contribution from the unknown terms  $L$ ,  $\kappa_k$ ,  $\nu^2$ , and  $\lambda_1$  in (3.19) and (3.33). It is likely that the combination of these terms with the  $\varepsilon^{-1}$  is greater than  $10k$  for small  $\varepsilon$ . However, the Bernstein inequalities used to derive the lower bounds on  $N$  in Corollaries 3.3 and 3.6 are also conservative. One can also assess if  $N$  is large enough a posteriori by examining the bootstrap intervals described below.

We form  $\hat{\mathbf{C}}$  using the samples of the gradient as in (3.1), and then compute its eigenvalue decomposition. We expect that computing the full eigendecomposition is much cheaper than computing the gradient samples. A function of a thousand variables produces  $\hat{\mathbf{C}}$  with dimension thousand-by-thousand. Full eigendecompositions for matrices this size are computed in seconds on modern laptops.

We suggest computing bootstrap intervals for the eigenvalues, which involves computing the eigendecompositions of several matrices the size of  $\hat{\mathbf{C}}$ . The bootstrap creates replicates by (i) sampling with replacement from the set of gradient samples, (ii) computing the replicate  $\hat{\mathbf{C}}^*$ , and (iii) computing its eigenvalue decomposition. The collection of eigenvalue replicates is used to estimate bounds on the true eigenvalues. Efron and Tibshirani use the bootstrap to get empirical density functions of estimated eigenvalues from a covariance matrix in section 7.2 of their book [10]. Chapter 3 of Jolliffe's book [19] also comments on the bootstrap approach for estimating eigenvalues and eigenvectors of a covariance matrix from independent samples. The bootstrap estimates of the standard error and confidence intervals for the eigenvalues may be biased, but this bias decreases as the number  $N$  of samples increases. Since these estimates may be biased, we refer to them as *bootstrap intervals* instead of *confidence intervals*.

Corollary 3.7 says that the error in the estimated subspace depends inversely on the gap between the eigenvalues scaled by the largest eigenvalue. The key to accurately approximating the subspace is to look for gaps in the eigenvalues; this is consistent with standard perturbation theory for eigenvector computations [28]. For example, if there is a larger gap between the third and fourth eigenvalues than between the second and third, then estimates of the three-dimensional subspace are more accurate than estimates of the two-dimensional subspace. This contrasts with heuristics for choosing the dimension of the subspace in (i) model reduction based on the proper orthogonal decomposition [2] and (ii) dimension reduction based on principal component analysis [19]. In these cases, one chooses the dimension of the subspace by a threshold on the magnitude of the eigenvalues—e.g., so that the sum of retained eigenvalues exceeds some proportion of the sum of all eigenvalues. To accurately approximate the active subspace, the most important quantity is the spectral gap, which indicates a separation between scales. To tease out the spectral gap, plot the estimated eigenvalues and their respective upper and lower bootstrap intervals; a gap between subsequent intervals offers confidence of a spectral gap and, hence, the presence of an active subspace. In Section 5, we show several examples of such plots (Figures 5.1, 5.2, 5.4, 5.5).

One should also consider the intent of the dimension reduction when choosing the dimension of the active subspace. For example, if the goal is to approximate a function of  $m$  variables by a surrogate model of  $n$  variables—as in [7]—then one may be limited to  $n$ 's small enough to permit surrogate construction. Suppose the largest  $n$  one is willing to use is  $n = n_{\max} = 5$ , but there is no gap between consecutive eigenvalues from  $\lambda_1$  to  $\lambda_6$ . Then subspace-based dimension reduction may not be an appropriate tool, and one should consider searching for other types of exploitable structure in the model.

Assuming we have chosen  $n$ , we wish to study the variability in the active subspace due to finite sampling; we again turn to the bootstrap. In particular, for each replicate  $\hat{\mathbf{W}}^*$  of the eigenvectors, we compute  $\text{dist}(\text{ran}(\hat{\mathbf{W}}_1), \text{ran}(\hat{\mathbf{W}}_1^*))$ . One can examine the bootstrap intervals of this quantity to study the stability of the subspace. Recall that the distance between subspaces is bounded above by 1, so a bootstrap interval whose values are close to 1 indicates a poorly approximated active subspace. Figures 5.1, 5.3, 5.4, and 5.5 show examples of plotting this metric for the stability of the subspace; the first two Figures also compare the measure of stability to the true error in the active subspace.

**4.1. A step-by-step procedure.** We summarize the practical approach to approximating the active subspace with bootstrap intervals. What follows is a modification of the procedure outlined at the beginning of Section 3 including our suggestions for parameter values. This procedure assumes the user has decided on the number  $k$  of eigenvalues to examine.

1. Choose  $N = \alpha k \log(m)$ , where  $\alpha$  is a multiplier between 2 and 10, and choose  $N_{\text{boot}}$  between 100 and 10000.
2. Draw  $N$  samples  $\{\mathbf{x}_j\}$  independently from  $\rho$ . For each  $\mathbf{x}_j$ , compute  $\nabla_{\mathbf{x}} f_j = \nabla_{\mathbf{x}} f(\mathbf{x}_j)$ .
3. Compute

$$\hat{C} = \frac{1}{N} \sum_{j=1}^N (\nabla_{\mathbf{x}} f_j)(\nabla_{\mathbf{x}} f_j)^T = \hat{\mathbf{W}} \hat{\Lambda} \hat{\mathbf{W}}^T. \quad (4.2)$$

4. **Bootstrap:** For  $i$  from 1 to  $N_{\text{boot}}$ , let  $\ell_1^i, \dots, \ell_N^i$  be  $N$  integers drawn randomly from  $\{1, \dots, N\}$  with replacement, and compute

$$\hat{\mathbf{C}}_i^* = \frac{1}{N} \sum_{j=1}^N (\nabla_{\mathbf{x}} f_{\ell_j^i}) (\nabla_{\mathbf{x}} f_{\ell_j^i})^T = \hat{\mathbf{W}}_i^* \hat{\Lambda}_i^* (\hat{\mathbf{W}}_i^*)^T. \quad (4.3)$$

The asterisk denotes a bootstrap replicate. Then compute the subspace distance

$$d_i^* = \text{dist}(\text{ran}(\hat{\mathbf{W}}), \text{ran}(\hat{\mathbf{W}}_i^*)). \quad (4.4)$$

5. Compute the intervals

$$\left[ \min_i \hat{\lambda}_{j,i}^*, \max_i \hat{\lambda}_{j,i}^* \right], \quad j = 1, \dots, k \quad (4.5)$$

where  $\hat{\lambda}_{j,i}^*$  is the  $j$ th diagonal from  $\hat{\Lambda}_i^*$  in (4.3). Also compute the mean, minimum, and maximum from the set  $\{d_i^*\}$  to estimate the subspace error.

6. Plot the eigenvalue bootstrap intervals and look for large gaps. Choose the dimension  $n$  of the active subspace corresponding to the largest eigenvalue gap. If there is no perceivable gap, then an active subspace may not be present in the first  $k - 1$  dimensions.

A few comments are in order. First, we assume the dimension  $m$  of  $\hat{\mathbf{C}}$  is small enough so that the eigendecompositions of  $\hat{\mathbf{C}}$  and its bootstrap replicates are much cheaper than the samples of the gradient. Such is the case when  $m$  is in the thousands (i.e.,  $f$  depends on thousands of input variables), and  $f$  and  $\nabla_{\mathbf{x}} f$  come from an expensive engineering simulation. Second, we choose the bootstrap to examine the variability because we assume that sampling more gradients is not feasible. If this is not the case, i.e., if one can cheaply evaluate many more gradient samples, then one can compute Monte Carlo estimates and central limit theorem confidence intervals of the eigenvalues in place of the bootstrap estimates. Lastly, we note that the elements of  $\mathbf{C}$  are multivariate integrals. If  $m$  is small enough (2 or 3) and evaluating  $\nabla_{\mathbf{x}} f$  is cheap enough, then more accurate numerical quadrature rules may perform better than the random sampling, i.e., greater accuracy for fewer samples. However, practical error estimates are more difficult to compute, since the error is due to bias instead of variance.

**5. Experiments.** We apply the procedures described in Section 4 to two models: (i) a quadratic function and (ii) a linear functional of the solution of a parameterized PDE. The quadratic model is simple enough to analytically derive the eigenpairs of the active subspace for thorough evaluation of the method. We study the same PDE model in [7, Section 5]. Gradients are available through adjoints, but the true active subspaces are not available. We support efforts for reproducible research [20, 9]; codes for the experiments in this section can be found at

<https://www.cs.purdue.edu/homes/dgleich/codes/compute-asm/compute-asm-code.tar.gz>

The PDE example uses the Random Field Simulation code (<http://www.mathworks.com/matlabcentral/fileexchange/276>) as well as the MATLAB PDE Toolbox.

- 5.1. A quadratic model.** Consider a quadratic function of  $m = 10$  variables,

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x}, \quad \mathbf{x} \in [-1, 1]^{10}, \quad (5.1)$$



where  $\mathbf{A}$  is symmetric and positive definite. We take  $\rho = 2^{-10}$  on the hypercube  $[-1, 1]^{10}$  and zero elsewhere. The gradient is  $\nabla_{\mathbf{x}}f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ , so

$$\mathbf{C} = \mathbf{A} \left( \int \mathbf{x} \mathbf{x}^T \rho d\mathbf{x} \right) \mathbf{A}^T = \frac{1}{3} \mathbf{A}^2. \quad (5.2)$$

The eigenvalues of  $\mathbf{C}$  are the squared eigenvalues of  $\mathbf{A}$  divided by 3, and the eigenvectors of  $\mathbf{C}$  are the eigenvectors of  $\mathbf{A}$ .

We study three different  $\mathbf{A}$ 's constructed from three choices for the eigenvalues: (1) exponential decay with a constant rate, (2) like the first but with a larger gap between the first and second eigenvalue, and (3) like the first with a larger gap between third and fourth eigenvalue. The three cases of eigenvalues for  $\mathbf{A}$  are shown in the top row of Figure 5.1. Each  $\mathbf{A}$  has the same eigenvectors, which we generate as an orthogonal basis from a random  $10 \times 10$  matrix.

To estimate the eigenvalues, we choose  $N$  as in (4.1) with the multiplier  $\alpha = 2$  and  $k = 6$  eigenvalues of interest, which yields  $N = 28$  evaluations of the gradient. The middle row of Figure 5.1 shows the bootstrap intervals for the first six eigenvalues along with the true eigenvalues of  $\mathbf{C}$ . The small bootstrap intervals suggest confidence in the estimates. The gaps are apparent in the last two cases. The bottom row of Figure 5.1 shows bootstrap intervals on the distance between the true  $k$ -dimensional active subspace and the subspace estimated with the  $N$  samples; the true distance is indicated by the circles. Notice that subspaces corresponding to the larger eigenvalue gap are much better approximated than the others. For example, the three-dimensional subspace is better approximated than the one- and two-dimensional subspaces for the third case.

Next we repeat the study using finite difference approximations of the gradient with step size  $h = 10^{-1}$ ,  $10^{-3}$ , and  $10^{-5}$ . The first of these step sizes is larger than would normally be used for such a model. We chose this large value to study the interplay between inaccurate gradients and the finite sample approximations of the eigenpairs. Figure 5.2 shows the true eigenvalues, their estimates, and the bootstrap intervals for all three cases and all three values of  $h$ ; the horizontal lines show the value of  $h$ . Eigenvalues that are smaller than  $h$  are estimated less accurately than those larger than  $h$ , which is not surprising since we are using first order finite differences. Also the gaps are much less noticeable in the estimates when finite difference parameter is not small enough to resolve the smaller eigenvalue in the pair defining the gap. In fact, this particular problem shows a large gap in the finite difference approximations when there is none in the true eigenvalues; see Figure 5.2(b,c) for examples of this phenomenon.

Figure 5.3 shows the distance between the true active subspace and the finite sample estimate with approximate gradients (circles). We use the bootstrap to estimate the error in the subspace as in Section 4. There is a strong bias in the estimates of the subspace error when the corresponding eigenvalues are not properly resolved. For instance, in Figure 5.3, the estimates of the error for subspaces of dimension 4 through 6 are biased for  $h = 10^{-3}$  and significantly biased for  $h = 10^{-1}$ . Compare this to the error in the last three eigenvalues for the smallest  $h = 10^{-5}$  in Figure 5.2.

**5.2. A parameterized PDE model.** In previous work [7], we exploited the active subspace in the following parameterized PDE model to efficiently construct a kriging surface. Here we perform a more careful study of the variation in the active subspace estimated with finite samples of the gradient. Consider the following linear

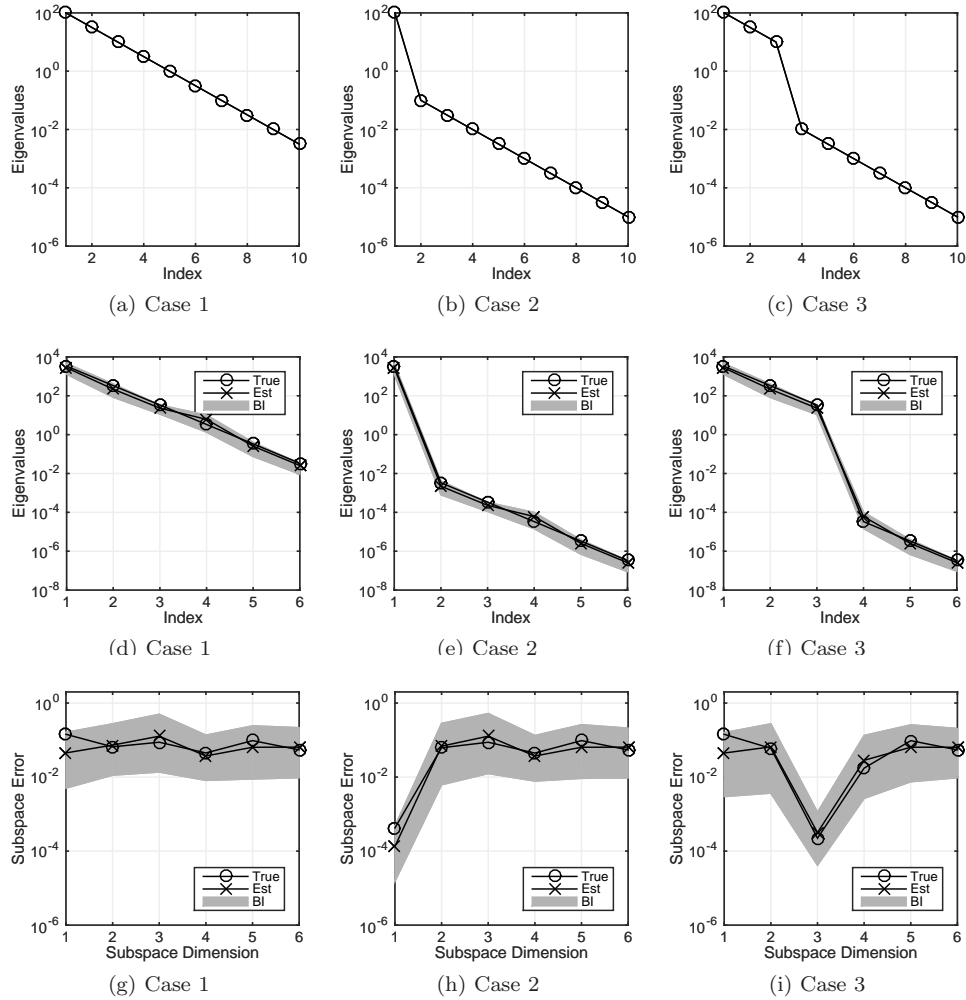


Fig. 5.1: The top row shows the eigenvalues of the three choices for  $\mathbf{A}$ . The second row shows the true and estimated eigenvalues along with the bootstrap intervals; eigenvalues are well approximated for all three cases. The third row shows the distance between the estimated subspace and the true subspace. In practice we do not have the true subspace, but we can estimate the distance with a bootstrap procedure as described in Section 4; the bootstrap intervals are shown, and the accuracy of the subspace estimates corresponds to the gaps in the eigenvalues of  $\mathbf{C}$ .

elliptic PDE with parameterized, variable coefficients. Let  $u = u(\mathbf{s}, \mathbf{x})$  satisfy

$$-\nabla_{\mathbf{s}} \cdot (a \nabla_{\mathbf{s}} u) = 1, \quad \mathbf{s} \in [0, 1]^2. \quad (5.3)$$

We set homogeneous Dirichlet boundary conditions on the left, top, and bottom of the spatial domain  $[0, 1]^2$ ; denote this boundary by  $\Gamma_1$ . The right side of the spatial domain, denoted  $\Gamma_2$ , has a homogeneous Neumann boundary condition. The log of

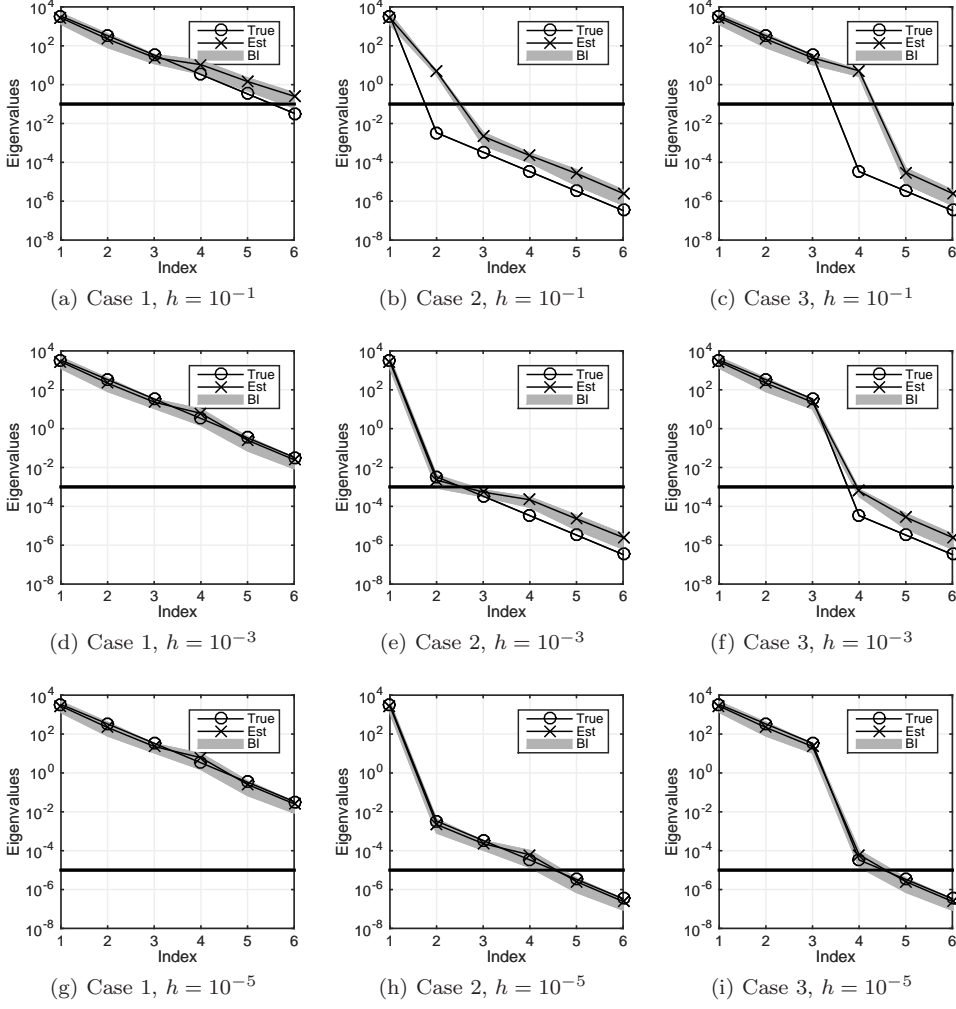


Fig. 5.2: Eigenvalues, estimates, and bootstrap intervals using finite difference gradients with  $h = 10^{-1}$  (top row),  $h = 10^{-3}$  (middle row), and  $h = 10^{-5}$  (bottom row). The horizontal black lines indicate the value of  $h$  in each plot. In general, estimates of eigenvalues smaller than  $h$  are less accurate than those larger than  $h$ .

the coefficients  $a = a(\mathbf{s}, \mathbf{x})$  is given by a truncated Karhunen-Loeve-type expansion

$$\log(a(\mathbf{s}, \mathbf{x})) = \sum_{i=1}^m x_i \gamma_i \phi_i(\mathbf{s}), \quad (5.4)$$

where the  $x_i$  are independent, identically distributed standard normal random variables, and the  $\{\phi_i(\mathbf{s}), \gamma_i\}$  are the eigenpairs of the correlation operator

$$\mathcal{C}(\mathbf{s}, \mathbf{t}) = \exp(-\beta^{-1} \|\mathbf{s} - \mathbf{t}\|_1). \quad (5.5)$$

We study the quality of the active subspace approximation for two correlation lengths,  $\beta = 1$  and  $\beta = 0.01$ . These correspond to *long* and *short* correlation lengths, respec-

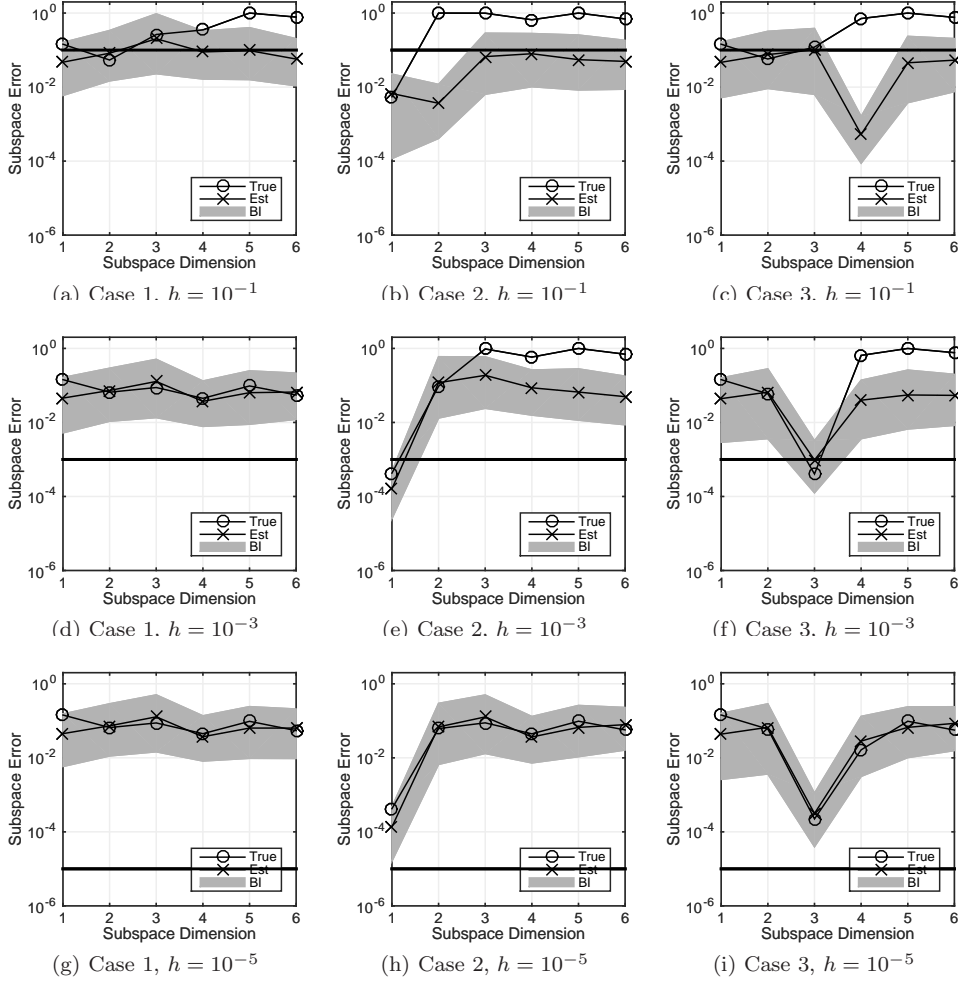


Fig. 5.3: The distance between the true active subspace and its finite sample approximation along with bootstrap intervals for  $h = 10^{-1}$  (top row),  $h = 10^{-3}$  (middle row), and  $h = 10^{-5}$  (bottom row). The subspaces are very poorly approximated when the finite difference step size is not small enough to resolve the eigenvalues corresponding the subspaces; compare to Figure 5.2. However, subspaces with a larger associated eigenvalue gap are generally approximated better than others.

tively, for the random field defining the log of the coefficients. We choose a truncation of the field  $m = 100$ , which implies that the parameter space  $\mathcal{X} = \mathbb{R}^{100}$  with  $\rho$  a standard Gaussian density function. Define the linear function of the solution

$$f(\mathbf{x}) = \frac{1}{|\Gamma_2|} \int_{\Gamma_2} u(\mathbf{s}, \mathbf{x}) ds. \quad (5.6)$$

This is the quantity of interest from the model (more precisely, its approximation with a finite element method). Given a value for the input parameters  $\mathbf{x}$ , we discretize the PDE with a standard linear finite element method using MATLAB's PDE Toolbox. The discretized domain has 34320 triangles and 17361 nodes; the eigenfunctions  $\phi_i$

from (5.4) are approximated on this mesh. We compute the gradient of the quantity of interest (5.6) using a discrete adjoint formulation. Further details appear in our previous work [7].

The top row of Figure 5.4 shows the estimates of the eigenvalues of  $\mathbf{C}$  along with the bootstrap intervals for  $\beta = 1$  in (5.5). The gap between the first and second eigenvalues is apparent and supported by the gap in the corresponding bootstrap intervals. We exploit this gap in [7] to construct an accurate univariate kriging surface of the active variable. The bottom row of Figure 5.4 shows the variance in the estimated subspace as computed with the bootstrap including the bootstrap intervals. The left column of Figure 5.4 uses the multiplier  $\alpha = 2$  when choosing the number  $N$  of gradient samples; the right column uses  $\alpha = 10$ . Notice the overall decrease in both the range of the bootstrap interval and the subspace error as we include more samples. Figure 5.5 shows the identical study with the short correlation length  $\beta = 0.01$  from (5.5).

**6. Summary and conclusions.** Consider a scalar-valued function of several variables. The average outer product of the gradient with itself is the central matrix in the development of active subspaces for dimension reduction. The dominant eigenvectors define the directions along which input perturbations change the output more, on average. We have analyzed a Monte Carlo method for approximating this matrix and its eigenpairs. We use recent theory developed for the eigenvalues of sums of random matrices to analyze the probability that the finite sample eigenvalue estimates deviate from the true eigenvalues, and we combine this analysis with results from matrix computations to derive results for the subspaces. We extend this analysis to quantities computed with samples of approximate gradients, e.g., finite differences. We also provide a practical computational approach that employs the bootstrap to reveal the error in the eigenvalues and the stability of the subspace.

Our analysis offers answers to the following important questions. First, how many gradient samples does one need for an accurate approximation of the first  $k$  eigenvalues? Precise theoretical bounds motivate a heuristic that chooses a number proportional to  $k$  times the log of the dimension  $m$ . Second, what can be said about the accuracy of the estimated subspace? The accuracy of the estimated subspace is directly related to gaps in the eigenvalues. Third, how does one judge the stability of the computed quantities? We propose to use bootstrap intervals for the eigenvalues and the stability of the subspace. Finally, how does this analysis change when gradients are not exact but approximate? Our theory shows that approximate gradients introduce a bias term in the error bounds that goes to zero as the approximate gradients become more accurate. The numerical examples suggest that this bias can produce inaccurate subspaces when the gradients are not well approximated.

**Acknowledgments.** The first author was partially supported by the U.S. Department of Energy Office of Science, Office of Advanced Scientific Computing Research, Applied Mathematics program under Award Number DE-SC-0011077. The second author would like to thank the Simon's Institute for Theory of Computing program on Big Data for the opportunity to learn about the randomized methods used in this paper and NSF CAREER award CCF-114975.

#### REFERENCES

- [1] ATHANASIOS C ANTOUNAS, *Approximation of large-scale dynamical systems*, vol. 6, Society for Industrial and Applied Mathematics, 2005.

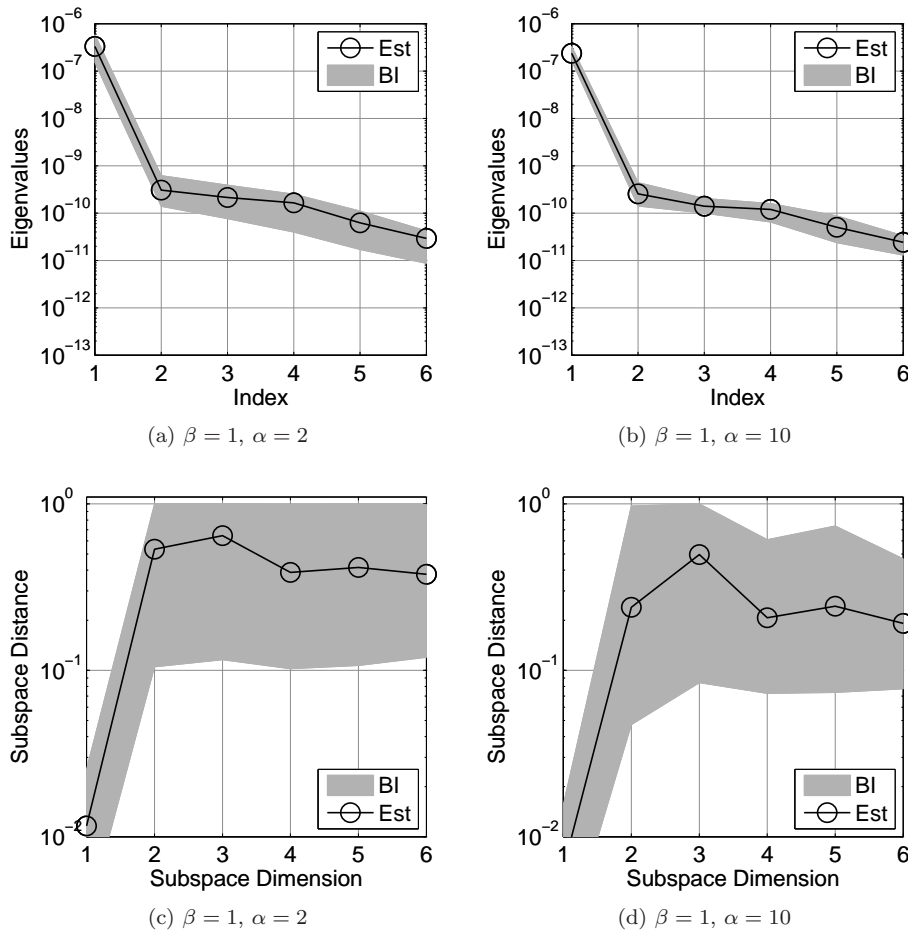


Fig. 5.4: The top row shows estimates of the eigenvalues of  $\mathbf{C}$  along with the bootstrap intervals for the quantity of interest (5.6) from the parameterized PDE model with the long correlation length  $\beta = 1$  from (5.5). The bottom row shows the estimates and bootstrap intervals on the distance between the estimated active subspace and the true active subspace. The left column is computed with the multiplier  $\alpha = 2$  when choosing  $N$ ; the right column uses  $\alpha = 10$ . The gap between the first and second eigenvalue is significant as judged by the gap between the bootstrap intervals.

- [2] G BERKOOZ, P HOLMES, AND J L LUMLEY, *The proper orthogonal decomposition in the analysis of turbulent flows*, Annual Review of Fluid Mechanics, 25 (1993), pp. 539–575.
- [3] ALFIO BORZI AND VOLKER SCHULZ, *Computational Optimization of Systems Governed by Partial Differential Equations*, SIAM, 2012.
- [4] ARTHUR E. BRYSON AND YU-CHI HO, *Applied Optimal Control: Optimization, Estimation, and Control*, Hemisphere Publishing Corporation, 1975.
- [5] HANS-JOACHIM BUNGARTZ AND MICHAEL GRIEBEL, *Sparse grids*, Acta Numerica, 13 (2004), pp. 147–269.
- [6] RUSSEL E. CAFLISCH, *Monte carlo and quasi-monte carlo methods*, Acta Numerica, 7 (1998), pp. 1–49.
- [7] P. CONSTANTINE, E. DOW, AND Q. WANG, *Active subspace methods in theory and practice: Applications to kriging surfaces*, SIAM Journal on Scientific Computing, 36 (2014),

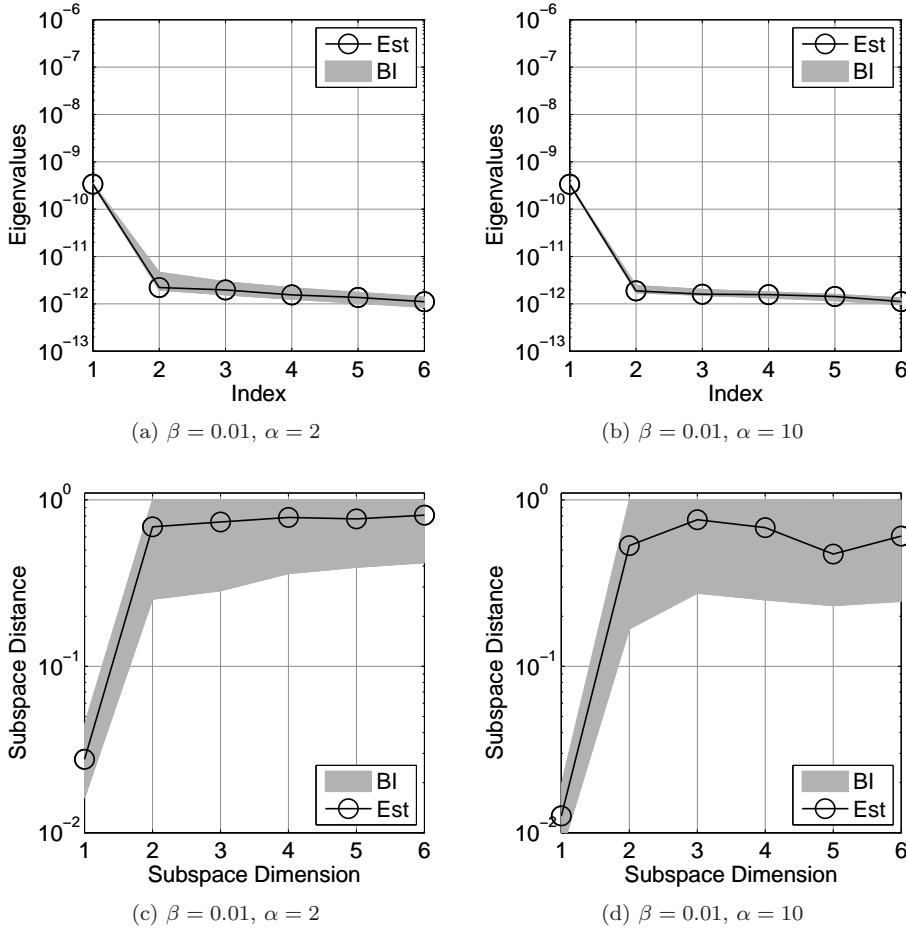


Fig. 5.5: The top row shows estimates of the eigenvalues of  $\mathcal{C}$  along with the bootstrap intervals for the quantity of interest (5.6) from the parameterized PDE model with the short correlation length  $\beta = 0.01$  from (5.5). The bottom row shows the estimates and bootstrap intervals on the distance between the estimated active subspace and the true active subspace. The left column is computed with the multiplier  $\alpha = 2$  when choosing  $N$ ; the right column uses  $\alpha = 10$ . The gap between the first and second eigenvalue is significant as judged by the gap between the bootstrap intervals.

pp. A1500–A1524.

[8] PAUL G CONSTANTINE, BRIAN ZAHARATOS, AND MARK CAMPANELLI, *Discovering an active subspace in a single-diode solar cell model*, arXiv preprint arXiv:1406.7607, (2014).

[9] DAVID L DONOHO, ARIAN MALEKI, INAM UR RAHMAN, MORTEZA SHAHRAM, AND VICTORIA STODDEN, *Reproducible research in computational harmonic analysis*, *Computing in Science & Engineering*, 11 (2009), pp. 8–18.

[10] BRADLEY EFRON AND ROBERT J TIBSHIRANI, *An Introduction to the Bootstrap*, vol. 57, CRC press, 1994.

[11] MASSIMO FORNASIER, KARIN SCHNASS, AND JAN VYBIRAL, *Learning functions of few arbitrary linear parameters in high dimensions*, *Foundations of Computational Mathematics*, 12 (2012), pp. 229–262.

[12] KENJI FUKUMIZU AND CHENLEI LENG, *Gradient-based kernel dimension reduction for regres-*

- sion, *Journal of the American Statistical Association*, 109 (2014), pp. 359–370.
- [13] ALEX GITTENS AND MICHAEL W. MAHONEY, *Revisiting the Nyström method for improved large-scale machine learning.*, in *ICML*, vol. 28, 2013, pp. 567–575.
- [14] ALEX GITTENS AND JOEL A TROPP, *Tail bounds for all eigenvalues of a sum of random matrices*, arXiv preprint arXiv:1104.4513, (2011).
- [15] GENE H GOLUB AND CHARLES F VAN LOAN, *Matrix Computations*, Johns Hopkins University Press, 3rd ed., 1996.
- [16] ANDREAS GRIEWANK, *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*, Society for Industrial and Applied Mathematics, 2000.
- [17] N. HALKO, P. G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, *SIAM Rev.*, 53 (2011), pp. 217–288.
- [18] MARIAN HRISTACHE, ANATOLI JUDITSKY, JORG POLZEHL, AND VLADIMIR SPOKOINY, *Structure adaptive approach for dimension reduction*, *The Annals of Statistics*, 29 (2001), pp. 1537–1566.
- [19] IAN JOLLIFFE, *Principal Component Analysis*, Wiley Online Library, 2005.
- [20] RANDALL J LEVEQUE, *Python tools for reproducible research on hyperbolic problems*, *Computing in Science & Engineering*, 11 (2009), pp. 19–27.
- [21] TRENT W. LUKACZYK, PAUL CONSTANTINE, FRANCISCO PALACIOS, AND JUAN J. ALONSO, *Active Subspaces for Shape Optimization*, American Institute of Aeronautics and Astronautics, 2014/02/16 2014.
- [22] JORGE J. MORÉ AND STEFAN M. WILD, *Estimating derivatives of noisy simulations*, *ACM Trans. Math. Softw.*, 38 (2012), pp. 19:1–19:21.
- [23] JORGE J. MORÉ AND STEFAN M. WILD, *Do you trust derivatives or differences?*, *Journal of Computational Physics*, 273 (2014), pp. 268 – 277.
- [24] ART B. OWEN, *Monte Carlo theory, methods and examples*, 2013. <http://statweb.stanford.edu/~owen/mc/>.
- [25] TRENT M. RUSSI, *Uncertainty Quantification with Experimental Data and Complex System Models*, PhD thesis, UC Berkeley, 2010.
- [26] A. SALTELLI, M. RATTO, T. ANDRES, F. CAMPOLONGO, J. CARIBONI, D. GATELLI, M. SAISANA, AND S. TARANTOLA, *Global Sensitivity Analysis: The Primer*, John Wiley & Sons, 2008.
- [27] A.M. SAMAROV, *Exploring regression structure using nonparametric functional estimation*, *Journal of the American Statistical Association*, 88 (1993), pp. 836–847.
- [28] G. STEWART, *Error and perturbation bounds for subspaces associated with certain eigenvalue problems*, *SIAM Review*, 15 (1973), pp. 727–764.
- [29] ALEX TOWNSEND AND LLOYD N. TREFETHEN, *Continuous analogues of matrix factorizations*, *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 471 (2014), p. 20140585.
- [30] JOEL A TROPP, *User-friendly tail bounds for sums of random matrices*, *Foundations of Computational Mathematics*, 12 (2012), pp. 389–434.
- [31] YINGCUN XIA, *A constructive approach to the estimation of dimension reduction directions*, *The Annals of Statistics*, 35 (2007), pp. 2654–2690.