

# Dimension Reduction via Gaussian Ridge Functions

Pranay Seshadri\*, Shaowu Yuchi<sup>†</sup>, and Geoffrey T. Parks<sup>‡</sup>

**Abstract.** Ridge functions have recently emerged as a powerful set of ideas for subspace-based dimension reduction. In this paper we begin by drawing parallels between ridge subspaces, sufficient dimension reduction and active subspaces, contrasting between techniques rooted in statistical regression and those rooted in approximation theory. This sets the stage for our new algorithm that approximates what we call a Gaussian ridge function—the posterior mean of a Gaussian process on a dimension-reducing subspace—suitable for both regression and approximation problems. To compute this subspace we develop an iterative algorithm that alternates between optimizing over the Stiefel manifold to compute the subspace and optimizing the hyperparameters of the Gaussian process. We demonstrate the utility of the algorithm on two analytical functions, where we obtain near exact ridge recovery, and a turbomachinery case study, where we compare the efficacy of our approach with three well-known sufficient dimension reduction methods: SIR, SAVE, CR. The comparisons motivate the use of the posterior variance as a heuristic for identifying the suitability of a dimension-reducing subspace.

**Key words.** Gaussian process regression, dimension reduction, manifold optimization, active subspaces, sufficient dimension reduction

**AMS subject classifications.** 60G15, 51M35

**1. Introduction.** Dimension reduction is essential for modern data analysis. Its utility spans both simulation informatics and statistical regression fields. The objective of dimension reduction is to obtain parsimoniously parameterized models—functions of few variables that can be used for predicting, understanding and visualizing the trends observed in high-dimensional data sets. These data sets may originate from a tailored *design of experiment* of computer models (see [31]) and are therefore deterministic, or they may originate from the observations of sensors and are therefore subject to measurement noise. Across both paradigms, techniques and ideas within *subspace-based dimension reduction* have proven to be extremely fruitful in inference. Application areas that have benefitted from subspace-based dimension reduction range from airfoil aerodynamics [21] and combustion [28] to economic forecasting (see page 4403 in [2]), turbomachinery aerothermodynamics [35], solar energy [10] and epidemiology [15].

To motivate an exposition of the relevant literature and to set the stage for our paper, we introduce the trivariate function  $y = \log(x_1 + x_2 + x_3)$ . In Figure 1, we generate scatter plots of the function on three different subspaces  $\mathbf{k}^T \mathbf{x}$ —i.e., linear combinations of the three variables  $x_1$ ,  $x_2$  and  $x_3$ . In (a)  $\mathbf{k}$  is one of canonical vectors, in (b) all the elements of  $\mathbf{k}$  are equivalent (and normalized), while in (c) we moderately perturb the entries in (b). Naturally, we find the subspace  $\mathbf{k}^T \mathbf{x}$  in (b) to be *optimal* in the sense that it most accurately characterizes the 3D function as a function of one variable. Scatter plots of the kind in Figure 1, where

---

\*Postdoctoral Fellow, Department of Engineering, University of Cambridge, U.K., and Group Leader, Data-Centric Engineering, The Alan Turing Institute, London, U. K., ps583@cam.ac.uk, pseshadri@turing.ac.uk

<sup>†</sup>Undergraduate Research Assistant, Department of Engineering, University of Cambridge, U. K.

<sup>‡</sup>Reader, Department of Engineering, University of Cambridge, U. K.

the horizontal axis is a subspace of the inputs, are known as *sufficient summary plots* [13] or shadow plots. These plots are useful in identifying low-dimensional structure in high-dimensional problems. However, as the selection of  $\mathbf{k}$  in Figure 1 illustrates, the impediment to realizing this identification is an appropriate choice of the subspace.

Enter *active subspaces* [7], a collection of ideas that facilitate subspace-based dimension reduction for deterministic computer models. If one can compute gradients of a function  $f$ , or even approximate them, then one can identify directions along which (on average)  $f$  exhibits the greatest variation, and conversely directions along which (on average)  $f$  is near constant<sup>1</sup>. Active subspaces owes their development to initial work by Samarov [34] and more recent work by Constantine et al. (see [6, 8, 7]). The latter references build the theory of active subspaces and also offer practical algorithms for their computation. Parallels between active subspaces and *ridge functions* have also been recently explored in the literature [9].

The topic of ridge functions [29] may be interpreted as a broad generalization of ideas within active subspaces. They are of the form  $f(\mathbf{x}) = g(\mathbf{M}^T \mathbf{x})$ , where  $\mathbf{M}$  is a tall matrix, implying that  $g$  is a function of fewer variables than  $f$ . Functions of this form are known as *generalized ridge functions*. When  $d = 1$ , Pinkus [29] defines these functions as *ridge functions*. The elements of  $\mathbf{M}$  and the function  $g$  are useful in identifying low-dimensional structures, in a similar vein to the example discussed above. In the case where  $\mathbf{M}$  is a vector, which implies that  $g$  is univariate function, Cohen et al. [5] provide estimates on how well  $f$  can be approximated using only point evaluations. Building on their work, Fornasier et al. [17] study the case where  $\mathbf{M}$  is a matrix and provide two randomized algorithms for approximating  $g(\mathbf{M}^T \mathbf{x})$ . The work of Tyagi and Cevher [37] is also similar in scope; they provide an algorithm that approximates functions of the form  $f(\mathbf{x}) = \sum_{i=1}^n g_i(\mathbf{M}_i^T \mathbf{x})$ . Clearly, there are similarities between ridge functions, projection pursuit regression (see page 390 in [22]) and, by association, single layer neural networks. The projection pursuit regression problem can be interpreted as an approximation with

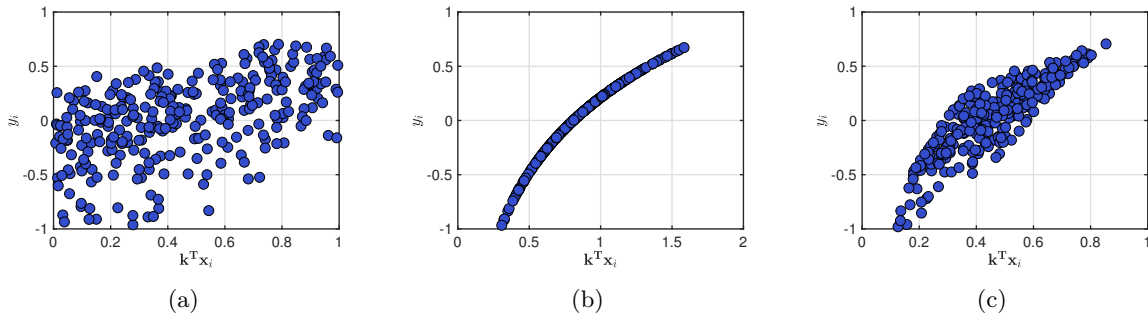
$$(1) \quad \sum_{i=1}^k g_i(\mathbf{m}_i^T \mathbf{x}),$$

where  $\mathbf{m}_i \in \mathbb{R}^d$  and the univariate function  $g_i$  are unknown; determined via a stepwise greedy algorithm following Friedman and Stuetzle [18]. It is worth emphasizing that although our function may not admit an analytical ridge structure, i.e.,  $f(\mathbf{x}) \neq g(\mathbf{M}^T \mathbf{x})$ , one may still be able to approximate  $f(\mathbf{x}) \approx g(\mathbf{M}^T \mathbf{x})$  extracting useful inference. We refer to this as a *ridge approximation*, that is obtained by evaluations of  $f$  and by minimizing a suitably constructed objective function to ascertain  $\mathbf{M}$  and  $g$ .

In introducing methods for approximation and statistical regression (projection pursuit regression), it is important to note the differences between both paradigms. In statistical regression, as alluded to previously,  $f(\mathbf{x})$  is not deterministic and therefore has a conditional density that is not a delta function [7]. Furthermore, the joint density of  $\mathbf{x}$  is not known in the regression setting. In comparison, when working with computer models—i.e., in the

---

<sup>1</sup>Assuming, of course,  $f$  admits an active subspace; there is no guarantee that every differentiable multivariate function  $f$  admits an active subspace.



**Figure 1.** Sufficient summary plots of the function  $y = \log(x_1 + x_2 + x_3)$  where the horizontal axis is given by  $\mathbf{k}^T(x_1, x_2, x_3)$ . In (a)  $\mathbf{k} = (1, 0, 0)^T$ , in (b)  $\mathbf{k} = (1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3})^T$ , while in (c)  $\mathbf{k} = (0.2, 0.5, 0.2)^T$ .

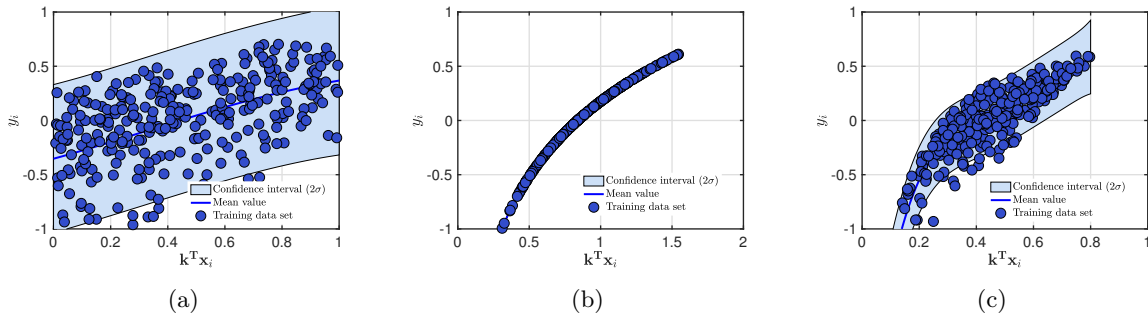
approximation paradigm—the joint density of  $\mathbf{x}$  is known and often prescribed according to some optimal design of experiment, for example, Latin hypercube sampling, D-optimal design, etc. Lastly, techniques for estimating gradients for unknown regression functions are notoriously expensive in high dimensions. Thus, non-gradient-based approaches are sought.

The field of sufficient dimension reduction is rich in ideas for achieving subspace-based dimension reduction within a statistical regression framework. The goal here is to replace the vector of covariates (inputs) with their projection onto a subspace assembled from the space of the covariates themselves [14]. This subspace must be constructed without loss of information based either on the conditional mean, conditional variance or, more generally, the conditional distribution of the outputs with respect to the inputs. The numerous methods for estimating these subspaces include sliced inverse regression (SIR) [24], sliced average variance estimation (SAVE) [11] and contour regression (CR) [23], to name but a few (see [27] for a detailed review).

We now return to the sufficient summary plots in Figure 1. Assume we did not know the true function  $f(x) = \log(x_1 + x_2 + x_3)$  and that subfigures (a), (b) and (c) were simply the outcome of any one of the aforementioned algorithms. How do we determine that the subspace in (b) is ideally suited for the data? The logic we pursue in the paper is as follows. Assume we utilize Gaussian process regression (see Chapter 2 in [33]) to estimate the mean and standard deviation of the response  $y_i$ , as shown in Figure 2. Then our criterion for selecting  $\mathbf{k}$  in (b) is clear: it reduces the Gaussian process yielded standard deviation. In other words, for the same value of  $\mathbf{x}_i$  there is only one unique  $y_i$ , ensuring that the deviation from the mean is near zero.

There are two key ideas that emerge from this discussion. The first is that the posterior mean of a Gaussian process computed on a dimension-reducing subspace is a good candidate for  $g$ . In fact, this idea has been previously studied in [36] and [25], albeit using different techniques than those presented here. The second idea is that the suitability of a dimension-reducing subspace is reflected by the posterior variance: should the posterior variance be too large, as in Figure 2(a), then one may need to opt for a different choice of  $\mathbf{k}$ .

In a nutshell, our objective in this paper is to offer an algorithm for computing a dimension-reducing subspace, under the constraint that  $g$  is the posterior mean of a Gaussian process



**Figure 2.** Sufficient summary plots of the function  $y = \log(x_1 + x_2 + x_3)$  where the horizontal axis is given by  $\mathbf{k}^T(x_1, x_2, x_3)$  with a Gaussian process regression response surface. In (a)  $\mathbf{k} = (1, 0, 0)^T$ , in (b)  $\mathbf{k} = (1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3})^T$ , while in (c)  $\mathbf{k} = (0.2, 0.5, 0.2)^T$ .

on this subspace. The remainder of this paper is structured as follows. In section 2 we survey the connections between ridge, active and sufficient dimension reduction subspaces. Then, motivated by the techniques in [36] and [9], we introduce an algorithm for computing a ridge subspace that uses tools from Gaussian process regression and manifold optimization. This is followed by numerical examples in section 4.

**2. Ridge functions.** If there exists a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  and a matrix  $\mathbf{M} \in \mathbb{R}^{d \times m}$ , with  $m \leq d$ , that satisfy

$$(2) \quad f(\mathbf{x}) = g(\mathbf{M}^T \mathbf{x}),$$

for a function  $g : \mathbb{R}^m \rightarrow \mathbb{R}$ , then  $f$  is a *ridge function* [29]. We call the subspace associated with the span of  $\mathbf{M}$  its *ridge subspace*, denoted by  $\mathcal{S}(\mathbf{M})$ . We also assume that the columns of  $\mathbf{M}$  are orthonormal, i.e.,  $\mathbf{M}^T \mathbf{M} = \mathbf{I}$ . The above definitions imply that the gradient of  $f$  is zero along directions that are orthogonal to  $\mathcal{S}(\mathbf{M})$ . In other words, if we replace  $\mathbf{x}$  with  $\mathbf{x} + \mathbf{h}$  where  $\mathbf{M}^T \mathbf{h} = 0$ , then it is trivial to see that  $f(\mathbf{x} + \mathbf{h}) = g(\mathbf{M}^T(\mathbf{x} + \mathbf{h})) = f(\mathbf{x})$ .

**2.1. Computing the optimal subspace.** So, given point evaluations of  $f$ , how does one compute  $\mathbf{M}$ ? Also, what are the properties of  $\mathbf{M}$ ? In their paper, Constantine et al. [9] draw our attention to the conditional and marginal densities defined on the subspace coordinates of the ridge subspace and its orthogonal complement. They present two main ideas.

The first is the *orthogonal invariance* associated with  $\mathbf{M}$ , i.e., we are not interested in estimating a particular  $\mathbf{M}$ , but rather the subspace of  $\mathbf{M}$ —i.e., the ridge subspace. To study this further, we place a few additional assumptions on (2). Let the joint probability density function  $\rho(\mathbf{x})$  and marginal densities  $\rho_i$  be related by

$$(3) \quad \rho(\mathbf{x}) = \prod_{i=1}^d \rho(x^{(i)}).$$

We assume that  $f$  is square integrable with respect to this probability density function

$$(4) \quad \int_{\rho} f^2(\mathbf{x}) \rho(\mathbf{x}) d\mathbf{x} < \infty,$$

where the domain of  $f$  is the support of  $\rho$ . Define  $\mathcal{S}(\mathbf{N})$  to be the orthogonal complement of the ridge subspace, where

$$(5) \quad \mathbf{N} = (\text{null}(\mathbf{M}^T)) \quad \text{with} \quad \mathbf{N} \in \mathbb{R}^{d \times (d-m)}.$$

Placing the columns of  $\mathbf{M}$  and  $\mathbf{N}$  together in a single matrix

$$(6) \quad \mathbf{G} = [ \mathbf{M} \quad \mathbf{N} ],$$

readily implies that  $\mathbf{G}\mathbf{G}^T = \mathbf{I}$ . We now write the full decomposition of  $f$  as

$$(7) \quad f(\mathbf{x}) = f(\mathbf{G}\mathbf{G}^T\mathbf{x}) = f(\mathbf{M}\mathbf{M}^T\mathbf{x} + \mathbf{N}\mathbf{N}^T\mathbf{x}) = f(\mathbf{M}\mathbf{u} + \mathbf{N}\mathbf{v}),$$

with the following subspace coordinates

$$(8) \quad \mathbf{u} \in \mathbb{R}^m \quad \text{and} \quad \mathbf{v} \in \mathbb{R}^{d-m},$$

where

$$(9) \quad \mathbf{u} = \mathbf{M}^T\mathbf{x} \quad \text{and} \quad \mathbf{v} = \mathbf{N}^T\mathbf{x}.$$

It should be noted here that if  $f$  admits a ridge function structure, then the gradient of  $f$  along the directions  $\mathbf{N}\mathbf{v}$  are expected to be zero.

For a fixed subspace coordinate  $\mathbf{u}$ , Constantine et al. [9] write the conditional expectation of  $f$  to be

$$(10) \quad \mathbb{E}(f|\mathbf{u}, \mathbf{M}) = \int f(\mathbf{M}\mathbf{u} + \mathbf{N}\mathbf{v}) \pi(\mathbf{v}|\mathbf{u}) d\mathbf{v}$$

where the conditional probability  $\pi(\mathbf{v}|\mathbf{u})$  can be expressed in terms of its marginal  $\pi(\mathbf{u})$  and  $\pi(\mathbf{u}, \mathbf{v})$  joint probabilities

$$(11) \quad \pi(\mathbf{v}|\mathbf{u}) = \frac{\pi(\mathbf{u}, \mathbf{v})}{\pi(\mathbf{u})} = \frac{\rho(\mathbf{M}\mathbf{u} + \mathbf{N}\mathbf{v})}{\int \rho(\mathbf{M}\mathbf{u} + \mathbf{N}\mathbf{v}) d\mathbf{v}}.$$

To demonstrate the orthogonal invariance associated with  $\mathbf{M}$  we replace it with  $\mathbf{M}\mathbf{Q}$  where  $\mathbf{Q} \in \mathbb{R}^{m \times m}$  is an orthogonal matrix. Plugging this into (10) yields

$$(12) \quad \mathbb{E}(f|\mathbf{u}, \mathbf{M}\mathbf{Q}) = \int f(\mathbf{M}\mathbf{Q}(\mathbf{M}\mathbf{Q})^T\mathbf{x} + \mathbf{N}\mathbf{v}) \pi(\mathbf{v}|\mathbf{u}) d\mathbf{v}$$

$$(13) \quad = \int f\left(\mathbf{M} \underbrace{\mathbf{Q}\mathbf{Q}^T}_{\mathbf{I}} \mathbf{M}^T\mathbf{x} + \mathbf{N}\mathbf{v}\right) \pi(\mathbf{v}|\mathbf{u}) d\mathbf{v}$$

$$(14) \quad = \int f(\mathbf{M}\mathbf{u} + \mathbf{N}\mathbf{v}) \pi(\mathbf{v}|\mathbf{u}) d\mathbf{v}.$$

It can be easily shown that the conditional probability also remains the same under the orthogonal transformation  $\mathbf{M}\mathbf{Q}$ . This demonstrates that the particular choice of  $\mathbf{M}$  does not

matter, but rather its subspace does. In [9], the authors draw our attention to Theorem 8.3 in Pinkus [29], which proves that  $\mathbb{E}(f|\mathbf{u}, \mathbf{M})$  is the *optimal ridge profile*—a function—in the  $L_2$  norm.

So, how do we compute this *optimal subspace* associated with this ridge profile? The second point made in [9] is that this subspace can be computed by solving the following quadratic form

$$(15) \quad \begin{aligned} \underset{\mathbf{M}}{\text{minimize}} \quad & \omega(\mathbf{M}) = \frac{1}{2} \int_{\rho} (f(\mathbf{x}) - \mathbb{E}(f|\mathbf{u}, \mathbf{M}))^2 \rho(\mathbf{x}) \, d\mathbf{x} \\ \text{subject to} \quad & \mathbf{M} \in \mathbb{G}(m, d). \end{aligned}$$

The constraint in (15) restricts  $\mathbf{M}$  to be from the subspace  $\mathbb{G}(m, d)$ , which denotes the space of  $m$ -dimensional subspaces of  $\mathbb{R}^d$ , i.e., the *Grassman manifold* (see page 30 in [1]). Recognizing (15) as a manifold optimization problem permits us to compute the gradients on the Grassman manifold. Following page 321 in [16], and denoting the gradient on the Grassman manifold by the symbol  $\nabla_{\mathbb{G}}$ , we have

$$(16) \quad \nabla_{\mathbb{G}} \omega(\mathbf{M}) = \int_{\rho} (f(\mathbf{x}) - \mathbb{E}(f|\mathbf{u}, \mathbf{M})) (-\nabla_{\mathbb{G}} \mathbb{E}(f|\mathbf{u}, \mathbf{M})) \rho(\mathbf{x}) \, d\mathbf{x}$$

$$(17) \quad = \int_{\rho} (\mathbb{E}(f|\mathbf{u}, \mathbf{M}) - f(\mathbf{x})) (\mathbf{I} - \mathbf{M}\mathbf{M}^T) \frac{\partial}{\partial \mathbf{M}} \mathbb{E}(f|\mathbf{u}, \mathbf{M}) \rho(\mathbf{x}) \, d\mathbf{x}.$$

A few remarks regarding this optimization problem are in order. First, to compute derivative inside the integral, we require  $f \in C^1(\mathbb{R}^d)$ , which denotes a class of real-valued functions with continuous first derivatives. Provided the gradients  $\frac{\partial}{\partial \mathbf{M}} \mathbb{E}(f|\mathbf{u}, \mathbf{M}) \in \mathbb{R}^{m \times d}$  can be determined, solutions to this optimization problem can be computed using an appropriate gradient-based optimizer. In subsection 3.3 we use tools from Gaussian processes to estimate these gradients. Second, (15) is not necessarily convex, therefore we are not guaranteed a unique global minimum during an optimization. However, as we demonstrate in section 4, near local optimal solutions can be computed.

Another salient point concerns the choice of the manifold itself. In using the Grassman manifold, we enforce the assumption of *homogeneity*, i.e.,  $\omega(\mathbf{M}) \equiv \omega(\mathbf{M}\mathbf{Q})$  implying—as stated previously—that the specific entries in  $\mathbf{M}$  do not matter, but rather its subspace does. In other words, our objective is invariant to any choice of the basis [16]. In 3 we discuss the ramifications of this assumption when using Gaussian process regression.

**2.2. Connections between ridge subspaces and sufficient dimension reduction.** One approach to estimate the ridge subspace is to develop a computational method that implements (15). In subsection 3.2 we detail such a procedure using Gaussian process regression. However, other algorithms that utilize ideas from sufficient dimension reduction theory can also be leveraged. The key is to understand the relationship between the subspace computed from techniques like SIR, SAVE and CR and the *ridge subspace*. As a prelude to exploring these connections, it is important to distinguish statistical regression from approximation.

In the regression setting we are given independent and identically distributed random variables  $(\mathbf{x}_i, y_i)$  from an *unknown* joint distribution  $\pi(\mathbf{x}, y)$ , where  $i = 1, \dots, N$  for some

*N.* Our objective is to characterize the conditional dependence of  $y$  on  $\mathbf{x}$  either through its expectation  $\mathbb{E}(y|\mathbf{x})$  or its variance  $\sigma^2(y|\mathbf{x})$ . This requires a model  $y = h(\mathbf{x}) + \epsilon$  for predicting  $y$  given  $\mathbf{x}$ , where  $\epsilon$  is a zero-mean random error that is independent of  $\mathbf{x}$ . This implies that for the same  $\mathbf{x}_*$  we have multiple values of  $y_*$ . Here  $h$  is a parametric function and will typically be obtained by penalizing the errors in prediction via an appropriate *loss function*, the most common one being the  $l_2$  norm error [22]. In the approximation setting, we are given input / output pairs  $(\mathbf{x}_i, y_i)$  obtained by evaluating a computer model under a known density  $\pi(\mathbf{x})$  via an appropriately chosen *design of experiment*. These inputs may be either design variables, boundary conditions or state variables. As the computer model, in general, does not have any random error associated with it, multiple evaluations for a particular  $\mathbf{x}_*$  will always yield the same  $y_*$ . In this context, our objective is to obtain a useful approximation  $y = f(\mathbf{x})$  valid for *all*  $\mathbf{x}$  over  $\pi(\mathbf{x})$ .

Now, despite these differences, ideas from sufficient dimension reduction can be used for approximation. Glaws et al. [19] detail the mathematical nuances involved in doing so for SIR and SAVE. They prove that the conditional independence of the inputs and outputs in the approximation setting—i.e., for a deterministic function—is equivalent to the function being a ridge function (see Theorem 6 in [19]). But is there any benefit in using SIR, SAVE and other sufficient dimension reduction methods for dimension reduction when approximating? Adraghi and Cook [2] suggest that moment-based sufficient dimension reduction techniques such as SIR and SAVE, which are designed to provide estimates of the minimum sufficient linear reduction, are not tailored for prediction. They further suggest that model-based inverse regression techniques, such as those presented in [12], may be more useful. In section 4 we test their suggestion by comparing the results of our proposed algorithm with SIR, SAVE, MAVE and CR.

**2.3. Connections between ridge subspaces and active subspaces.** Active subspaces [7] refer to a set of ideas for parameter-based dimension reduction in the *approximation* setting (see subsection 2.2). To compute the active subspace of a function  $f$ , we require that  $f$  be differentiable and Lipschitz continuous—implying that the norm of its gradient  $\nabla f_{\mathbf{x}}$  is bounded. Then one assembles the covariance matrix  $\mathbf{C}$

$$(18) \quad \mathbf{C} = \int \nabla f_{\mathbf{x}} \nabla f_{\mathbf{x}} \rho(\mathbf{x}) d\mathbf{x},$$

also known as the outer product of the gradient matrix<sup>2</sup>. What follows is an eigenvalue decomposition of  $\mathbf{C}$ , where the first  $m$  eigenvectors are used to define the *active* subspaces and remaining  $d - m$  eigenvectors the *inactive* subspaces; the theorem below connects active subspaces and ridge subspaces.

**Theorem 1.** *If we assume that  $\mathbf{C}$  admits the eigendecomposition*

$$(19) \quad \mathbf{C} = \begin{bmatrix} \mathbf{M} & \mathbf{N} \end{bmatrix} \begin{bmatrix} \mathbf{\Lambda}_1 & \\ & \mathbf{\Lambda}_2 \end{bmatrix} \begin{bmatrix} \mathbf{M} & \mathbf{N} \end{bmatrix}^T \quad \text{where } \mathbf{\Lambda}_1 \in \mathbb{R}^{m \times m},$$

---

<sup>2</sup>Note that this is a symmetric positive semidefinite matrix.

then

$$(20) \quad \left( \int (f(\mathbf{x}) - \mathbb{E}(f|\mathbf{M}^T \mathbf{x}))^2 \rho(\mathbf{x}) d\mathbf{x} \right)^{\frac{1}{2}} \leq C (\text{trace}(\Lambda_2))^{\frac{1}{2}},$$

where  $C$  is a constant that depends on  $\rho(\mathbf{x})$ .

The proof follows from Theorem 4.4 in [7]. It is important to emphasize that, in practice, the right-hand side of (20) is difficult to compute because we do not have access to the eigenvalues associated with the covariance matrix that yields eigenvectors  $[\mathbf{M}, \mathbf{N}]$ . Furthermore, while the decay in the eigenvalues above can be used to set the size of  $m$  when using active subspaces, when approximating ridge subspaces, one has to opt for other criteria—for instance, the authors in [36] utilize a Bayesian information criterion to determine  $m$ .

**3. Gaussian processes.** The definition of a Gaussian process is a “collection of random variables, any finite number of which have a joint Gaussian distribution” [33]. It is completely defined by its mean and covariance functions. Restating the definition in the form we need, we have

$$(21) \quad f(\mathbf{x}) \sim \mathcal{GP}(\mu(\mathbf{u}), \mathbf{k}(\mathbf{u})),$$

where  $\mu$  is the mean function and  $\mathbf{k}$  represents a parameterized covariance matrix, albeit with a slight abuse of notation. It is important to note that we define our Gaussian process on the reduced coordinates  $\mathbf{u} = \mathbf{M}^T \mathbf{x}$ .

In this section we have two main goals. First, to define a Gaussian ridge function and to outline its properties. Second, to provide an algorithm that computes a Gaussian ridge function. To ease our exposition, we assume the existence of a set of:

- A set of input/output training data pairs:

$$(22) \quad \hat{\mathbf{f}} = \begin{pmatrix} \hat{f}_1 \\ \vdots \\ \hat{f}_{N_{train}} \end{pmatrix}, \quad \hat{\mathbf{U}} = (\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_N)^T \quad \text{and} \quad \hat{\mathbf{u}}_i = \mathbf{M}^T \hat{\mathbf{x}}_i,$$

- A set of input/output testing data pairs:

$$(23) \quad \tilde{\mathbf{f}} = \begin{pmatrix} \tilde{f}_1 \\ \vdots \\ \tilde{f}_{N_{test}} \end{pmatrix}, \quad \tilde{\mathbf{U}} = (\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_N)^T \quad \text{and} \quad \tilde{\mathbf{u}}_i = \mathbf{M}^T \tilde{\mathbf{x}}_i,$$

**3.1. Gaussian ridge functions (posterior mean).** Succinctly stated, a Gaussian ridge function is the posterior mean of the Gaussian process in (21), written as

$$(24) \quad \bar{g}(\mathbf{u}) = \mathbf{k}(\mathbf{u}, \hat{\mathbf{U}}) \beta.$$



It is a linear sum of kernel functions with coefficients  $\boldsymbol{\beta}$ —an outcome of the representer theorem (see page 132 of [33]). Before we define these coefficients, we first focus on the covariance term on the right-hand side in (24). In what follows, we employ the squared exponential covariance kernel

$$(25) \quad k(\mathbf{u}_i, \mathbf{u}_j) = \sigma_f^2 \exp\left(-\frac{1}{2}(\mathbf{u}_i - \mathbf{u}_j)^T \boldsymbol{\Gamma}^{-1}(\mathbf{u}_i - \mathbf{u}_j)\right),$$

where the diagonal matrix  $\boldsymbol{\Gamma}$  is given by

$$(26) \quad \boldsymbol{\Gamma} = \begin{bmatrix} l_1^2 & & 0 \\ & \ddots & \\ 0 & & l_m^2 \end{bmatrix}.$$

The constant  $\sigma_f^2$  is known as the *signal variance* and constants  $l_1, \dots, l_m$  are the *correlation lengths* along each of the  $m$  coordinates of  $\mathbf{u}$ . The parameterized covariance matrix  $\mathbf{k}(\mathbf{u}, \hat{U}) \in \mathbb{R}^{1 \times N}$  in (24) can now be written as

$$(27) \quad \mathbf{k}(\mathbf{u}, \hat{U}) = (k(\mathbf{u}, \hat{\mathbf{u}}_1), \dots, k(\mathbf{u}, \hat{\mathbf{u}}_N)).$$

The vector of coefficients  $\boldsymbol{\beta} \in \mathbb{R}^N$  is given by

$$(28) \quad \boldsymbol{\beta} = (\mathbf{K} + \sigma_n \mathbf{I})^{-1} \hat{\mathbf{f}},$$

where the  $(i, j)$ -th entry of the matrix  $\mathbf{K} \in \mathbb{R}^{N \times N}$  is given by

$$(29) \quad \mathbf{K}(i, j) = k(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j),$$

in other words it depends only on the training data set. Once again, for notational convenience, we define the posterior mean values at the training and testing points by

$$(30) \quad \bar{\mathbf{g}}_{train} = \begin{pmatrix} \bar{g}(\hat{\mathbf{u}}_1) \\ \vdots \\ \bar{g}(\hat{\mathbf{u}}_{N_{train}}) \end{pmatrix}, \quad \bar{\mathbf{g}}_{test} = \begin{pmatrix} \bar{g}(\tilde{\mathbf{u}}_1) \\ \vdots \\ \bar{g}(\tilde{\mathbf{u}}_{N_{test}}) \end{pmatrix}.$$

Typically, the cost of constructing a Gaussian process is dominated by  $\mathcal{O}(N_{train}^3)$  linear solve operation (see (28)). By defining our Gaussian process over a reduced set of coordinates, there are two instances where we reduce the computational cost:

1. we reduce the number of coordinates where our kernel has to be evaluated, which happens  $\mathcal{O}(N_{train}^2)$  times;
2. we reduce the number of hyperparameters that need to be optimized.

The first instance is dominated by the linear solve, while the second is harder to analyze and will depend on whether there is a hyperparameter associated with each dimension  $d$ .

To aid the optimization strategy we describe in the forthcoming section, we are interested in obtaining the gradients of (24); this yields

$$(31) \quad \begin{aligned} \frac{\partial \bar{g}}{\partial \mathbf{u}} &= \frac{\partial \mathbf{k}(\mathbf{u}, \hat{U})}{\partial \mathbf{u}} \boldsymbol{\beta} \\ &= -\boldsymbol{\Gamma}^{-1} \tilde{U} \left( \mathbf{k}(\mathbf{u}, \hat{U})^T \odot \boldsymbol{\beta} \right), \end{aligned}$$

where the symbol  $\odot$  indicates an element-wise product. In practice, any differentiable kernel may be used to compute the gradient of the Gaussian ridge function. Now, in using a Gaussian ridge function, we make the assertion that it is approximately the expectation of  $f$  projected on  $\mathbf{u} = \mathbf{M}^T \mathbf{x}$ , i.e.,

$$(32) \quad \mathbb{E}(f|\mathbf{u}, \mathbf{M}) \approx \bar{g}(\mathbf{u}).$$

Its gradients, can then be computed via

$$(33) \quad \begin{aligned} \frac{\partial}{\partial \mathbf{M}} \bar{g}(\mathbf{u}) &= \frac{\partial \bar{g}}{\partial \mathbf{u}} \frac{\partial (\mathbf{M}^T \mathbf{x})}{\partial \mathbf{M}} \\ &= -\boldsymbol{\Gamma}^{-1} \tilde{U} \left( \mathbf{k}(\mathbf{u}, \hat{U})^T \odot \boldsymbol{\beta} \right) \mathbf{x}. \end{aligned}$$

It is clear that the quality of our approximation in (32) is contingent upon a suitable choice of the hyperparameters

$$(34) \quad \boldsymbol{\theta} = \{\sigma_f^2, \sigma_n^2, l_1, \dots, l_m\}.$$

A well-worn heuristic to estimate these hyperparameters, which in itself is a non-convex optimization problem, is to maximize the marginal likelihood (see page 113 [33])

$$(35) \quad \begin{aligned} &\underset{\boldsymbol{\theta}}{\text{maximize}} \quad \log p(\mathbf{u}, \boldsymbol{\theta}) \\ &\text{subject to} \quad \log p(\mathbf{u}, \boldsymbol{\theta}) = -\frac{1}{2} \mathbf{f}^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{f} - \frac{1}{2} \log |\mathbf{K} + \sigma_n^2 \mathbf{I}| - \frac{N}{2} \log(2\pi). \end{aligned}$$

The first term in  $\log p(\mathbf{u}, \boldsymbol{\theta})$ , which depends solely on the covariance matrix and values of  $f$ , is the *data-fit* term. The second term, which depends only upon the covariance function, is known as the *complexity penalty*, while the third term is simply a normalizing constant. The problem of computing the maximum in (35) can be readily solved via a gradient-based optimizer; formulas for the gradients of  $\log p$  with respect to the hyperparameters are given in section 5.9 of [33].

**3.2. The utility of the posterior variance.** In our presentation of Gaussian ridge functions so far we have not discussed the posterior variance. It is given by

$$(36) \quad \mathbb{V}[\mathcal{G}\mathcal{P}] = k(\mathbf{u}, \mathbf{u}) - \mathbf{k}(\mathbf{u}, \hat{U})^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{u}, \hat{U}),$$

and can be very useful in quantifying the suitability of a particular choice of  $\mathbf{M}$ . Recall the plots in Figures 1 and 2. Our perception of a successful dimension reduction strategy through a sufficient summary plot is intimately tied to the variance in  $f$  for a fixed coordinate in the dimension-reducing subspace. Here, we argue that the expectation of the posterior variance,  $\mathbb{E}[\mathbb{V}(g)]$ , can be used to gauge the overall effectiveness of the subspace approximant  $\tilde{\mathbf{M}}$ .

We make the above argument for two reasons. First, different dimension reduction methods—tailored for the same projection  $\mathbb{R}^d \rightarrow \mathbb{R}^m$ —will likely yield different subspaces. In [14] a clear distinction is made between dimension reduction techniques that identify a *central subspace* from those that find a *central mean subspace*; an algorithm for computing the *central variance subspace* is introduced in [39]. Our hypothesis is that regardless of the subspace sought, to infer the suitability of a particular  $\tilde{\mathbf{M}}$  as a ridge subspace we must be able to compare different methods that by definition will yield different subspaces. To this end, while a metric based on the subspace angle  $\phi$  (see page 329 of [20])—i.e.,

$$(37) \quad \begin{aligned} \text{dist}(\mathbf{M}, \tilde{\mathbf{M}}) &= \left\| \mathbf{M}\mathbf{M}^T - \tilde{\mathbf{M}}\tilde{\mathbf{M}}^T \right\|_2 \\ &= \sin(\phi) \end{aligned}$$

—is useful for comparing the *distance* between the *ridge subspace*  $\mathbf{M}$  and its approximant  $\tilde{\mathbf{M}}$ , it is limited to comparisons with one technique / subspace. The second point we articulate concerns cases when  $m \geq 3$ , where it is difficult to visually perceive how suitable a choice of  $\bar{g}$  is. Here, the measure  $\mathbb{E}[\mathbb{V}(g)]$  can be used (i) to quantify the error in the approximation, and (ii) to guide further computer simulations.

**3.3. Computing the Gaussian ridge via manifold optimization.** Recall the definition of the correlation lengths associated with the Gaussian process regression model in (26). If we were to set the correlation lengths to be the same, i.e.,  $l_1 = l_2 = \dots = l_m$ , then we satisfy our requirement of homogeneity, and can optimize (15) using the gradient computation in (17). While this may be a potential solution strategy it does require us to modify the maximum likelihood criterion in (35). To circumvent this, we propose to alter the constraint in (15); instead of optimizing over the Grassman manifold, we optimize over the Stiefel manifold,

$$(38) \quad \begin{aligned} \underset{\mathbf{M}}{\text{minimize}} \quad \omega(\mathbf{M}) &= \frac{1}{2} \int_{\rho} (f(\mathbf{x}) - \bar{g}(\mathbf{M}^T \mathbf{x}))^2 \rho(\mathbf{x}) d\mathbf{x} \\ \text{subject to} \quad \mathbf{M} &\in \text{St}(m, d). \end{aligned}$$

Here  $\text{St}(m, d)$  is the set of all  $d \times m$  orthogonal matrices

$$(39) \quad \text{St} := \left\{ \mathbf{M} \in \mathbb{R}^{d \times m} : \mathbf{M}^T \mathbf{M} = \mathbf{I}_m \right\};$$

it is an embedded submanifold of  $\mathbb{R}^{d \times m}$  and is compact [1]. Denoting the gradients on the Stiefel manifold by  $\tilde{\nabla}$ , gradients of  $\omega$  can be computed via

$$(40) \quad \nabla_{\text{St}} \omega(\mathbf{M}) = \int_{\rho} (f(\mathbf{x}) - \bar{g}(\mathbf{M}^T \mathbf{x})) (-\tilde{\nabla}_{\text{St}} \bar{g}(\mathbf{M}^T \mathbf{x})) \rho(\mathbf{x}) d\mathbf{x}$$

$$(41) \quad = \int_{\rho} (\bar{g}(\mathbf{M}^T \mathbf{x}) - f(\mathbf{x})) \left( \frac{\partial}{\partial \mathbf{M}} \bar{g}(\mathbf{M}^T \mathbf{x}) - \mathbf{M} \frac{\partial}{\partial \mathbf{M}} \bar{g}(\mathbf{M}^T \mathbf{x})^T \mathbf{M} \right) \rho(\mathbf{x}) d\mathbf{x},$$

(see 2.53 of [16]). We use this expression in our algorithm for estimating  $\mathbf{M}$  using a Gaussian ridge function. Ours is an iterative approach—optimizing over the Stiefel manifold for estimating  $\mathbf{M}$  and then over the space of the hyperparameters for determining  $\boldsymbol{\theta}$ , where the correlation lengths along each of the  $m$  directions can be different. We remark here that our approach is similar to the alternating approach from Algorithm 2 in Constantine et al. [9] that uses polynomial approximations.

Our algorithmic workflow is given as follows. First, the data set must be split into training and testing data sets, where the former has  $N_{train}$  input/output pairs while the latter has  $N_{test}$  pairs. It is difficult to provide heuristics on what values to set for  $N_{test}$  and  $N_{train}$ , however, we do require that both testing and training sets are distributed with respect to the measure  $\rho(\mathbf{x})$ .

For our algorithm, we require an initial choice of the ridge  $\mathbf{M}$ , which is readily computed by applying a QR factorization on a random normal matrix and then taking the first  $m$  columns of  $\mathbf{Q}$  (see lines 3–5). This trick ensures that the columns of  $\mathbf{M}$  are orthogonal. It also yields samples from the Stiefel manifold with uniform probability. As we do not know which local minima this algorithm will converge to, by sampling with uniform probability we avoid the possibility that the initialization will bias the algorithm towards a particular minima. It also implies that each call of our algorithm will yield a different result, by virtue of the fact that the initial starting point—i.e.,  $\mathbf{M}_0$ —is different. Following this, it is easy to project the  $d$ -dimensional data  $\mathbf{x}$  to the  $m$ -dimensional space (line 7).

---

**Algorithm 1** Gaussian ridge function approximation.

---

**Data:** Input/output pairs  $(\mathbf{x}_j, f_j)$  where  $j = 1, \dots, N_{train} + N_{test}$

**Result:** Matrix  $\mathbf{M}$

1 Split the data into training data set  $(\hat{\mathbf{x}}_i, \hat{f}_i)$  where  $i = 1, \dots, N_{train}$ .

2 Assemble the testing data set  $(\tilde{\mathbf{x}}_i, \tilde{f}_i)$  where  $i = 1, \dots, N_{test}$ .

3 Initialize a random normal matrix  $\mathbf{A} \in \mathbb{R}^{d \times m}$ .

4 Compute its QR factorization  $\mathbf{A} = (\mathbf{Q}_1 \ \mathbf{Q}_2) \mathbf{R}$ , where  $\mathbf{Q}_1 \in \mathbb{R}^{d \times m}$  and  $\mathbf{M} := \mathbf{Q}_1$ .

5 **while**  $|\Delta r| \leq \epsilon$  **do**

6     Compute:

$$\hat{\mathbf{u}}_i = \mathbf{M}^T \hat{\mathbf{x}}_i \quad \text{and} \quad \tilde{\mathbf{u}}_i = \mathbf{M}^T \tilde{\mathbf{x}}_i.$$

7     Solve the optimization problem using the training data  $(\hat{\mathbf{u}}_i, \hat{f}_i)$ :

$$\boldsymbol{\theta}_* = \operatorname{argmax} \log p(\hat{\mathbf{u}}, \boldsymbol{\theta}).$$

8     Solve the Stiefel manifold optimization problem using  $\boldsymbol{\theta}_*$  on the testing data:

$$\mathbf{M}^* = \operatorname{argmin} \frac{1}{2N_{test}} \left\| \tilde{\mathbf{f}} - \tilde{\mathbf{g}}(\mathbf{M}, \boldsymbol{\theta}_*) \right\|_2^2.$$

9     Set:  $\mathbf{M} = \mathbf{M}^*$ .

10 **end**

---

Next, we train a Gaussian process model using  $(\hat{\mathbf{u}}_i, \hat{\mathbf{f}}_i)$  and solve the optimization problem in (35) (line 7) to obtain the hyperparameters  $\boldsymbol{\theta}$ . These values are then fed back into the Gaussian process posterior mean—i.e., our approximation of  $\mathbb{E}(f|\mathbf{u}, \mathbf{M})$ —for prediction. In line 8, we optimize over the Stiefel manifold using the hyperparameters  $\boldsymbol{\theta}$ . This requires discretizing the integral in (38), resulting in

$$(42) \quad \begin{aligned} r &= \frac{1}{2N_{test}} \sum_{i=1}^{N_{test}} \left( \tilde{f}_i - \bar{g}(\mathbf{M}^T \tilde{\mathbf{x}}_i, \boldsymbol{\theta}_*) \right)^2 \\ &= \frac{1}{2N_{test}} \left\| \tilde{\mathbf{f}} - \tilde{\mathbf{g}} \right\|_2^2. \end{aligned}$$

The gradients of this function on the Stiefel manifold can then be obtained using

$$(43) \quad \nabla_{\text{St}r}(\mathbf{M}) = \frac{1}{N_{test}} \sum_{i=1}^{N_{test}} \left( \tilde{f}_i - \bar{g}(\tilde{\mathbf{u}}_i, \boldsymbol{\theta}) \right) \left\{ \frac{\partial \bar{g}(\tilde{\mathbf{u}}_i)}{\partial \mathbf{M}} - \mathbf{M} \left( \frac{\partial \bar{g}(\tilde{\mathbf{u}}_i)}{\partial \mathbf{M}} \right)^T \mathbf{M} \right\},$$

where  $\partial \bar{g}(\tilde{\mathbf{u}}_i) / \partial \mathbf{M}$  is determined by plugging in the testing data into (33). Both the values of  $r$  and  $\nabla_{\text{St}r} \in \mathbb{R}^{d \times m}$  are provided to the manifold optimizer in step 8. This process is then repeated till the difference in residuals  $|\Delta r|$  between successive iterations is below a chosen value  $\epsilon$ .

**4. Numerical studies.** Our codes for Algorithm 1 can be found at <https://www.github.com/~psesh/Gaussian-Ridges>. The codes are in MATLAB and require the `gpml` [32] toolbox for Gaussian processes and `manopt` [4], a set of utilities for manifold optimization. For the manifold optimization subroutine in our codes, we use the preconditioned Riemannian conjugate gradients algorithm of Boumal and Absi (see page 224 of [3]). For optimizing the hyperparameters in our Gaussian process models, we use the in-built (within `gpml`) conjugate gradient approach based on the Polak-Ribière method (see [30]). Initial values for all our hyperparameters were set to 0.

**4.1. An analytical linear ridge.** In this section, we study the outcome of the algorithm provided in subsection 3.3 for discovering the ridge subspace associated with the simple problem

$$(44) \quad \begin{aligned} f(\mathbf{x}) &= g(\mathbf{M}^T \mathbf{x}) \quad \text{where } \mathbf{M} = \begin{bmatrix} \mathbf{m}_1 & \mathbf{m}_2 \end{bmatrix}, \\ \mathbf{M} &\in \mathbb{R}^{d \times 2}, \quad \mathbf{x} \in \mathbb{R}^d \quad \text{and } g: \mathbb{R}^2 \rightarrow \mathbb{R}, \end{aligned}$$

with  $\mathbf{x} \in [-1, 1]^d$ . Here  $g$  is a bivariate function given by

$$(45) \quad g(y_1, y_2) = y_1 + y_2.$$

We explore the result of the conjugate gradient optimizer by varying the number of training and testing points along with the number of dimensions  $d$ . As our optimization problem is non-convex, we expect different outcomes based on the initial value provided to the optimizer. As a result, we run each optimization numerous times—using the same training and testing

data-set. In other words, each optimization trial begins with a random orthogonal matrix (see Step 3 in Algorithm 1) and then aims to minimize our objective function (see Step 8 in Algorithm 1). In the graphs below, we denote our optimized solution as

$$(46) \quad \mathbf{M}^* = \begin{bmatrix} \mathbf{m}_1^* & \mathbf{m}_2^* \end{bmatrix}, \quad \text{where } \mathbf{M}^* \in \mathbb{R}^{d \times 2}.$$

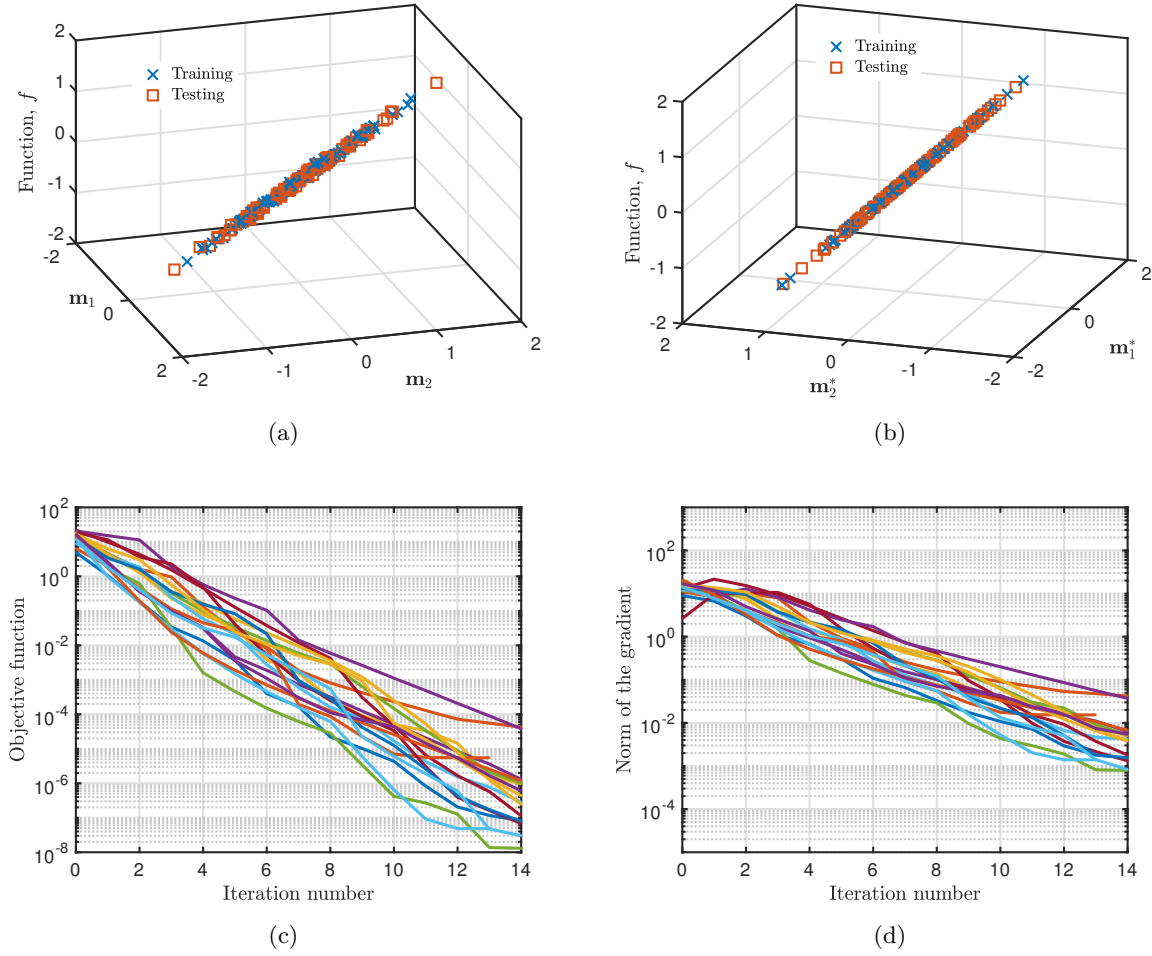
**4.1.1. Case  $d = 10$ .** To begin, consider the isolated case where we set  $N_{test} = 50$  and  $N_{train} = 50$ . Figure 3(a) illustrates the sufficient summary plot on  $\mathbf{M}$ , while in (b) we plot the sufficient summary plot on the optimizer-yielded  $\mathbf{M}^*$ , taken from a trial that yielded the lowest value of the objective function; clearly our algorithm offers a suitable approximation to this ridge function. We repeat this experiment 20 times, each time specifying a different initial orthogonal matrix. In Figure 3(c) we plot the change in the objective function over the number of iterations and in (d) we plot the norm of the gradient per iteration.

As alluded to previously, it is difficult to offer bounds on the number of testing and training samples one should use in Algorithm 3.3. In Figure 3(d) we plot the probability of success for different combination of testing and training samples for this problem. We define an optimization trial to be successful if the objective function attains a value  $\leq 0.005$ . Furthermore, as the outcome of any optimization is dependent upon the initial point, we repeat each experiment with 20 trials and report the average values in the individual circles in Figure 4(a). Circles that are black have, on average, a high chance of recovering the ridge subspace. For this example problem, it is clear that beyond  $N_{test} \geq 50$  and  $N_{train} \geq 50$ , ridge recovery is obtained.

**4.1.2. Case  $d = 100$ .** A similar plot for the case where  $d = 100$  is shown in Figure 4(b). Here we observe that on average, beyond  $N_{train} = 200$  we observe high recovery probabilities. Thus, we repeat the same set of experiments with  $N_{test} = 300$  and  $N_{train} = 300$  and plot the results in Figure 5. Once again our algorithm is able to identify the ridge structure associated with this problem. It should be noted that in even with 300 testing and training samples, we do encounter solutions that do not converge, as shown in Figure 5(c) and (d). Our strategy for finding a suitable solution is to repeat a single experiment—fixing the number of training and testing samples—20 times and then select the solution that yields the lowest value of the objective function. This is one of the key shortcomings of our approach, i.e., the need to run several random trials and then select a solution. Selecting one of the optimization trials with poor convergence results sufficient summary plots of the form Figure 5(e). Here it is clear that the ridge structure is not identified.

**4.2. A bivariate normal distribution.** In this example, we seek to approximate the function

$$(47) \quad f(\mathbf{M}^T \mathbf{x}) = \frac{1}{\sqrt{|\boldsymbol{\Sigma}|} (2\pi)^2} \exp\left(-\frac{1}{2} \mathbf{u} \boldsymbol{\Sigma}^{-1} \mathbf{u}^T\right),$$

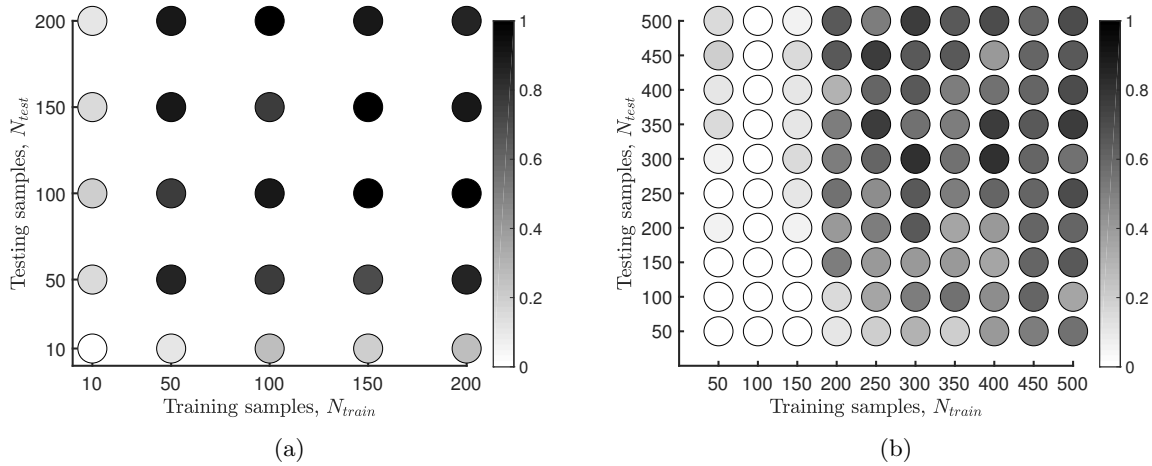


**Figure 3.** Results of the analytical linear ridge problem with  $d = 10$ ,  $N_{test} = 100$  and  $N_{train} = 100$  with 20 repetitions: (a) Sufficient summary plot on  $\mathbf{M}$ ; (b) sufficient summary plot on  $\mathbf{M}^*$ ; (c) objective function vs. iterations; (d) gradient norm vs. iterations.

where

$$(48) \quad \mathbf{M} = \begin{bmatrix} -0.5425 & 0.0654 \\ -0.6784 & -0.4931 \\ -0.2381 & -0.2367 \\ -0.1655 & 0.4643 \\ -0.4017 & 0.6935 \end{bmatrix}, \quad \text{and} \quad \mathbf{\Sigma} = \begin{bmatrix} 0.25 & 0.30 \\ 0.30 & 1.0 \end{bmatrix},$$

where identically distributed independent samples are taken from a multivariate normal distribution with zero mean and a variance of unity, i.e.,  $\mathbf{x} \in \mathcal{N}(0, 1)^5$ . Thus, our 5D function is actually a multivariate Gaussian distribution over a 2D subspace. In Figure 6 we compare the original function in (e) with two different approximations, shown in (a,b) and (c,d). These approximations use a different number of training and testing pairs, starting with 50 in (a, b)



**Figure 4.** Average probability of success values from 20 trials, for different training and testing samples for the case where (a)  $d = 10$ ; (b)  $d = 100$  for the analytical linear ridge problem. Darker colored circles denote a greater probability of success.

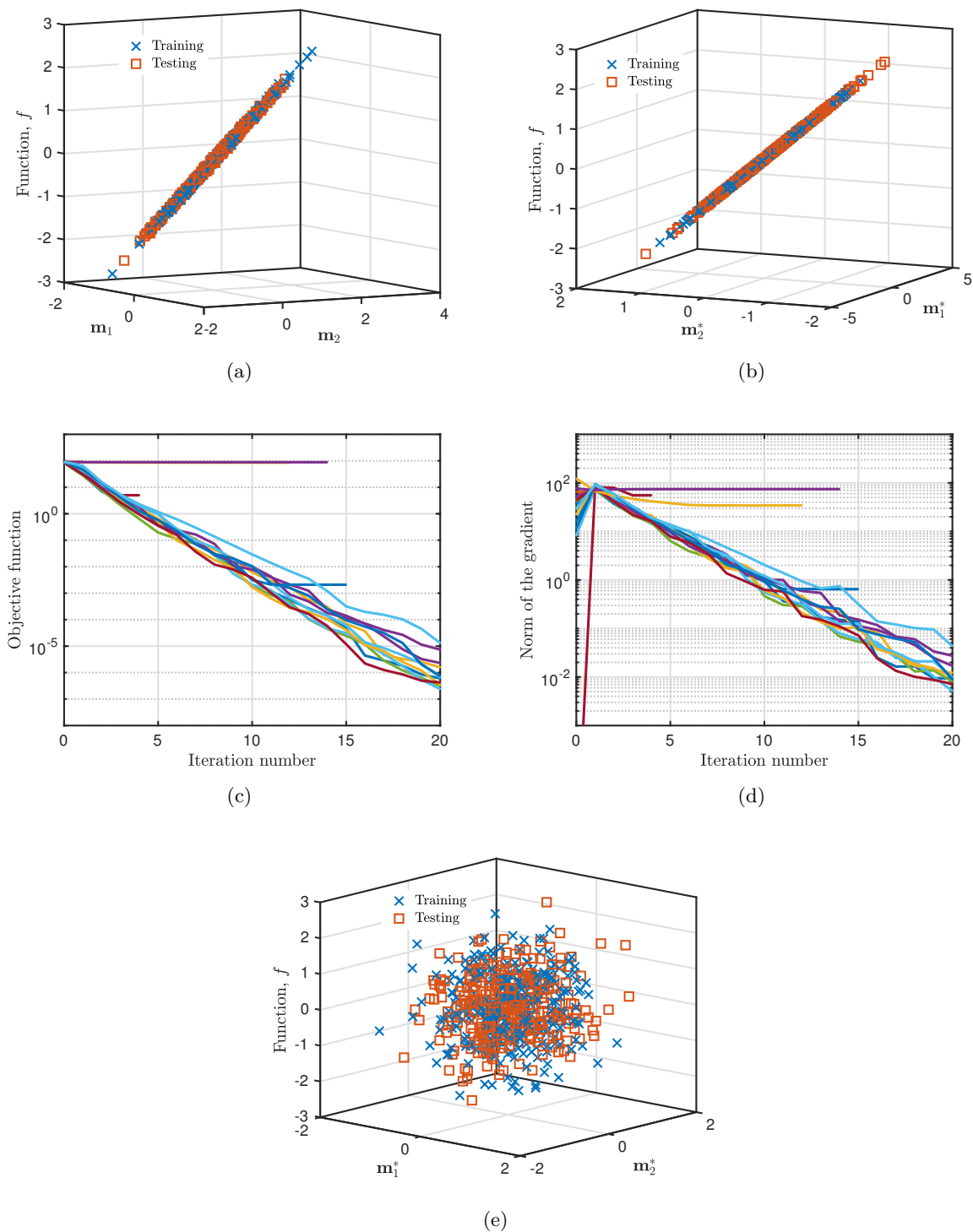
and 100  $N_{test}$  and  $N_{train}$  samples each in (c,d). As before, we examine a family of 20 random trials using the same testing and training input/output pairs, and then select the solution that yields the lowest value of the objective function. Once again, it is clear that our solution strategy is able to roughly identify the ridge structure associated with the problem.

**4.3. Turbomachinery case study.** In this subsection, we use the turbomachinery case study considered in [35] to compare the efficacy of our algorithm with a few sufficient dimension reduction techniques. We make two key claims in this subsection for this turbomachinery case study: (a) the standard deviation obtained from a Gaussian process may be used to make statements on the accuracy associated with the ridge approximation; and (b) our algorithm achieves a lower predictive variance compared to four well-known sufficient dimension reduction methods: SIR, SAVE and CR. To begin, we briefly summarize the turbomachinery case presented here.

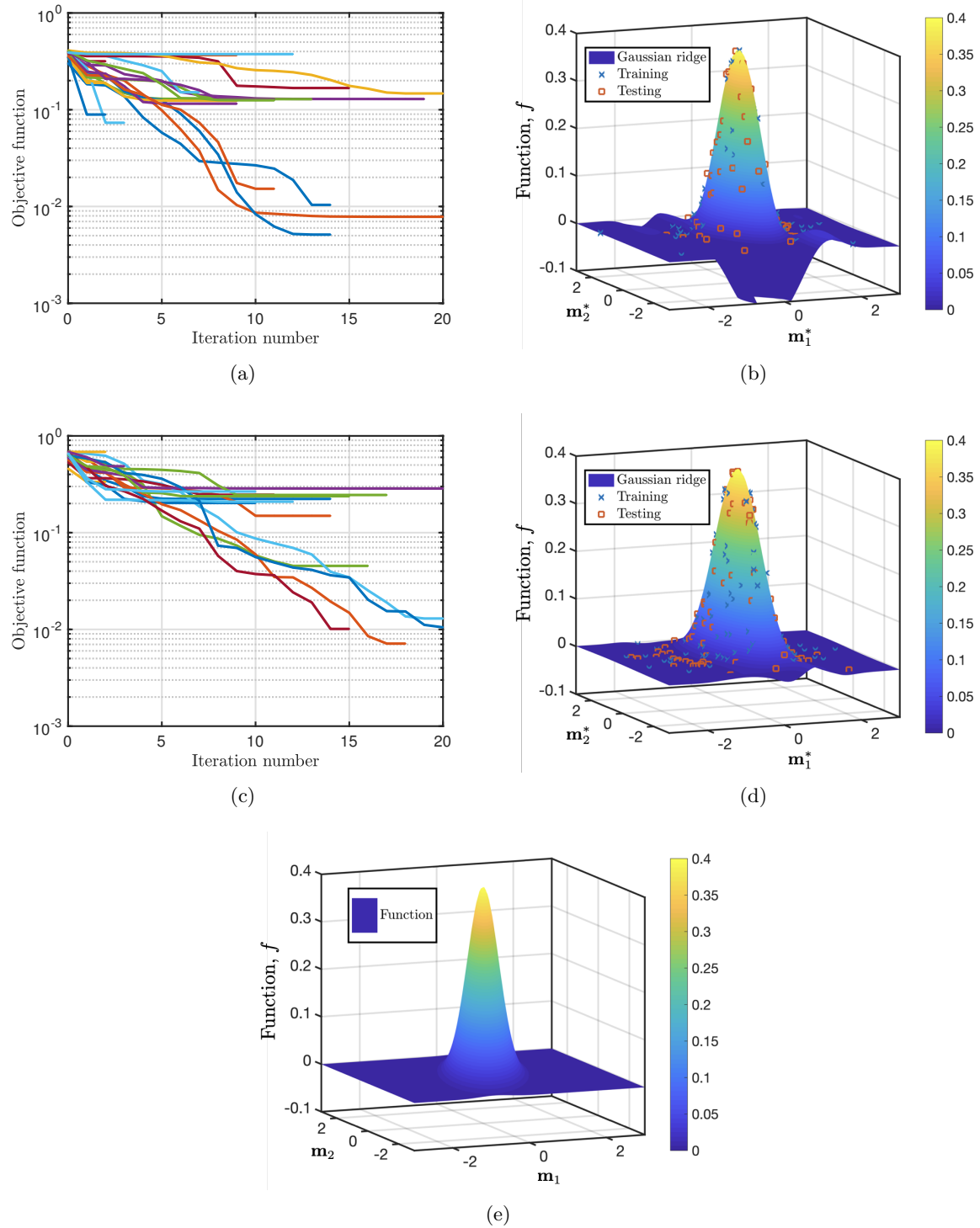
In their work investigating the use of active subspaces for learning turbomachinery aerodynamic pedigree rules of design, the authors of [35] parameterize an aero-engine fan blade with 25 design variables—five degrees of freedom defined at five spanwise locations. Once the geometry is generated, the mesh generation and flow solver codes described in [35] are used to solve the Reynolds average Navier Stokes equations (RANS) for the particular blade design; a post-processing routine is then utilized for estimating the efficiency—our *quantity of interest*. In the aforementioned paper, the authors run a design of experiment (DOE) with 548 different designs for estimating the active subspaces via a global quadratic model, which required estimating 351 polynomial coefficients. In what follows, we demonstrate the efficacy of our algorithm with far fewer samples.

We randomly (uniformly) select  $N = 300$  samples from the above DOE: input/output pairs  $(\mathbf{x}_i \in \mathbb{R}^{25}, f_i \in \mathbb{R})$  for  $i = 1, \dots, N$ . These samples are provided as inputs to SIR, SAVE and CR for estimating a dimension-reducing subspace. The computed subspaces  $\hat{\mathbf{M}} \in \mathbb{R}^{d \times 2}$  are





**Figure 5.** Results of the analytical linear ridge problem with  $d = 100$ ,  $N_{test} = 300$  and  $N_{train} = 300$  with 20 repetitions: (a) Sufficient summary plot on  $\mathbf{M}$ ; (b) sufficient summary plot on  $\mathbf{M}^*$  taken from the trial with the lowest value of the objective function; (c) objective function vs. iterations; (d) gradient norm vs. iterations; (e) Sufficient summary plot on  $\mathbf{M}^*$  taken from the trial with a high value of the objective function.



**Figure 6.** Results of the bivariate normal distribution problem with  $d = 5$  with the true function in (e) and approximations with (a, b)  $N_{test} = 50$ ,  $N_{train} = 50$  and (c, d)  $N_{test} = 100$ ,  $N_{train} = 100$ .

then used for generating sufficient summary plots and for fitting a Gaussian process regression response. To ensure parity, all methods use the same 300 samples.

To ascertain how appropriate the choice of a particular  $\tilde{\mathbf{M}}$  is, we use all 548 DOE points in the sufficient summary plots and in the Gaussian process regression (even though the subspaces  $\tilde{\mathbf{M}}$  are estimated with strictly 300 samples). Figures 7 (a, b and c) are the computed sufficient summary plots for SIR, SAVE and CR. Here the vertical axis represents the normalized efficiency, while the horizontal axis for a particular design  $\mathbf{x}_i$  is given by

$$(49) \quad \mathbf{u}_1 = \left( \tilde{\mathbf{M}}(:, 1) \right)^T \mathbf{x} \quad \text{and} \quad \mathbf{u}_2 = \left( \tilde{\mathbf{M}}(:, 2) \right)^T \mathbf{x},$$

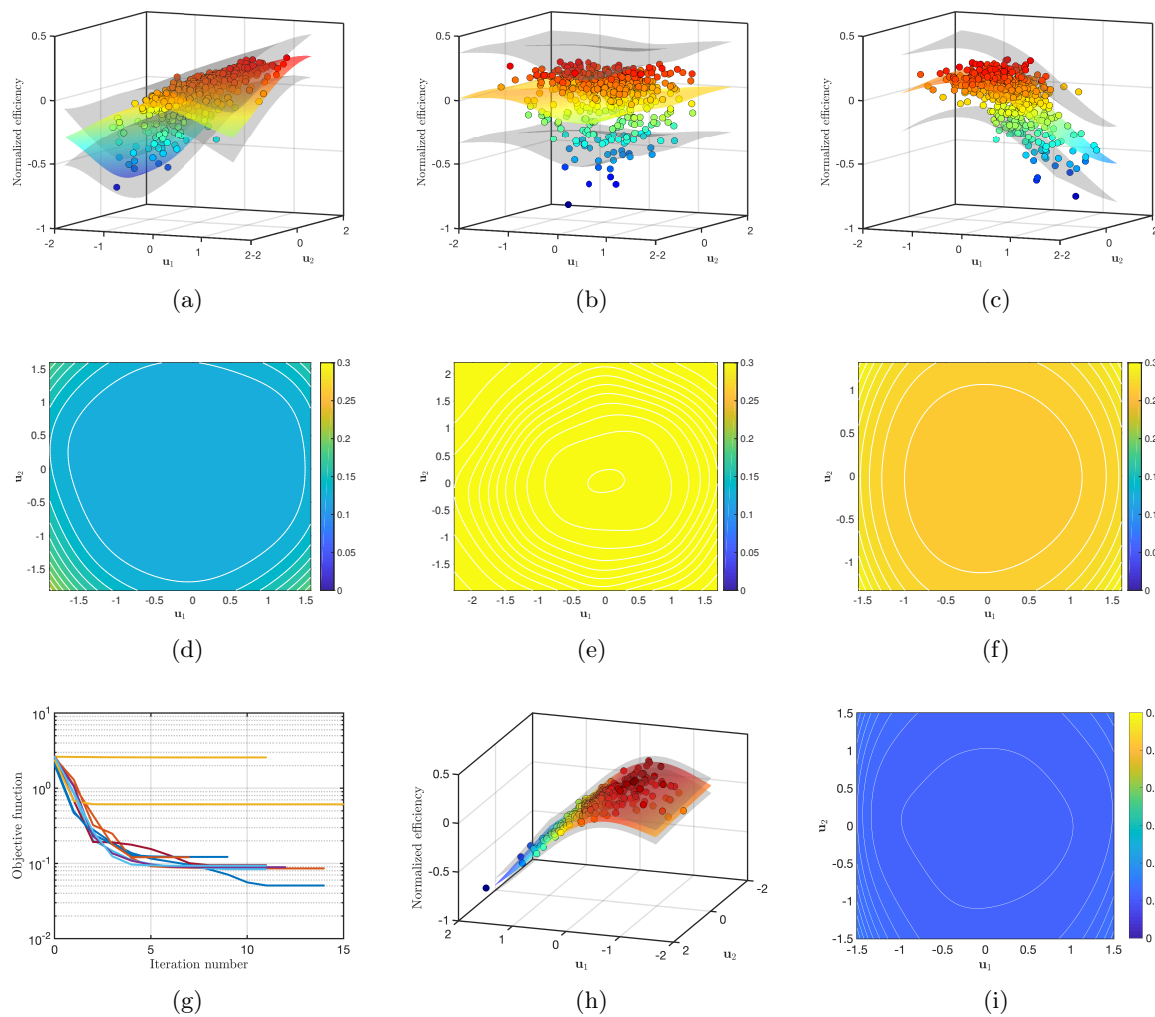
where the notation  $(:, k)$  denotes the  $k$ -th column of  $\tilde{\mathbf{M}}$ . Figures 7 (d, e and f) show contours of the  $2\sigma$  values obtained from the Gaussian process response surface. It is clear from these plots that SIR yields lower  $2\sigma$  values on average. SAVE and CR fail to identify any low-dimensional structure. In contrast, the result of our algorithm is shown in Figure 7(g-i). The  $2\sigma$  contours shown in Figure 7(i) are by far the lowest, and one can clearly see that all the data lies on a quadratic manifold. This result was obtained by selecting the subspace that yielded the lowest value of the objective function after 20 trials (see Figure 7(g)).

**5. Conclusion.** In this paper we introduce a new algorithm for ridge function approximation tailored for subspace-based dimension reduction. Given point evaluations of a model, our algorithm approximates the ridge subspace  $\mathbf{M}$  by minimizing a quadratic form of the function and its best ridge approximation. Throughout this paper, we assume our ridge function is the posterior mean of a Gaussian process. Our results for the examples considered demonstrate that: (i) our algorithm is able to obtain near local minima recovery; and (ii) our algorithm offers better results compared with the four sufficient dimension reduction techniques considered in this paper.

**6. Acknowledgements.** The authors are grateful to Paul Constantine for insightful discussions. The first author would like to acknowledge the support of a Rolls-Royce fellowship. The second author would like to acknowledge the financial support of Magdalene College, Cambridge.

## REFERENCES

- [1] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization Algorithms on Matrix Manifolds*, Princeton University Press, Princeton, NJ, 2009.
- [2] K. P. ADRAGNI AND R. D. COOK, *Sufficient dimension reduction and prediction in regression*, Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, 367 (2009), pp. 4385–4405, <https://doi.org/10.1098/rsta.2009.0110>.
- [3] N. BOUMAL AND P.-A. ABSIL, *Low-rank matrix completion via preconditioned optimization on the grassmann manifold*, Linear Algebra and its Applications, 475 (2015), pp. 200–239.
- [4] N. BOUMAL, B. MISHRA, P.-A. ABSIL, AND R. SEPULCHRE, *Manopt, a matlab toolbox for optimization on manifolds*, Journal of Machine Learning Research, 15 (2014), pp. 1455–1459, <http://jmlr.org/papers/v15/boumal14a.html>.
- [5] A. COHEN, I. DAUBECHIES, R. DEVORE, G. KERKYACHARIAN, AND D. PICARD, *Capturing ridge functions in high dimensions from point queries*, Constructive Approximation, 35 (2012), pp. 225–243, <https://doi.org/10.1007/s00365-011-9147-6>.



**Figure 7.** Turbomachinery case study results for estimating the ridge subspace followed by Gaussian process regression for estimating the mean and  $2\sigma$  response surfaces. Sufficient summary plots for SIR, SAVE and CR are shown in (a, b and c); corresponding  $2\sigma$  contours are shown in (d, e and f). Optimization convergence for 20 random trails for our algorithm is shown in (g), sufficient summary plot in (h) and the corresponding  $2\sigma$  plot in (i).

- [6] P. CONSTANTINE AND D. GLEICH, *Computing active subspaces with Monte Carlo*, arXiv preprint arXiv:1408.0545, (2014).
- [7] P. G. CONSTANTINE, *Active Subspaces: Emerging Ideas for Dimension Reduction in Parameter Studies*, SIAM, Philadelphia, PA, 2015.
- [8] P. G. CONSTANTINE, E. DOW, AND Q. WANG, *Active subspace methods in theory and practice: applications to kriging surfaces*, SIAM Journal on Scientific Computing, 36 (2014), pp. A1500–A1524, <https://doi.org/10.1137/130916138>.
- [9] P. G. CONSTANTINE, A. EFTEKHARI, J. HOKANSON, AND R. A. WARD, *A near-stationary subspace for ridge approximation*, Computer Methods in Applied Mechanics and Engineering, 326 (2017), pp. 402 – 421, <https://doi.org/10.1016/j.cma.2017.07.038>.
- [10] P. G. CONSTANTINE, B. ZAHARATOS, AND M. CAMPANELLI, *Discovering an active subspace in a single-*

- diode solar cell model*, Statistical Analysis and Data Mining: The ASA Data Science Journal, 8 (2015), pp. 264–273, <https://doi.org/10.1002/sam.11281>.
- [11] R. D. COOK, *Save: a method for dimension reduction and graphics in regression*, Communications in Statistics – Theory and Methods, 29 (2000), pp. 2109–2121, <https://doi.org/10.1080/03610920008832598>.
- [12] R. D. COOK, *Rejoinder: Fisher Lecture: Dimension Reduction in Regression*, Statistical Science, 22 (2007), pp. 40–43, <https://doi.org/10.1214/088342307000000078>.
- [13] R. D. COOK, *Regression Graphics: Ideas for Studying Regressions Through Graphics*, vol. 482 of Wiley Series in Probability and Statistics, John Wiley & Sons, New York, 2009.
- [14] R. D. COOK AND L. NI, *Sufficient dimension reduction via inverse regression: A minimum discrepancy approach*, Journal of the American Statistical Association, 100 (2005), pp. 410–428, <https://doi.org/10.1198/016214504000001501>.
- [15] P. DIAZ, P. CONSTANTINE, K. KALMBACH, E. JONES, AND S. PANKAVICH, *A modified SEIR model for the spread of Ebola in Western Africa and metrics for resource allocation*, Applied Mathematics and Computation, 324 (2018), pp. 141 – 155, <https://doi.org/10.1016/j.amc.2017.11.039>.
- [16] A. EDELMAN, T. A. ARIAS, AND S. T. SMITH, *The geometry of algorithms with orthogonality constraints*, SIAM Journal on Matrix Analysis and Applications, 20 (1998), pp. 303–353, <https://doi.org/10.1137/S0895479895290954>.
- [17] M. FORNASIER, K. SCHNASS, AND J. VYBIRAL, *Learning functions of few arbitrary linear parameters in high dimensions*, Foundations of Computational Mathematics, 12 (2012), pp. 229–262, <https://doi.org/10.1007/s10208-012-9115-y>.
- [18] J. H. FRIEDMAN AND W. STUETZLE, *Projection pursuit regression*, Journal of the American statistical Association, 76 (1981), pp. 817–823.
- [19] A. T. GLAWS, P. G. CONSTANTINE, AND R. D. COOK, *Inverse regression for ridge recovery*, arXiv preprint arXiv:1702.02227, (2017).
- [20] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, Johns Hopkins University Press, Baltimore, MD, 4 ed., 2013.
- [21] Z. GREY AND P. CONSTANTINE, *Active subspaces of airfoil shape parameterizations*, in 58th AIAA/ASCE/AHS/ASC Structures, Structural Dynamics, and Materials Conference, American Institute of Aeronautics and Astronautics, 2017, <https://doi.org/10.2514/6.2017-0507>.
- [22] T. HASTIE, R. TIBSHIRANI, AND J. FRIEDMAN, *The Elements of Statistical Learning*, Springer Series in Statistics, Springer-Verlag, New York, 2 ed., 2009.
- [23] B. LI AND S. WANG, *On directional regression for dimension reduction*, Journal of the American Statistical Association, 102 (2007), pp. 997–1008, <https://doi.org/10.1198/016214507000000536>.
- [24] K. C. LI, *Sliced inverse regression for dimension reduction*, Journal of the American Statistical Association, 86 (1991), pp. 316–327, <https://doi.org/10.2307/2290563>.
- [25] X. LIU AND S. GUILLAS, *Dimension reduction for Gaussian process emulation: an application to the influence of bathymetry on tsunami heights*, SIAM/ASA Journal on Uncertainty Quantification, 5 (2017), pp. 787–812, <https://doi.org/10.1137/16M1090648>.
- [26] B. F. LOGAN AND L. A. SHEPP, *Optimal reconstruction of a function from its projections*, Duke Math. J., 42 (1975), pp. 645–659, <https://doi.org/10.1215/S0012-7094-75-04256-8>, <https://doi.org/10.1215/S0012-7094-75-04256-8>.
- [27] Y. MA AND L. ZHU, *A review on dimension reduction*, International Statistical Review, 81 (2013), pp. 134–150, <https://doi.org/10.1111/j.1751-5823.2012.00182.x>.
- [28] L. MAGRI, M. BAUERHEIM, F. NICLOUD, AND M. P. JUNIPER, *Stability analysis of thermo-acoustic nonlinear eigenproblems in annular combustors. Part II. Uncertainty quantification*, Journal of Computational Physics, 325 (2016), pp. 411–421, <https://doi.org/10.1016/j.jcp.2016.08.043>.
- [29] A. PINKUS, *Ridge Functions*, vol. 205 of Cambridge Tracts in Mathematics, Cambridge University Press, Cambridge, UK, 2015, <https://doi.org/10.1017/CBO9781316408124>.
- [30] E. POLAK, *Optimization: Algorithms and Consistent Approximations*, vol. 124 of Applied Mathematical Sciences, Springer-Verlag, New York, 2012.
- [31] F. PUKELSHEIM, *Optimal Design of Experiments*, vol. 50 of Classics in Applied Mathematics, SIAM, Philadelphia, PA, 1993, <https://doi.org/10.1137/1.9780898719109>.
- [32] C. E. RASMUSSEN AND H. NICKISCH, *Gaussian processes for machine learning (GPML) toolbox*, Jour-

- nal of Machine Learning Research, 11 (2010), pp. 3011–3015, <http://www.jmlr.org/papers/v11/rasmussen10a.html>.
- [33] C. E. RASMUSSEN AND C. K. WILLIAMS, *Gaussian Processes for Machine Learning*, MIT Press, Cambridge, MA, 1 ed., 2006, <http://www.gaussianprocess.org/gpml/>.
- [34] A. M. SAMAROV, *Exploring regression structure using nonparametric functional estimation*, Journal of the American Statistical Association, 88 (1993), pp. 836–847, <https://doi.org/10.2307/2290772>.
- [35] P. SESHADRI, S. SHAHPAR, P. CONSTANTINE, G. PARKS, AND M. ADAMS, *Turbomachinery active subspace performance maps*, Journal of Turbomachinery, 140 (2018), pp. 041003–11, <http://dx.doi.org/10.1115/1.4038839>.
- [36] R. TRIPATHY, I. BILIONIS, AND M. GONZALEZ, *Gaussian processes with built-in dimensionality reduction: Applications to high-dimensional uncertainty propagation*, Journal of Computational Physics, 321 (2016), pp. 191–223, <https://doi.org/10.1016/j.jcp.2016.05.039>.
- [37] H. TYAGI AND V. CEVHER, *Learning non-parametric basis independent models from point queries via low-rank methods*, Applied and Computational Harmonic Analysis, 37 (2014), pp. 389–412, <https://doi.org/10.1016/j.acha.2014.01.002>.
- [38] Y. XIA, H. TONG, W. LI, AND L.-X. ZHU, *An adaptive estimation of dimension reduction space*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 64 (2002), pp. 363–410, <https://doi.org/10.1111/1467-9868.03411>.
- [39] L.-P. ZHU AND L.-X. ZHU, *Dimension reduction for conditional variance in regressions*, Statistica Sinica, (2009), pp. 869–883, <http://www.jstor.org/stable/24308860>.