

Compressed Gaussian Process for Manifold Regression

Rajarshi Guhaniyogi

*Department of Applied Mathematics & Statistics
University of California
Santa Cruz, CA 95064, USA*

RGUHANII@UCSC.EDU

David B. Dunson

*Department of Statistical Science
Duke University
Durham, NC 27708-0251, USA*

DUNSON@DUKE.EDU

Editor: Francois Caron

Abstract

Nonparametric regression for large numbers of features (p) is an increasingly important problem. If the sample size n is massive, a common strategy is to partition the feature space, and then separately apply simple models to each partition set. This is not ideal when n is modest relative to p , and we propose an alternative approach relying on random compression of the feature vector combined with Gaussian process regression. The proposed approach is particularly motivated by the setting in which the response is conditionally independent of the features given the projection to a low dimensional manifold. Conditionally on the random compression matrix and a smoothness parameter, the posterior distribution for the regression surface and posterior predictive distributions are available analytically. Running the analysis in parallel for many random compression matrices and smoothness parameters, model averaging is used to combine the results. The algorithm can be implemented rapidly even in very large p and moderately large n nonparametric regression, has strong theoretical justification, and is found to yield state of the art predictive performance.

Keywords: Compressed regression; Gaussian process; Gaussian random projection; Large p ; Manifold regression.

1. Introduction

With recent technological progress, it is now routine in many disciplines to collect data containing large numbers of features, ranging from thousands to millions. To account for complex nonlinear relationships between the features and the response, nonparametric regression models are employed. For example,

$$y = \mu_0(\mathbf{x}) + \epsilon, \quad \epsilon \sim N(0, \sigma^2),$$

where $\mathbf{x} \in \mathcal{R}^p$, $\mu_0(\cdot)$ is the unknown regression function and ϵ is a residual. When p is large, estimating μ_0 can lead to a statistical and computational curse of dimensionality. One strategy for combatting this curse is dimensionality reduction via variable selection or (more broadly) subspace learning, with the high-dimensional features replaced with their projection to a d -dimensional subspace or manifold with $d \ll p$. In many applications,

the relevant information about the high-dimensional features can be encoded in such low dimensional coordinates.

There is a vast frequentist literature on subspace learning for regression, typically employing a two stage approach. In the first stage, a dimensionality reduction technique is used to obtain lower dimensional features that can “faithfully” represent the higher dimensional features. Examples include principal components analysis and more elaborate methods that accommodate non-linear subspaces, such as isomap (Tenenbaum et al., 2000) and Laplacian eigenmaps (Belkin and Niyogi, 2003; Guerrero et al., 2011). Once lower dimensional features are obtained, the second stage uses these features in standard regression and classification procedures as if they were observed initially. Such two stage approaches rely on learning the manifold structure embedded in the high dimensional features, which adds unnecessary computational burden when inferential interest lies mainly in prediction.

Another thread of research focuses on prediction using divide-and-conquer techniques. As the number of features increases, the problem of finding the best splitting attribute becomes intractable, so that CART (Breiman et al., 1984), MARS and multiple tree models, such as Random Forest (Breiman, 2001), cannot be efficiently applied. A much simpler approach is to apply high dimensional clustering techniques, such as metis, cover trees and spectral clustering. Once the observations are clustered into a few groups, simple models (glm, Lasso etc) are fitted in each cluster (Zhang et al., 2013). Such methods are sensitive to clustering, do not characterize predictive uncertainty, and may lack efficiency, an important consideration outside the $n \gg p$ setting. There is also a recent literature on scaling up sparse optimization methods, such as Lasso, to large p and n settings relying on algorithms that can exploit multiple processors in a distributed manner e.g., (Boyd et al., 2011). However, such methods are yet to be developed for non-linear manifold regression, which is the central focus of this article.

This naturally motivates Bayesian models that simultaneously learn the mapping to the lower-dimensional subspace along with the regression function in the coordinates on this subspace, providing a characterization of predictive uncertainties. Tokdar et al. (2010) proposes a logistic Gaussian process approach, while Reich et al. (2011) use finite mixture models for sufficient dimension reduction. Page et al. (2013) propose a Bayesian nonparametric model for learning of an affine subspace in classification problems. These approaches have the disadvantages of being limited to linear subspaces, lacking scalability beyond a few dozen features and having potential sensitivity to features corrupted with noise. There is also a literature on Bayesian methods that accommodate non-linear subspaces, ranging from Gaussian process latent variable models (GP-LVMs) (Lawrence, 2005) for probabilistic nonlinear PCA to mixture factor models Chen et al. (2010). However, such methods similarly face barriers in scaling up to large p and/or n . There is a heavy computational price for learning the number of latent variables, the distribution of the latent variables, and the mapping functions while maintaining identifiability restrictions.

Recently, Yang and Dunson (2013) show that this computational burden can be largely bypassed by using usual Gaussian process (GP) regression without attempting to learn the mapping to the lower-dimensional subspace. They showed that when the features lie on a d -dimensional manifold embedded in the p -dimensional feature space with $d \ll p$ and the regression function is not highly smooth, the optimal rate can be obtained using GP regression with a squared exponential covariance in the original high-dimensional feature

space. This is an exciting theoretical result, which provides motivation for the approach in this article, which is focused on scalable Bayesian nonparametric regression in large p settings. For broader applicability than Yang and Dunson (2013), we accommodate features that are contaminated by noise and hence do not lie exactly on a low-dimensional manifold. In addition, we facilitate computational efficiency by bypassing MCMC and reducing matrix inversion bottlenecks via random projections. Sensitivity to the random projection and to tuning parameters is reduced through the use of Bayesian model averaging. The proposed approach that accommodates all these features is coined as the *compressed Gaussian process* (CGP).

Snelson and Ghahramani (2012) also considered manifold regression for big data, comprising feature vectors via pre-multiplying with a short and fat projection matrix. Their approach involves estimating a total of $(M + p)m$ parameters in a feature compression matrix and input points, with M the number of input points, leading to intractability as p increases. We demonstrate substantial advantages of our random compression approach in Section 5 in terms of computational scalability and predictive performance. In addition, SG lacks theory guarantees, while we show that CGP has a minimax optimal adaptive convergence rate dependent only on the true manifold dimension (assumed small). Calandra et al. (2014) instead use a neural network-like mapping of the input space, requiring non-convex optimization in high-dimensions. Scaling to moderate n , such as $n \sim 5,000 - 10,000$, is problematic. Other manifold regression methods (see Bickel and Li, 2007; Aswani et al., 2011) either lack scalability even for moderate p and n , or fail to characterize predictive uncertainties.

Section 2 proposes the model and computational approach in large p settings. Section 3 describes extensions to moderately large n , and Section 4 develops theoretical justification. Section 5 contains simulation examples relative to state-of-the-art competitors. Section 6 presents an image data application, and Section 7 concludes the paper with a discussion.

2. Compressed Gaussian process regression

This section details out compressed Gaussian process model with the associated prior and posterior distributions of the parameters.

2.1 Model

For subjects $i = 1, \dots, n$, let $y_i \in \mathcal{Y}$ denote a response with associated features $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})' = (z_{i1}, \dots, z_{ip})' + (\delta_{i1}, \dots, \delta_{ip})' = \mathbf{z}_i + \boldsymbol{\delta}_i$, $\mathbf{z}_i \in \mathcal{M}$, $\boldsymbol{\delta}_i \in \mathcal{R}^p$, where \mathcal{M} is a d -dimensional manifold embedded in the ambient space \mathcal{R}^p . We assume that the response $y \in \mathcal{Y}$ is continuous. The measured features do not fall exactly on the manifold \mathcal{M} but are corrupted by noise. We assume a compressed nonparametric regression model

$$y_i = \mu(\boldsymbol{\Psi} \mathbf{x}_i) + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2), \quad (1)$$

with the residuals modeled as Gaussian with variance σ^2 , though other distributions including heavy-tailed ones can be accommodated. $\boldsymbol{\Psi}$ is an $m \times p$ matrix that compresses p -dimensional features to dimension m . Following a Bayesian approach, we choose a prior distribution for the regression function μ and residual variance σ^2 , while randomly generating $\boldsymbol{\Psi}$ following precedence in the literature on feature compression (Maillard and Munos,

2009; Fard et al., 2012; Guhaniyogi and Dunson, 2013). These earlier approaches differ from ours in focusing on parametric regression. We independently draw elements $\{\Psi_{ij}\}$ of Ψ from $N(0, 1)$, and then normalize the rows using Gram-Schmidt orthogonalization.

We assume that $\mu \in \mathcal{H}_s$ is a continuous function belonging to \mathcal{H}_s , a Holder class with smoothness s . To allow μ to be unknown, we use a Gaussian process (GP) prior, $\mu \sim \text{GP}(0, \sigma^2 \kappa)$ with the covariance function chosen to be squared exponential

$$\kappa(\mathbf{x}_i, \mathbf{x}_j; \lambda) = \exp(-\lambda \|\mathbf{x}_i - \mathbf{x}_j\|^2), \quad (2)$$

with λ a smoothness parameter and $\|\cdot\|^2$ the Euclidean norm. To additionally allow the residual variance σ^2 and smoothness λ to be unknown, we let

$$\sigma^2 \sim \text{IG}(a, b), \quad \lambda^d \sim \text{Ga}(a_0, b_0),$$

with $\text{IG}()$ and $\text{Ga}()$ denoting the inverse-gamma and gamma densities, respectively. The powered gamma prior for λ is motivated by the result of van der Vaart and van Zanten (2009) showing minimax adaptive rates of $n^{-s/(2s+p)}$ for a GP prior with squared exponential covariance and powered gamma prior. This is the optimal rate for nonparametric regression in the original p -dimensional ambient space. The rate can be improved to $n^{-s/(2s+d)}$ when $\mathbf{x}_i \in \mathcal{M}$, with \mathcal{M} a d -dimensional manifold. Yang and Dunson (2013) shows that a GP prior with powered gamma prior on the smoothness can achieve this rate. In practice, replacing the powered gamma prior for λ with a gamma prior has essentially no impact on the results in examples we have considered.

In many applications, features may not lie exactly on \mathcal{M} due to noise and corruption in the data. We apply random compression in (1) to de-noise the features, obtaining $\Psi \mathbf{x}_i$ much more concentrated near a lower-dimensional subspace than the original \mathbf{x}_i . With this enhanced concentration, the theory in Yang and Dunson (2013) suggests excellent performance for an appropriate GP prior. In addition to de-noising, this approach has the major advantage of bypassing estimation of a geodesic distance along the unknown manifold \mathcal{M} between any two data points \mathbf{x}_i and $\mathbf{x}_{i'}$.

2.2 Posterior form

Let $\boldsymbol{\mu} = (\mu(\Psi \mathbf{x}_1), \dots, \mu(\Psi \mathbf{x}_n))'$ and $\mathbf{K}_1 = (\kappa(\Psi \mathbf{x}_i, \Psi \mathbf{x}_j; \lambda))_{i,j=1}^n$. The prior distribution on $\boldsymbol{\mu}, \sigma^2$ induces a normal-inverse gamma (NIG) prior on $(\boldsymbol{\mu}, \sigma^2)$,

$$(\boldsymbol{\mu} | \sigma^2) \sim N(\mathbf{0}, \sigma^2 \mathbf{K}_1), \quad \sigma^2 \sim \text{IG}(a, b),$$

leading to a NIG posterior distribution for $(\boldsymbol{\mu}, \sigma^2)$ given $\mathbf{y}, \Psi \mathbf{x}, \lambda$. In the special case in which $a, b \rightarrow 0$, we obtain Jeffrey's prior and the posterior distribution is

$$\boldsymbol{\mu} | \mathbf{y} \sim t_\nu(\mathbf{m}, \boldsymbol{\Sigma}) \quad (3)$$

$$\sigma^2 | \mathbf{y} \sim \text{IG}(a_1, b_1), \quad (4)$$

where $a_1 = n/2$, $b_1 = \mathbf{y}'(\mathbf{K}_1 + \mathbf{I})^{-1} \mathbf{y}/2$, $\mathbf{m} = [\mathbf{I} + \mathbf{K}_1^{-1}]^{-1} \mathbf{y}$, $\boldsymbol{\Sigma} = (2b_1/n) [\mathbf{I} + \mathbf{K}_1^{-1}]^{-1}$, and $t_\nu(\mathbf{m}, \boldsymbol{\Sigma})$ denotes a multivariate- t distribution with ν degrees of freedom, mean \mathbf{m} and covariance $\boldsymbol{\Sigma}$.

Hence, the exact posterior distribution of $(\boldsymbol{\mu}, \sigma^2)$ conditionally on $(\boldsymbol{\Psi}, \lambda)$ is available analytically. The predictive of $\mathbf{y}^* = (y_1^*, \dots, y_{n_{pred}}^*)'$ given $\mathbf{X}^* = (\mathbf{x}_1^{*'}, \dots, \mathbf{x}_{n_{pred}}^{*'})'$ and $\boldsymbol{\Psi}, \lambda$ for new n_{pred} subjects marginalizing out $(\boldsymbol{\mu}, \sigma^2)$ over their posterior distribution is available analytically as

$$\mathbf{y}^* | \mathbf{x}_1^*, \dots, \mathbf{x}_{n_{pred}}^*, \mathbf{y} \sim t_n(\mu_{pred}, \sigma_{pred}^2), \quad (5)$$

where $\mathbf{K}_{pred} = \{\kappa(\mathbf{x}_i^*, \mathbf{x}_j^*; \lambda)\}_{i,j=1}^{n_{pred}}$, $\mathbf{K}_{pred,1} = \{\kappa(\mathbf{x}_i^*, \mathbf{x}_j; \lambda)\}_{i=1, j=1}^{i=n_{pred}, j=n}$, $\mathbf{K}_{1,pred} = \mathbf{K}_{pred,1}'$, $\mu_{pred} = \mathbf{K}_{pred,1}(\mathbf{I} + \mathbf{K}_1)^{-1} \mathbf{y}$, $\sigma_{pred}^2 = (2b_1/n) \left[\mathbf{I} + \mathbf{K}_{pred} - \mathbf{K}_{pred,1} \{\mathbf{I} + \mathbf{K}_1\}^{-1} \mathbf{K}_{1,pred} \right]$.

2.3 Model averaging

The approach described in the previous section can be used to obtain a posterior distribution for $\boldsymbol{\mu}$ and a predictive distribution for $\mathbf{y}^* = (y_1^*, \dots, y_{n_{pred}}^*)'$ given \mathbf{X}^* for a new set of n_{pred} subjects *conditionally* on the $m \times p$ random projection matrix $\boldsymbol{\Psi}$ and the scaling parameter λ . To accomplish robustness with respect to the choice of $(\boldsymbol{\Psi}, \lambda)$ and the subspace dimension m , following Guhaniyogi and Dunson (2013), we propose to generate s random matrices having different m , s and λ from the marginal posterior distribution, $(\boldsymbol{\Psi}^{(l)}, \lambda^{(l)})$, $l = 1, \dots, s$, and then use model averaging to combine the results. To make matters more clear, let \mathcal{M}_l , $l = 1, \dots, s$, represent (1) with m_l number of rows. Corresponding to the model \mathcal{M}_l , we denote $\boldsymbol{\Psi}, \lambda, \boldsymbol{\mu}$ and σ^2 by $\boldsymbol{\Psi}^{(l)}, \lambda^{(l)}, \boldsymbol{\mu}^{(l)}$ and $\sigma^{2(l)}$ respectively. Given $\boldsymbol{\Psi}^{(l)}$, we draw a few $\lambda_1, \dots, \lambda_k$ randomly from $U(3/dmax, 3/dmin)$ where $dmax = \max_{i,j} \|\mathbf{x}_i - \mathbf{x}_j\|^2$ and $dmin = \min_{i,j} \|\mathbf{x}_i - \mathbf{x}_j\|^2$. Next we use the fact that the marginal posterior distribution of $\lambda | \boldsymbol{\Psi}^{(l)}, \mathbf{y}$ is given by

$$f(\lambda | \mathbf{y}, \boldsymbol{\Psi}^{(l)}) \propto \frac{1}{|\mathbf{K}_1 + \mathbf{I}|^{\frac{1}{2}}} \frac{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}{\left[\mathbf{y}' (\mathbf{K}_1 + \mathbf{I})^{-1} \mathbf{y} \right]^{\frac{n}{2}} (\sqrt{2\pi})^n} \times \pi(\lambda),$$

where $\pi(\lambda)$ is the prior distribution of λ . Clearly, a discrete approximation of $\lambda | \boldsymbol{\Psi}^{(l)}, \mathbf{y}$ is given by $\sum_{i=1}^k w_i \delta_{\lambda_i}$, where $w_i = \frac{f(\lambda_i | \mathbf{y}, \boldsymbol{\Psi}^{(l)})}{\sum_{j=1}^k f(\lambda_j | \mathbf{y}, \boldsymbol{\Psi}^{(l)})}$ and δ_{λ_i} is the Dirac Delta function at λ_i .

Finally, $\lambda^{(l)}$ is drawn from $\sum_{i=1}^k w_i \delta_{\lambda_i}$. Although Section 4 shows minimax optimality of CGP with $\lambda^d \sim \text{Gamma}(a, b)$, we use $d = 1$ in practical implementations with no practical loss in cases we have considered.

Let $\mathcal{M} = \{\mathcal{M}_1, \dots, \mathcal{M}_s\}$ denote the set of models corresponding to different random projections, $\mathcal{D} = \{(y_i, \mathbf{x}_i), i = 1, \dots, n\}$ denote the observed data, and \mathbf{y}^* denote the data for future subjects with features \mathbf{X}^* . Then, the predictive density of \mathbf{y}^* given \mathbf{X}^* is

$$f(\mathbf{y}^* | \mathbf{X}^*, \mathcal{D}) = \sum_{l=1}^s f(\mathbf{y}^* | \mathbf{X}^*, \mathcal{M}_l, \mathcal{D}) P(\mathcal{M}_l | \mathcal{D}), \quad (6)$$

where the predictive density of \mathbf{y}^* given \mathbf{X}^* under projection \mathcal{M}_l is given in (9) and the posterior probability weight on projection \mathcal{M}_l is

$$P(\mathcal{M}_l | \mathcal{D}) = \frac{P(\mathcal{D} | \mathcal{M}_l) P(\mathcal{M}_l)}{\sum_{h=1}^s P(\mathcal{D} | \mathcal{M}_h) P(\mathcal{M}_h)}.$$

Assuming equal prior weights for each random projection, $P(\mathcal{M}_l) = 1/s$. In addition, the marginal likelihood under \mathcal{M}_l is

$$P(\mathcal{D} | \mathcal{M}_l) = \int P(\mathcal{D} | \mathcal{M}_l, \boldsymbol{\mu}^{(l)}, \sigma^{2(l)}) \pi(\boldsymbol{\mu}^{(l)}, \sigma^{2(l)}). \quad (7)$$

After a little algebra, one observes that for (1) with $(\boldsymbol{\mu} | \sigma^2) \sim N(\mathbf{0}, \sigma^2 \mathbf{K}_1)$, $\pi(\sigma^2) \propto \frac{1}{\sigma^2}$,

$$P(\mathcal{D} | \mathcal{M}_l) = \frac{1}{|\mathbf{K}_1 + \mathbf{I}|^{\frac{1}{2}}} \frac{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}{\left[\mathbf{y}' (\mathbf{K}_1 + \mathbf{I})^{-1} \mathbf{y} \right]^{\frac{n}{2}} (\sqrt{2\pi})^n}.$$

Plugging in the above expressions in (6), one obtains the posterior predictive distribution as a weighted average of t densities. Given that the computation over different sets of $\boldsymbol{\Psi}$, λ are not dependent on each other, the calculations are embarrassingly parallel with a trivial expense for combining. The main computational expense comes from the inversion of an $n \times n$ matrix under the l th random projection. There is a vast literature on obtaining rapid approximations to such inversions under low rank assumptions. In the next section, we describe one such approach for enabling scaling to moderate n . Other recent methods can be easily substituted to scale to very large or massive n .

3. Scaling to moderately large n

Fitting (1) using model averaging requires computing inverses and determinants of covariance matrices of the order $n \times n$. In problems with even moderate n , this adds a heavy computational burden of the order of $O(n^3)$. Additionally, as dimension increases, matrix inversion becomes more unstable with the propagation of errors due to finite machine precision. This problem is further exacerbated if the covariance matrix is nearly rank deficient.

To address such issues, existing solutions rely on approximating $\mu(\cdot)$ by another process $\tilde{\mu}(\cdot)$, which is more tractable computationally. One popular approach constructs $\tilde{\mu}(\cdot)$ as a finite basis approximation via kernel convolution (Higdon, 2002) or kalman filtering (Wikle and Cressie, 1999). Alternatively, one can let $\tilde{\mu}(\cdot) = \mu(\cdot)\eta(\cdot)$, where $\eta(\cdot)$ is a Gaussian process having compactly supported correlation function that essentially makes the covariance matrix of $(\tilde{\mu}(\mathbf{x}_1), \dots, \tilde{\mu}(\mathbf{x}_n))$ sparse (Kaufman et al., 2008), facilitating inversion through efficient sparse solvers.

Banerjee et al. (2008) proposes a low rank approach that imputes $\mu(\cdot)$ conditionally on a few knot-points, closely related to subset of regressor methods in machine learning (Smola and Schölkopf, 2000). Subsequently, Finley et al. (2009) in statistics and Snelson and Ghahramani (2006) in machine learning report bias in both variance and length-scale parameter estimation which affects predictive estimates for the proposed approaches (Banerjee et al., 2008; Smola and Schölkopf, 2000). They also suggest possible remedies for bias adjustments. To avoid sensitivity to knot selection in the low rank approaches, Banerjee et al. (2013) approximates $\mu(\cdot)$ using $\tilde{\mu}(\cdot) = E[\mu(\cdot) | \boldsymbol{\Phi} \boldsymbol{\mu}(\mathbf{X})] + \epsilon_{\boldsymbol{\Phi}}(\cdot)$, with $\boldsymbol{\Phi}$ an $m \times n$, $m \ll n$ random matrix with $\Phi_{ij} \sim N(0, 1)$. $\epsilon_{\boldsymbol{\Phi}}(\mathbf{x})$ are independent feature specific noises with $\epsilon_{\boldsymbol{\Phi}}(\mathbf{x}) \sim N(\mathbf{0}, \text{var}(\mu(\mathbf{x})) - \text{var}(\tilde{\mu}(\mathbf{x})))$, which are introduced for bias correction similar to Finley et al. (2009). There is a parallel literature on nearest neighbor Gaussian processes which is built upon approximating a multivariate high dimensional Gaussian distribution

by a product of lower dimensional conditional distributions. Such an idea was first pursued by Vecchia (1988) and Stein et al. (2004), and has recently gained traction in the computer experiments literature (Gramacy and Apley, 2015) and in spatial geo-statistics (Emery, 2009; Stroud et al., 2014; Datta et al., 2014). Some of the recent versions of the this idea are found to be amenable to parallel computations as well.

We adapt Banerjee et al. (2013) from usual GP regression to our compressed manifold regression setting. In particular, let

$$\mathbf{y} = \tilde{\mu}_{\Phi}(\Psi\mathbf{x}) + \epsilon_{\Phi}(\Psi\mathbf{x}) + \epsilon, \quad \epsilon \sim N(0, \sigma^2), \quad (8)$$

where $\tilde{\mu}_{\Phi}(\Psi\mathbf{x}) = E[\mu(\Psi\mathbf{x}) | \Phi\mu(\mathbf{X}\Psi')]$, $\epsilon_{\Phi}(\Psi\mathbf{x}) | \sigma^2 \sim N(0, \sigma_{\epsilon}^2(\mathbf{x}))$,

$\sigma_{\epsilon}^2(\mathbf{x}) = \sigma^2 \left[\kappa(\Psi\mathbf{x}, \Psi\mathbf{x}; \lambda) - (\Phi\mathbf{k}_{\mathbf{x}})' \{ \Phi\mathbf{K}_1\Phi' \}^{-1} (\Phi\mathbf{k}_{\mathbf{x}}) \right]$ and

$\mathbf{k}_{\mathbf{x}} = (\kappa(\Psi\mathbf{x}, \Psi\mathbf{x}_1; \lambda), \dots, \kappa(\Psi\mathbf{x}, \Psi\mathbf{x}_n; \lambda))'$. Denoting $\mathbf{H}_1 = \text{diag}(\mathbf{K}_1 - \mathbf{K}_1\Phi'(\Phi\mathbf{K}_1\Phi')^{-1}\Phi\mathbf{K}_1) + \mathbf{I}$ and $\mathbf{H}_2 = \mathbf{K}_1\Phi'(\Phi\mathbf{K}_1\Phi')^{-1}\Phi$, marginal posterior distributions of μ and σ^2 are available in analytical forms

$$\mu | \mathbf{y} \sim t_n(\mathbf{m}_{RGP}, \Sigma_{RGP}), \quad \sigma^2 | \mathbf{y} \sim IG(a_2, b_2),$$

where $a_2 = n/2$, $b_2 = \mathbf{y}'(\mathbf{H}_1 + \mathbf{H}_2\mathbf{K}_1)^{-1}\mathbf{y}/2$, $\mathbf{m}_{RGP} = [\mathbf{H}_2'\mathbf{H}_1^{-1}\mathbf{H}_2 + \mathbf{K}_1^{-1}]^{-1}\mathbf{H}_2'\mathbf{H}_1^{-1}\mathbf{y}$, $\Sigma_{RGP} = (2b_2/n)[\mathbf{H}_2'\mathbf{H}_1^{-1}\mathbf{H}_2 + \mathbf{K}_1^{-1}]^{-1}$. Owing to the special structure of Σ_{RGP} and \mathbf{m}_{RGP} , $n \times n$ matrix inversion can be efficiently achieved by Sherman-Woodbury-Morrison matrix inversion technique.

Attention now turns to prediction from (8). The predictive of $\mathbf{y}^* = (y_1^*, \dots, y_{n_{pred}}^*)'$ given $\mathbf{X}^* = (\mathbf{x}_1^*, \dots, \mathbf{x}_{n_{pred}}^*)'$ and Ψ, λ for new n_{pred} subjects marginalizing out (μ, σ^2) over their posterior distribution is available analytically as

$$\mathbf{y}^* | \mathbf{x}_1^*, \dots, \mathbf{x}_{n_{pred}}^*, \mathbf{y} \sim t_n(\mu_{pred}, \sigma_{pred}^2), \quad (9)$$

where $\mathbf{K}_{pred} = \{\kappa(\mathbf{x}_i^*, \mathbf{x}_j^*; \lambda)\}_{i,j=1}^{n_{pred}}$, $\mathbf{K}_{pred,1} = \{\kappa(\mathbf{x}_i^*, \mathbf{x}_j; \lambda)\}_{i=1, j=1}^{i=n_{pred}, j=n}$, $\mathbf{K}_{1,pred} = \mathbf{K}'_{pred,1}$, $\mathbf{H}_3 = \mathbf{I} + \text{diag}(\mathbf{K}_{pred} - \mathbf{K}_{pred,1}\Phi'(\Phi\mathbf{K}_1\Phi)^{-1}\Phi\mathbf{K}_{1,pred})$, $\mu_{pred} = \mathbf{K}_{pred,1}\mathbf{H}_2'(\mathbf{H}_1 + \mathbf{H}_2\mathbf{K}_1)^{-1}\mathbf{y}$, $\sigma_{pred}^2 = (2b_1/n)[\mathbf{H}_3 + \mathbf{K}_{pred,1}\Phi'(\Phi\mathbf{K}_1\Phi')^{-1}\Phi\mathbf{K}_{1,pred} - \mathbf{K}_{pred,1}\mathbf{H}_2'(\mathbf{H}_1 + \mathbf{H}_2\mathbf{K}_1)^{-1}\mathbf{H}_2\mathbf{K}_{1,pred}]$. Evaluating the above expression requires inverting matrices of order $m_{\Phi} \times m_{\Phi}$. Model averaging is again employed to limit sensitivity over the choices of Ψ, λ . Following similar calculations as in Section 2.3, model averaging weights are found to be

$$P(\mathcal{D} | \mathcal{M}_l) = \frac{1}{|\mathbf{H}_2\mathbf{K}_1 + \mathbf{H}_1|^{\frac{1}{2}}} \frac{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}{\left[\mathbf{y}'(\mathbf{H}_2\mathbf{K}_1 + \mathbf{H}_1)^{-1}\mathbf{y} \right]^{\frac{n}{2}} (\sqrt{2\pi})^n}.$$

Model averaging is performed on a wide interval of possible m values determined by the ‘‘compressed sample size’’ m_{Φ} and p , analogous to Section 2.3.

Although we focus in this article on using the Banerjee et al. (2013) approach within CGP for scaling to moderately large n , alternative low rank or scalable approximations to Gaussian processes can be substituted essentially without complication. For example, there has been a recent emphasis on methods that break the data into exhaustive and mutually

exclusive subsets (Parikh and Boyd, 2011), run computation separately for each subset and then combine the results; such methods have been applied to GPs (Deisenroth and Ng, 2015) and have complexity that scales as $O\left(\left(\frac{n}{K}\right)^3\right)$, where K is the number of subsets. Choosing K large enough, with this approach, one can compute CGP with moderate sized subset in each processor followed by combining inferences from different subsets. This can be further reduced by using low rank approximations to the GPs within each subset.

An important question that remains is how much information is lost in compressing the high-dimensional feature vector to a much lower dimension? In particular, one would expect to pay a price for the huge computational gains in terms of predictive performance or other metrics. We address this question in two ways. First we argue satisfactory theoretical performance in prediction in a large p asymptotic paradigm in Section 4. Then, we will consider practical performance in finite samples using simulated and real data sets.

4. Convergence analysis

This section provides theory supporting the excellent practical performance of the proposed method. In our context the feature vector \mathbf{x} is assumed to be $\mathbf{x} = \mathbf{z} + \boldsymbol{\delta}$, $\mathbf{z} \in \mathcal{M}$, $\boldsymbol{\delta} \in \mathcal{R}^p$. Compressing the feature vector results in compressing \mathbf{z} and the noise followed by their addition, $\boldsymbol{\Psi}\mathbf{x} = \boldsymbol{\Psi}\mathbf{z} + \boldsymbol{\Psi}\boldsymbol{\delta}$. The following two directions are used to argue that compression results in near optimal inference.

- (A) When features lie on a manifold a two stage estimation procedure (compression followed by a Gaussian process regression) leads to optimal convergence properties. This is used to show that using $\{\boldsymbol{\Psi}\mathbf{z}_i\}_{i=1}^n$ as features in the Gaussian process regression yields the optimal rate of convergence.
- (B) Noise compression through $\boldsymbol{\Psi}$ mitigates the deleterious effect of noise in \mathbf{x} on the resulting performance.

Let $\mu_0(\cdot)$ and $\mu(\cdot)$ be the true and the fitted regression functions respectively. Define $\rho(\mu, \mu_0)^2 = \frac{1}{n} \sum_{i=1}^n (\mu(\mathbf{x}_i) - \mu_0(\mathbf{x}_i))^2$ as the distance between μ , μ_0 under a fixed design. When the design is random, let $\rho(\mu, \mu_0)^2 = \int_{\mathcal{M}} (\mu(\mathbf{x}) - \mu_0(\mathbf{x}))^2 F(d\mathbf{x})$, where F is the marginal distribution of the features. Denote $\Pi(\cdot | y_1, \dots, y_n)$ to be the posterior distribution given y_1, \dots, y_n . Then the interest lies in the rate at which the posterior contracts around μ_0 under the metric $\rho(\cdot, \cdot)$. This calls for finding a sequence $\{\zeta_n\}_{n \geq 1}$ of lower bounds such that

$$\Pi(\rho(\mu, \mu_0) > \zeta_n | y_1, \dots, y_n) \rightarrow 0, \text{ as } n \rightarrow \infty. \quad (10)$$

Definition: Given two manifolds \mathcal{M} and \mathcal{N} , a differentiable map $f : \mathcal{M} \rightarrow \mathcal{N}$ is called a diffeomorphism if it is a bijection and its inverse $f^{-1} : \mathcal{N} \rightarrow \mathcal{M}$ is differentiable. If these functions are r times continuously differentiable, f is called a C^r -diffeomorphism.

Our analysis builds on the following result (Theorem 2.3 in Yang and Dunson (2013)).

Theorem 1 *Assume \mathcal{M} is a d dimensional C^{r_1} compact sub-manifold of \mathcal{R}^p . Let $G : \mathcal{M} \rightarrow \mathcal{R}^p$ be the embedding map so that $G(\mathcal{M}) \simeq \mathcal{M}$. Further assume $T : \mathcal{R}^p \rightarrow \mathcal{R}^m$ is a dimensionality reducing map s.t. the restriction $T_{\mathcal{M}}$ of T on $G(\mathcal{M})$ is a C^{r_2} -diffeomorphism*

onto its image. Then for any $\mu_0 \in \mathcal{C}^s$ with $s \leq \min\{2, r_1 - 1, r_2 - 1\}$, a Gaussian process prior on μ with features $\{T(\mathbf{z}_i)\}_{i=1}^n$, $\mathbf{z}_i \in \mathcal{M}$, leads to a posterior contraction rate at least $\zeta_n = n^{-s/(2s+d)} \log(n)^{d+1}$.

This is a huge improvement upon the minimax optimal adaptive rate of $n^{-s/(2s+p)}$ without the manifold structure in the features. We use the above result in our context. Define the linear transformation $T(\mathbf{z}) = \mathbf{\Psi}\mathbf{z}$. Using properties of random projection matrix, we have that, given $\kappa \in (0, 1)$, if the projected dimension $m > O(\frac{m}{\kappa^2} \log(Cp\kappa^{-1}) \log(\phi_n^{-1}))$ then with probability greater than $1 - \phi_n$, the following relationship holds for every point $\mathbf{z}_i, \mathbf{z}_j \in \mathcal{M}$,

$$(1 - \kappa) \sqrt{\frac{m}{p}} \|\mathbf{z}_i - \mathbf{z}_j\| < \|T(\mathbf{z}_i) - T(\mathbf{z}_j)\| < (1 + \kappa) \sqrt{\frac{m}{p}} \|\mathbf{z}_i - \mathbf{z}_j\|, \quad (11)$$

implying that T is a diffeomorphism onto its image with probability greater than $(1 - \phi_n)$. Define $\mathcal{A}_n = \{\text{Equation 11 holds}\}$ so that $P(\mathcal{A}_n) > 1 - \phi_n$.

$$\begin{aligned} \Pi(d(\mu, \mu_0) > \zeta_n | y_1, \dots, y_n) &= \Pi(d(\mu, \mu_0) > \zeta_n | y_1, \dots, y_n, \mathcal{A}_n) P(\mathcal{A}_n) \\ &\quad + \Pi(d(\mu, \mu_0) > \zeta_n | y_1, \dots, y_n, \mathcal{A}'_n) P(\mathcal{A}'_n) \\ &< \Pi(d(\mu, \mu_0) > \zeta_n | y_1, \dots, y_n, \mathcal{A}_n) + P(\mathcal{A}'_n) \\ &< \Pi(d(\mu, \mu_0) > \zeta_n | y_1, \dots, y_n, \mathcal{A}_n) + \phi_n. \end{aligned}$$

On \mathcal{A}_n , T is a diffeomorphism. Therefore, Theorem 1 implies that with features $\{T(\mathbf{z}_i)\}_{i=1}^n$ $\Pi(d(\mu, \mu_0) > \zeta_n | y_1, \dots, y_n, \mathcal{A}_n) \rightarrow 0$. Finally, assuming $\phi_n \rightarrow 0$ yields $\Pi(d(\mu, \mu_0) > \zeta_n | y_1, \dots, y_n) \rightarrow 0$ with features $\{T(\mathbf{z}_i)\}_{i=1}^n$. This proves (A).

Let $\mathbf{\Psi}^{(l)}$ be the l -th row of $\mathbf{\Psi}$, $l = 1, \dots, m$. Denote $\mathbf{\Delta} = [\boldsymbol{\delta}_1 : \dots : \boldsymbol{\delta}_n] \in \mathcal{R}^{p \times n}$ and assume \mathbf{z}_i is the i -th row of $\mathbf{\Delta}$. Using Lemma 2.9.5 in Van der Vaart and Wellner (1996), we obtain

$$\sqrt{p} \sum_{j=1}^p \Psi_{lj} \mathbf{z}_j \rightarrow N(\mathbf{0}, \text{Cov}(\mathbf{z}_1)).$$

Therefore, $\sum_{j=1}^p \Psi_{lj} \mathbf{z}_j = O_p(p^{-1/2})$, reducing the magnitude of noise in the original features. Hence (B) is proved. Thus, even if noise exists, asymptotic performance of $\{T(\mathbf{x}_i)\}_{i=1}^n$ will be similar to $\{T(\mathbf{z}_i)\}_{i=1}^n$ in the GP regression (which by (A) has ‘‘optimal’’ asymptotic performance).

5. Simulation Examples

We assess the performance of compressed Gaussian process (CGP) regression in a number of simulation examples. We consider various numbers of features (p) and level of noise in the features (τ) to study their impact on the performance. In all the simulations out of sample predictive performance of the proposed CGP regression was compared to that of uncompressed Gaussian process (GP), BART (Bayesian Additive Regression Trees) Chipman et al. (2010), RF (Random Forests) Breiman (2001) and TGP (Treed Gaussian process) Gramacy and Lee (2008). Unfortunately, with massive number of features, traditional BART, RF and TGP are computationally prohibitive. Therefore, we consider compressed versions in

which we generate a single projection matrix to obtain a single set of compressed features, running the analysis with compressed features instead of original features. This idea leads to compressed versions of random forest (CRF), Bayesian additive regression tree (CBART) and Treed Gaussian process (CTGP). These methods entail faster implementation when the number of features is massive.

As a default in these analyses, we use $m = 60$, which seems to be a reasonable choice of upper bound for the dimension of the linear subspace to compress to. In addition, we implement two stage GP (2GP) where the p -dimensional features are projected into smaller dimension by using Laplacian eigenmap (Belkin and Niyogi, 2003; Guerrero et al., 2011) in the first stage and then a GP with projected features is fitted in the second stage. We also compared Lasso and partial least square regression (PLSR) to indicate advantages of our proposed method over linear regularizing methods. However, in presence of strong nonlinear relationship between the response and the features, Lasso and PLSR perform poorly and hence results for them are omitted.

When n is moderately large ~ 5000 , to bypass heavy computational price associated with CGP for inverting an $n \times n$ matrix, we employ a low rank approximation of the compressed Gaussian process as described in Section 3. As an uncompressed competitor of CGP in settings with moderately large n , efficient Gaussian random projection technique Banerjee et al. (2013) is implemented. This is also referred to as the GP to avoid needless confusion. Along with GP, CBART and CRF are included as competitors. CTGP with moderately large n poses heavy computational burden and is, therefore, omitted.

As a more scalable competitor, we employ the popular two stage technique of clustering the massive sample into a number of clusters followed by fitting simple model such as Lasso in each of these clusters. To facilitate clustering of high dimensional features in the first stage, we use the spectral clustering algorithm (Ng et al., 2001) described in Algorithm 1. Once observations are clustered, separate Lasso is fitted in each of these clusters. Hence-

Algorithm 1 Spectral Clustering Algorithm

Input: features $\mathbf{x}_1, \dots, \mathbf{x}_n$ and the number of clusters required $n.clust$.

- Form the affinity matrix $\mathbf{A} \in \mathcal{R}^{n \times n}$ defined by $A_{ij} = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma^2)$ if $i \neq j$, $A_{ii} = 0$, for some judicious choice of σ^2 .
 - Define \mathbf{D} to be the diagonal matrix whose (i, i) -th entry is the sum of the elements in the i -th row of \mathbf{A} . Construct $\mathbf{L} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$.
 - Find $\mathbf{s}_1, \dots, \mathbf{s}_{n.clust}$ be the eigenvectors corresponding to the $n.clust$ largest eigenvalues of \mathbf{L} . Form the matrix $\mathbf{S} = [\mathbf{s}_1 : \dots : \mathbf{s}_{n.clust}] \in \mathcal{R}^{n \times n.clust}$ by stacking the eigenvectors in column.
 - Normalize so that each row of \mathbf{S} has unit norm.
 - Now treating each row of \mathbf{S} as a point in $\mathcal{R}^{n.clust}$ cluster them into $n.clust$ clusters via $K - means$ clustering.
 - Finally assign \mathbf{x}_i in cluster j if the i th row of \mathbf{S} goes to cluster j .
-

forth, we refer to this procedure as distributed supervised learning (DSL). Along with the above methods, for large moderately n , we also implement the Bayesian analogue of sparse

Gaussian process with dimension reduction (Snelson and Ghahramani, 2012), referred to as the SG method.

The model averaging step in CGP requires choosing a window over the possible values of m . When n is small, we adopt the choice suggested in Guhaniyogi and Dunson (2013) to have a window of $[\lceil 2\log(p) \rceil, \min(n, p)]$, which implies that the number of possible models to be averaged across is $s = \min(n, p) - \lceil 2\log(p) \rceil + 1$. When n is moderately large, we choose the window of $[\lceil 2\log(p) \rceil, \min(m_{\Phi}, p)]$. The number of rows of Φ is fixed at $m_{\Phi} = 100$ for the simulation study with moderately large n . However, changing m_{Φ} moderately does not alter the performance of CGP.

5.1 Manifold Regression on Swiss Roll

To provide some intuition for our model, we start with a concrete example where the distribution of the response is a nonlinear function of the coordinates along a swissroll, which is embedded in a high dimensional ambient space. To be more specific, we sample manifold coordinates, $t \sim U(\frac{3\pi}{2}, \frac{9\pi}{2})$, $h \sim U(0, 3)$. A high dimensional feature $\mathbf{x} = (x_1, \dots, x_p)$ is then sampled following

$$x_1 = t \cos(t) + \delta_1, x_2 = h + \delta_2, x_3 = t \sin(t) + \delta_3, x_i = \delta_i, i \geq 4, \delta_1, \dots, \delta_p \sim N(0, \tau^2).$$

Finally responses are simulated to have nonlinear and non-monotonic relationship with the features

$$y_i = \sin(5\pi t) + h^2 + \epsilon_i, \epsilon_i \sim N(0, 0.02^2). \quad (12)$$

Clearly, \mathbf{x} and y are conditionally independent given θ, h , which is the low-dimensional signal manifold. In particular, \mathbf{x} lives on a (*noise corrupted*) swissroll embedded in a p -dimensional ambient space (see Figure 1(a)), but y is only a function of coordinates along the swissroll \mathcal{M} (see Figure 1(b)).

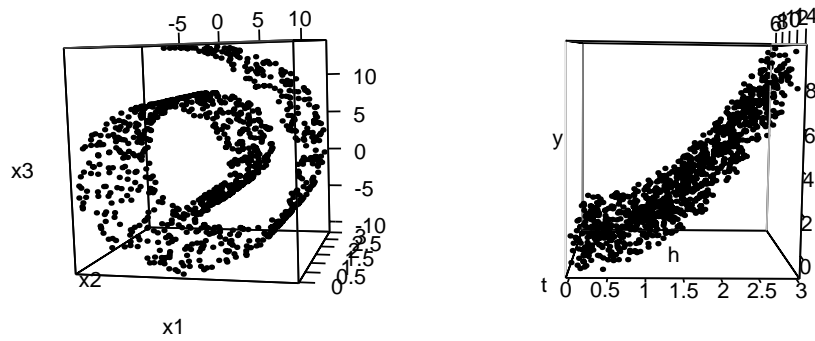
The geodesic distance between two points in a swiss roll can be substantially different from their Euclidean distance in the ambient space \mathcal{R}^p . For example, in Figure 1(c) two points joined by the line segment have much smaller Euclidean distance than geodesic distance. Theorem 1 in Section 4 guarantees optimal performance when the compact submanifold \mathcal{M} is sufficiently smooth, so that the locally Euclidean distance serves as a good approximation of the geodesic distance. The Swiss roll presents a challenging set up for CGP, since points on \mathcal{M} that are close in a Euclidean sense can be quite far in a geodesic sense.

To assess the impact of the number of features (p) and noise levels of the features (τ) on the performance of CGP, a number of simulation scenarios are considered in Table 1. For each of these simulation scenarios, we generate multiple datasets and present predictive inference such as mean squared prediction error (MSPE), coverage and lengths of 95% predictive intervals (PI) averaged over all replicates.

In our experiments, \mathbf{y} and \mathbf{X} are centered. To implement LASSO, we use `glmnet` (Friedman et al., 2009) package in R with the optimal tuning parameter selected through 10 fold cross validation. CRF, CBART and CTGP in R using `randomForest` (Liaw and Wiener, 2002), `BayesTree` (Chipman et al., 2009) and `tgp` (Gramacy, 2007) packages, respectively.

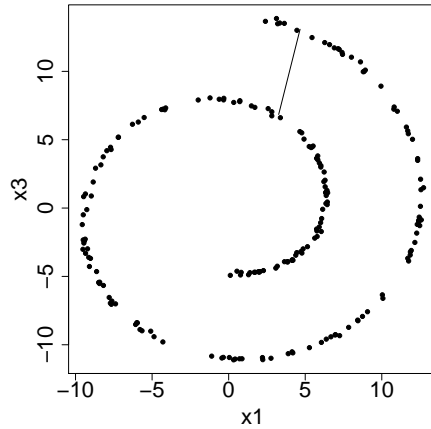
Simulation	sample size (n)	no. of features (p)	noise in the features (τ)
1	100	10,000	0.02
2	100	20,000	0.02
3	100	10,000	0.05
4	100	20,000	0.05
5	100	10,000	0.10
6	100	20,000	0.10

Table 1: Different Simulation settings for CGP.



(a) noise corrupted swiss roll

(b) response vs. x_1, x_2



(c) swiss roll shown in 2d

Figure 1: Simulated features and response on a *noisy* Swiss Roll, $\tau = 0.05$

5.1.1 MSPE RESULTS

Predictive MSE for each of the simulation settings averaged over 50 simulated datasets is shown in Table 2. Subscripted values represent bootstrap standard errors for the averaged

MSPEs, calculated by generating 50 bootstrap datasets resampled from the MSPE values, finding the average MSPE of each, and then computing their standard error.

Table 2 shows that feeding randomly compressed features into any of the nonparametric methods leads to good predictive performance, while Lasso fails to improve much upon the null model (not shown here). For both $p = 10,000$ and $20,000$, when the swiss roll is corrupted with low noise, CGP performs significantly better than GP, while CBART and CRF provide competitive performance with GP. Increasing noise in the features results in deteriorating performances for all the competitors. CGP is an effective tool to reduce the effect of noise in the features, but at a *tipping point* (depending on n) noise distorts the manifold too much, and CGP starts performing similarly to GP. CRF and CTGP perform much worse than CGP in high noise scenarios, while CBART produces competitive performance. Two stage GP (2GP) performs much worse than all the other competitors; perhaps the two stage procedure is considerably more sensitive to noise. Increasing number of features does not alter MSPE for CGP significantly in presence of low noise, consistent with asymptotic results showing posterior convergence rates depend on the intrinsic dimension of \mathcal{M} instead of p when features are concentrated close to \mathcal{M} . In the next section, we will study these aspects with increasing sample size and noise in the features.

	Noise in the feature			
		.02	.05	.10
$p = 10000$	CGP	4.09 _{0.08}	5.49 _{0.08}	7.03 _{0.11}
	GP	4.71 _{0.10}	5.63 _{0.10}	7.09 _{0.12}
	CRF	4.13 _{0.11}	6.23 _{0.09}	7.44 _{0.11}
	CBART	3.73 _{0.13}	6.14 _{0.10}	7.34 _{0.12}
	CTGP	4.24 _{0.14}	7.13 _{0.11}	7.72 _{0.14}
	2GP	5.72 _{0.15}	6.55 _{0.13}	7.85 _{0.16}
$p = 20000$	CGP	4.43 _{0.07}	6.21 _{0.10}	7.28 _{0.13}
	GP	4.86 _{0.07}	6.25 _{0.12}	7.18 _{0.12}
	CRF	5.06 _{0.11}	6.81 _{0.11}	7.47 _{0.13}
	CBART	4.84 _{0.15}	6.77 _{0.11}	7.33 _{0.11}
	CTGP	5.59 _{0.11}	7.40 _{0.11}	7.51 _{0.15}
	2GP	6.05 _{0.10}	6.69 _{0.13}	7.09 _{0.19}

Table 2: Performance comparisons for competitors in terms of mean squared prediction errors (MSPE)

5.1.2 COVERAGE AND LENGTH OF PIS

To assess if CGP is well calibrated in terms of uncertainty quantification, we compute coverage and length of 95% predictive intervals (PI) of CGP along with all the competitors. Although most frequentist methods such as CRF are unable to provide such coverage probabilities in producing point estimates, we present a measure of predictive uncertainty for those methods following the popular two stage plug-in approach, (i) estimate the regression function in the first stage; (ii) construct 95% PI based on the normal distribution centered

on the predictive mean from the regression model with variance equal to the estimated variance in the residuals. Boxplots for coverage probabilities in all the simulation cases are presented in Figure 2. Figure 3 presents median lengths of the 95% predictive intervals.

Both these Figures demonstrate that in all the simulation scenarios CGP, uncompressed GP, 2GP and CBART result in predictive coverage of around 95%, while CRF suffers from severe under-coverage. The gross under-coverage of CRF is attributed to the overly narrow predictive intervals. Additionally, CTGP shows some under-coverage, with shorter predictive intervals than CGP, GP, 2GP or CBART. CGP turns out to be an excellent choice among all the competitors in fairly broad simulation scenarios. We consider larger sample sizes and high noise scenarios in the next subsection.

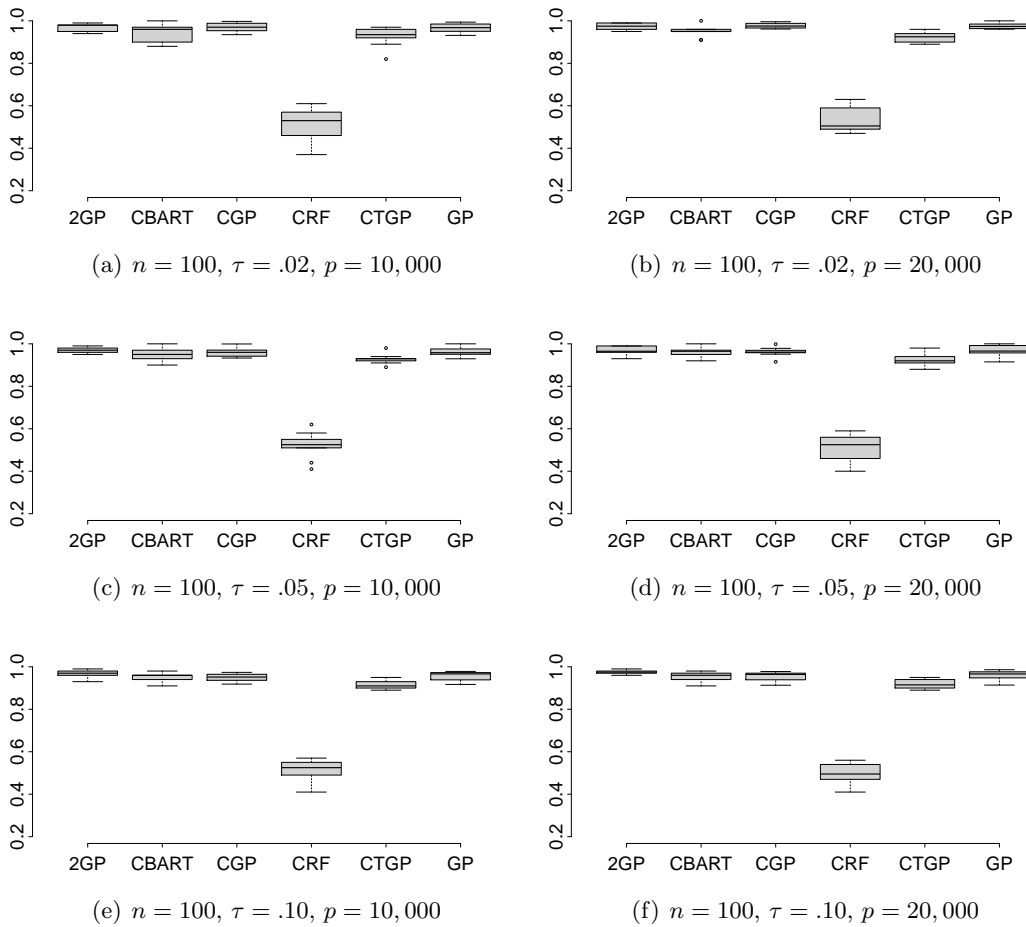


Figure 2: coverage of 95% PI's for CGP, GP, CBART, CTGP, CRF, 2GP

5.2 Manifold Regression on Swiss roll for Larger Samples

To assess how the relative performance of CGP changes for larger sample size, we implement manifold regression on swiss roll using methodologies developed in Section 3. For this simulation example, a data generation scheme similar to Section 5.1 is used. Ideally, larger

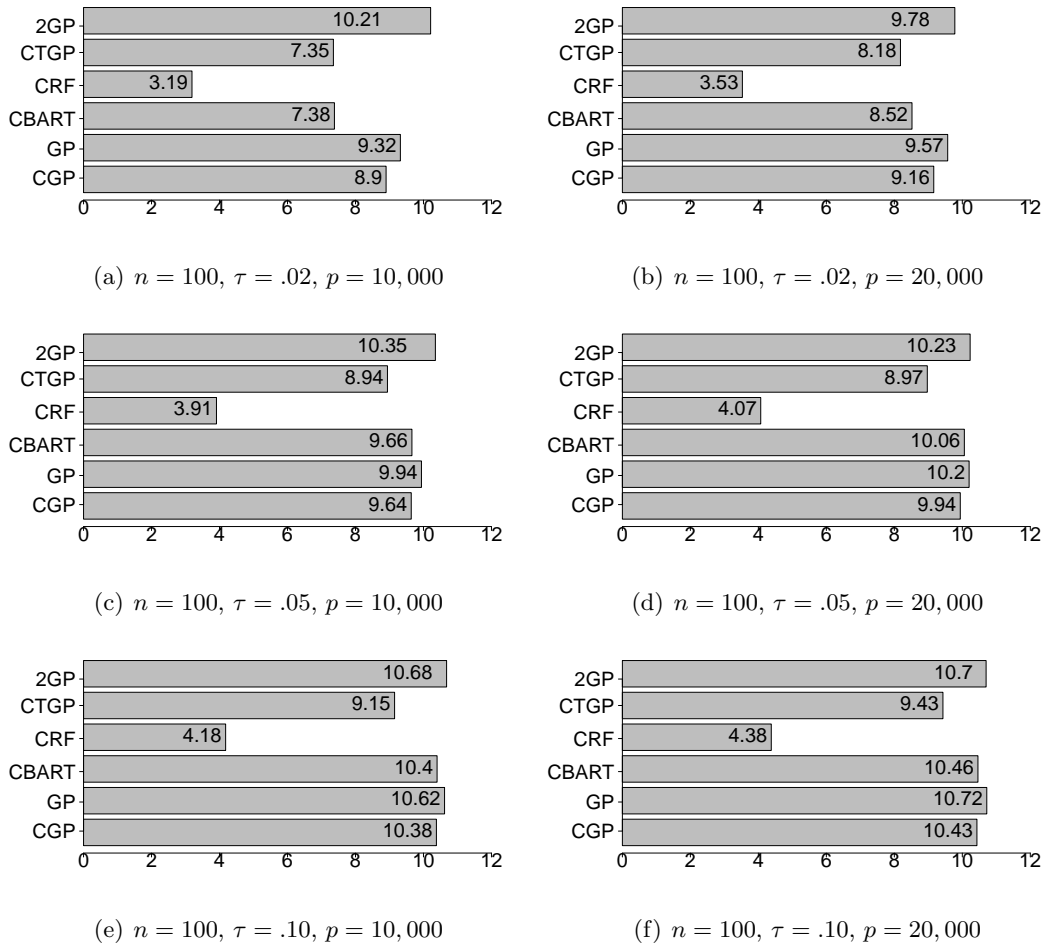


Figure 3: lengths of 95% PI's for CGP, GP, CBART, CRF, CTGP, 2GP

sample size should lead to better predictive performance. Therefore, one would expect more accurate prediction even with higher degree of noise in the features for larger sample size, as long as there is sufficient signal in the data. To accommodate higher signal than in Section 5.1, we simulate manifold coordinates as $t \sim U(\frac{3\pi}{2}, \frac{9\pi}{2})$, $h \sim U(0, 5)$ and sample responses as per (12). We also increase noise variability in the features for all the simulation settings. Simulation scenarios are described in Table 3.

MSPE of all the competing methods are calculated along with their bootstrap standard errors and presented in Table 4. Results in Table 4 provide more evidence supporting our conclusion in Section 5.1. With smaller noise variance, CGP along with other compressed methods outperform uncompressed GP and 2GP. However, when τ exceeds a certain limit, the manifold structure is more and more distorted, with performance of all the competitors worsening. In particular with increasing noise, performance of CGP and GP start becoming more comparable. On the other hand, SG method faces computational issues for $p \sim 10000 - 20000$ features. Therefore, we select only 500 features without disrupting the noisy manifold structure. Even with many fewer features, SG performs worse than CGP with

Simulation	sample size (n)	no. of features (p)	noise in the features (τ)
1	5,000	10,000	.03
2	5,000	20,000	.03
3	5,000	10,000	.06
4	5,000	20,000	.06
5	5,000	10,000	.10
6	5,000	20,000	.10

Table 3: Different Simulation settings for CGP for large n .

MSPE 46.3, 52.2 for $\tau = 0.03, 0.1$ respectively. Our investigation shows that the performance of SG is quite competitive when p is less than a few dozen. However, as p increases over a few hundreds, SG starts performing poorly. This is perhaps due to the fact that SG estimates a large number of poorly identifiable parameters resulting in inaccurate estimation. CGP with random compression of high dimensional features remarkably reduces the number of parameters to be estimated. Comparing results from the last section it is quite evident that with large samples, CGP is able to perform well even with very large number of features and moderate variance of noise in the features. This shows the effectiveness of CGP for large p and moderately large n when features are close to lying on a low-dimensional manifold.

In all the simulation scenarios, DSL is the best performer in terms of MSPE, consistent with the routine use of DSL in large scale settings. However, the performance is extremely sensitive to the choice of clusters. In real data applications often inaccurate clustering leads to suboptimal performance, as will be seen in the data analysis. Additionally, we are not just interested in obtaining a point prediction approach, but want to obtain methods that provide an accurate characterization of predictive uncertainty. With this in mind, we additionally examine coverage probabilities and lengths of 95% predictive intervals (PIs). Boxplots for coverage probabilities of 95% PI's are presented in Figure 4. Figure 5 presents

		Noise in the feature		
		.03	.06	.10
$p = 10,000$	CGP	0.56 _{0.06}	1.06 _{0.03}	2.18 _{0.08}
	GP	2.05 _{0.32}	2.37 _{0.35}	3.35 _{0.42}
	CRF	1.05 _{0.10}	2.16 _{0.11}	3.52 _{0.09}
	CBART	0.69 _{0.07}	1.72 _{0.11}	2.79 _{0.13}
	DSL	0.50 _{0.07}	0.52 _{0.03}	0.50 _{0.03}
	2GP	3.78 _{0.31}	3.95 _{0.41}	4.05 _{0.38}
$p = 20,000$	CGP	1.17 _{0.048}	2.11 _{0.107}	2.57 _{0.222}
	GP	1.98 _{0.418}	2.33 _{0.321}	2.78 _{0.330}
	CRF	1.46 _{0.070}	2.76 _{0.224}	3.88 _{0.224}
	CBART	1.22 _{0.092}	2.53 _{0.151}	3.84 _{0.192}
	DSL	0.48 _{0.015}	0.45 _{0.014}	0.57 _{0.078}
	2GP	3.84 _{0.581}	4.10 _{0.370}	4.53 _{0.481}

Table 4: $MSPE \times 0.1$ along with the bootstrap $sd \times 0.1$ for all the competitors

lengths of 95% prediction intervals for all the competitors. As expected, CGP, GP, 2GP

and CBART demonstrate better performance in terms of coverage. However, in low noise cases CGP and CBART achieve similar coverage with a two fold reduction in the length of PIs compared to GP or 2GP. CRF, like in the previous section, shows under-coverage with narrow predictive intervals. The predictive interval for CGP is found to be marginally wider than CBART with comparable coverage. With high noise, it becomes intractable to recover the manifold structure and hence performance is affected for all the competitors. It is observed that with high noise all approaches tend to have wider predictive intervals. DSL presents overly narrow predictive intervals (not shown here) yielding severe under-coverage.

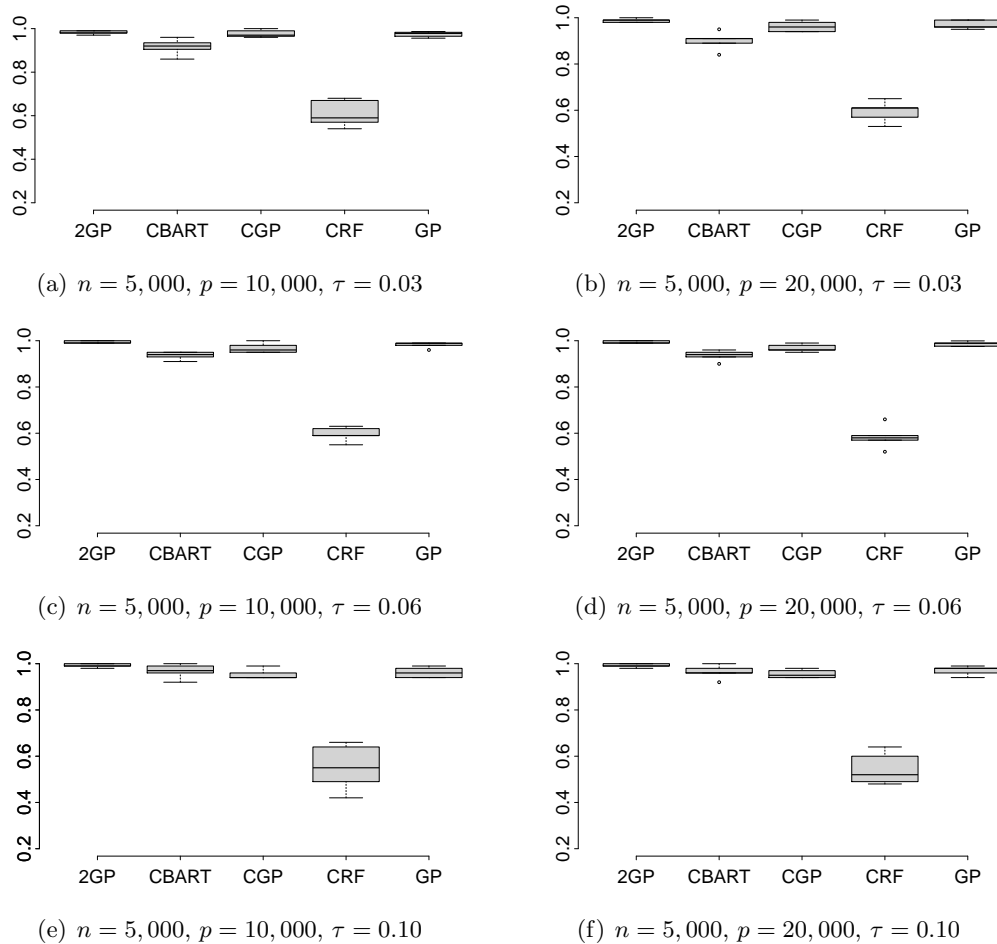


Figure 4: coverage of 95% PI's for CGP, GP, CRF, CBART, 2GP

5.3 Computation Time

One of the major motivations in developing CGP was to improve computational scalability to large p settings. Clearly, the computational time for nonparametric estimation methods such as BART, TGP or RF applied to the original data will become notoriously prohibitive for large p , and hence we focus on comparisons with more scalable methods. The approach of applying BART, RF and TGP to the compressed features, which is employed in CBART,

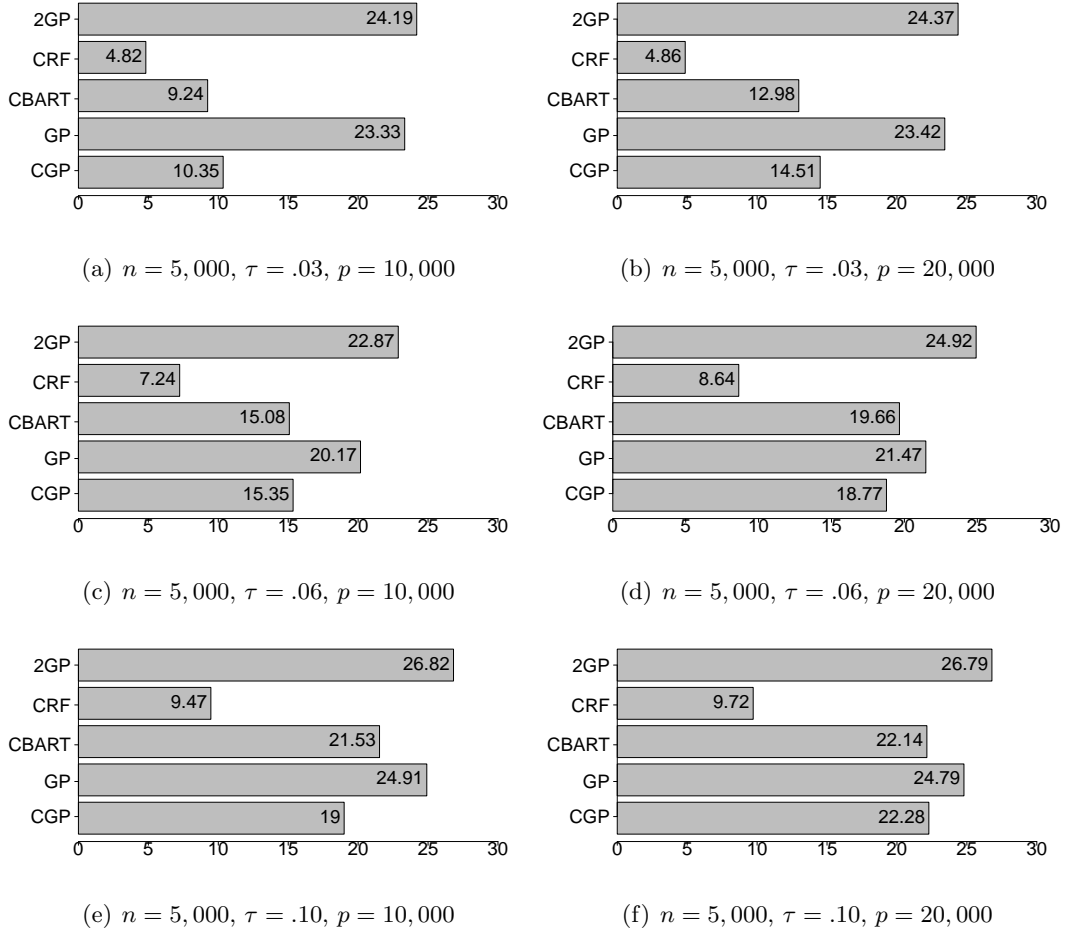


Figure 5: lengths of 95% PI's for CGP, GP, CBART, CRF, 2GP

CRF and CTGP respectively, is faster to implement. Using R code in a standard server, the computing time for 2,000 iterations of CBART for $n = 100$ and $p = 10,000, 20,000$ are only 7.21, 8.36 seconds, while CGP has run time of 7.48, 8.05 seconds, respectively. Increasing n moderately, we find CBART and CGP have similar run time. CRF is a bit faster than both of them, while CTGP has run time 37.64, 38.33 seconds for $p = 10,000, 20,000$ respectively. For moderate n , 2GP is found to have similar run time as CBART.

With large n , CTGP is impractically slow and hence omitted in the comparison. GP needs to calculate and store a distance matrix of p features. Apart from the storage bottleneck, the computational complexity is $O(n^2p)$. CGP instead proposes calculating and storing a distance matrix of m compressed features, with a computational complexity of $O(n^2m)$. Computation time for CGP additionally depends on a number of factors, (i) Gram Schmidt orthogonalization of m rows of $m \times p$ matrices, (ii) inverting an $m_{\Phi} \times m_{\Phi}$ matrix, (iii) multiplying $n \times p$ and $p \times m$ matrices. Along with these three steps, one requires multiplying $m_{\Phi} \times n$ matrix Φ with $n \times n$ matrix K_1 at each MCMC iteration that incurs a computation complexity of order n^2m_{Φ} . Typically the computation complex-

ity is dominated by n^2 and hence scaling with sample size is computationally feasible for about $n \sim 10000$ observations. For much larger n , one can resort to distributed GP based approaches as mentioned in Section 3. On the other hand, SPGP with dimensionality reduction (SG method) introduces exorbitantly large number of parameters even for moderate p .

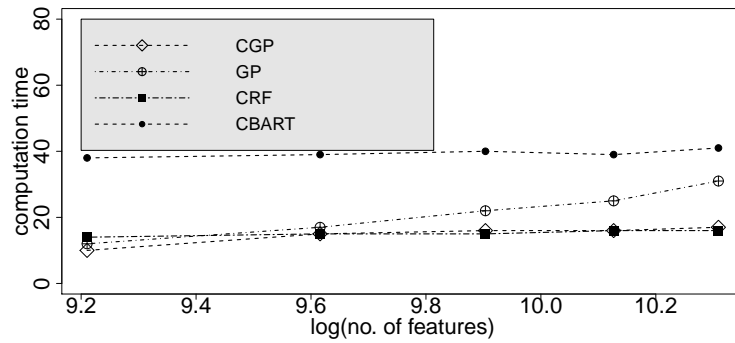
Figure 6 shows the computational speed comparison between CGP, GP, CBART and CRF for various n and p . Computational speed is recorded assuming existence of a number of processors on which parallelization can be executed. As n increases, CGP enjoys substantial computational advantage over competitors. The computational advantage is especially notable over CBART and GP. Run times of DSL are also recorded for $n = 5,000$ and $p = 10,000, 15,000, 20,000, 25,000, 30,000$ and they are 449, 599, 737, 945, 1158 seconds, respectively. Alternatively, 2GP involves creating adjacency matrices followed by an eigen-decomposition of an $n \times n$ matrix. Both these steps are computationally demanding. We find 2GP takes 602, 723, 856, 983, 1108 seconds to run for $n = 5000$ and $p = 10000, 15000, 20000, 25000, 30000$, respectively. Therefore, CGP can outperform even a simple two stage estimation procedure such as DSL in terms of computational speed.

6. Application to Face Images

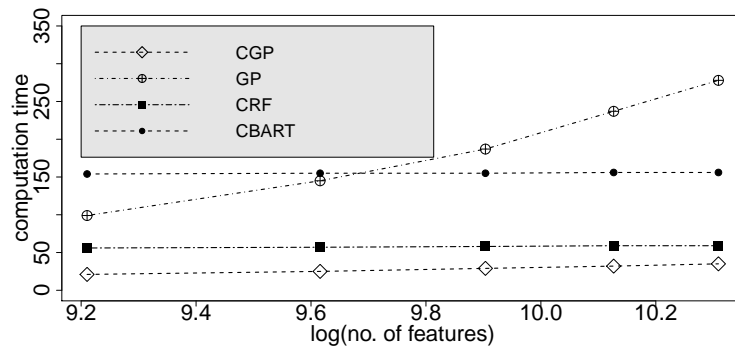
In our simulation examples, the underlying manifold is three dimensional and can be directly visualized. In this section we present an application in which both the dimension and the structure of the underlying manifold is unknown. The dataset consists of 698 images of an artificial face and is referred to as the *Isomap face data* (Tenenbaum et al., 2000). A few such representative images are presented in Figure 7. Each image is labeled with three different variables: illumination-, horizontal- and vertical-orientation. Two dimensional projections of the images are presented in the form of 64×64 pixel matrices. Intuitively, a limited number of additional features are needed for different views of the face. This is confirmed by the recent work of Levina and Bickel (2004); Aswani et al. (2011) where the intrinsic dimensionality is estimated to be small from these images. More details about the dataset can be found in <http://isomap.stanford.edu/datasets.html>.

We apply CGP and all the competitors to the dataset to assess relative performances. To set up the regression problem, we consider horizontal pose angles (vary in $[-75^0, 75^0]$) of the images, after standardization, as the responses. The features are taken $64 \times 64 = 4096$ dimensional vectorized images for each sample. To simulate more realistic situations, $N(0, \tau^2)$ noise is added to each pixel of the images, with varying τ , to make predictive inference more challenging from the noisy images. We carry out random splitting of the data into $n = 648$ training cases and $n_{pred} = 50$ test cases and run all the competitors to obtain predictive inference in terms of MSPE, length and coverage of 95% predictive intervals. To avoid spurious inference due to small validation set, this experiment is repeated 50 times.

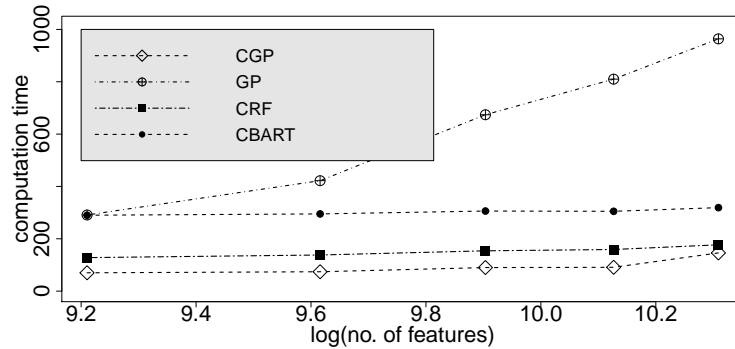
Table 5 presents MSPE for all the competing methods averaged over 50 experiments along with their standard errors computed using 100 bootstrap samples. Note that, because of the standardization, the null model yields MSPE 1. It is clear from Table 5 that CGP along with its compressed competitors explain a lot of variation in the response. DSL and 2GP are the worst performers in terms of MSPE. This is consistent with our experience that, in the presence of a complex and unknown manifold structure along with noise, DSL



(a) $n = 1,000$



(b) $n = 3,000$



(c) $n = 5,000$

Figure 6: Computational time in seconds for CGP, GP, CBART, CRF against log of the number of features.

can be unreliable relative to CGP which tends to be more robust to the type of manifold and noise level. GP also performs much worse than CGP and other compressed competitors especially in presence of small amount of noise in the features. As the noise in the features increases, performance of CGP and GP are found to be comparable. On the other hand



Figure 7: Representative images from the Isomap face data.

τ	CGP	GP	CBART	CRF	DSL	2GP
0.03	0.07 _{0.004}	0.85 _{0.054}	0.07 _{0.005}	0.06 _{0.009}	0.70 _{0.010}	0.98 _{0.001}
0.06	0.08 _{0.008}	0.75 _{0.043}	0.09 _{0.008}	0.10 _{0.012}	0.78 _{0.015}	0.94 _{0.022}
0.10	0.09 _{0.003}	0.68 _{0.041}	0.11 _{0.006}	0.11 _{0.004}	0.83 _{0.024}	0.98 _{0.001}

Table 5: MSPE and standard error (computed using 100 bootstrap samples) for all the competitors over 50 replications

SG implemented with only a subset of 500 features yields much worse performance (MSPE 0.97, 0.98 for $\tau = 0.1, 0.03$) respectively.

To see how well calibrated these methods are, Figure 8 provides coverage probabilities along with the lengths of predictive intervals for all the competitors. It is evident from the Figure that CGP, CBART, GP and 2GP yield excellent coverage. However, for CGP and CBART this coverage is achieved with much narrower predictive intervals compared to GP and 2GP. On the other hand, both CRF and DSL produce extremely narrow predictive intervals resulting in severe under-coverage. In fact for $\tau = 0.03, 0.06, 0.10$, length of 95% predictive intervals for DSL are 0.13, 0.19, 0.21 respectively. Therefore, both in terms of MSPE and predictive coverage, CGP does a good job. More importantly, these results serve as a testimony of the robust performance demonstrated by compressed Bayesian non-parametric methods (CGP being one of them) even in the presence of unknown and complex manifold structure in the features.

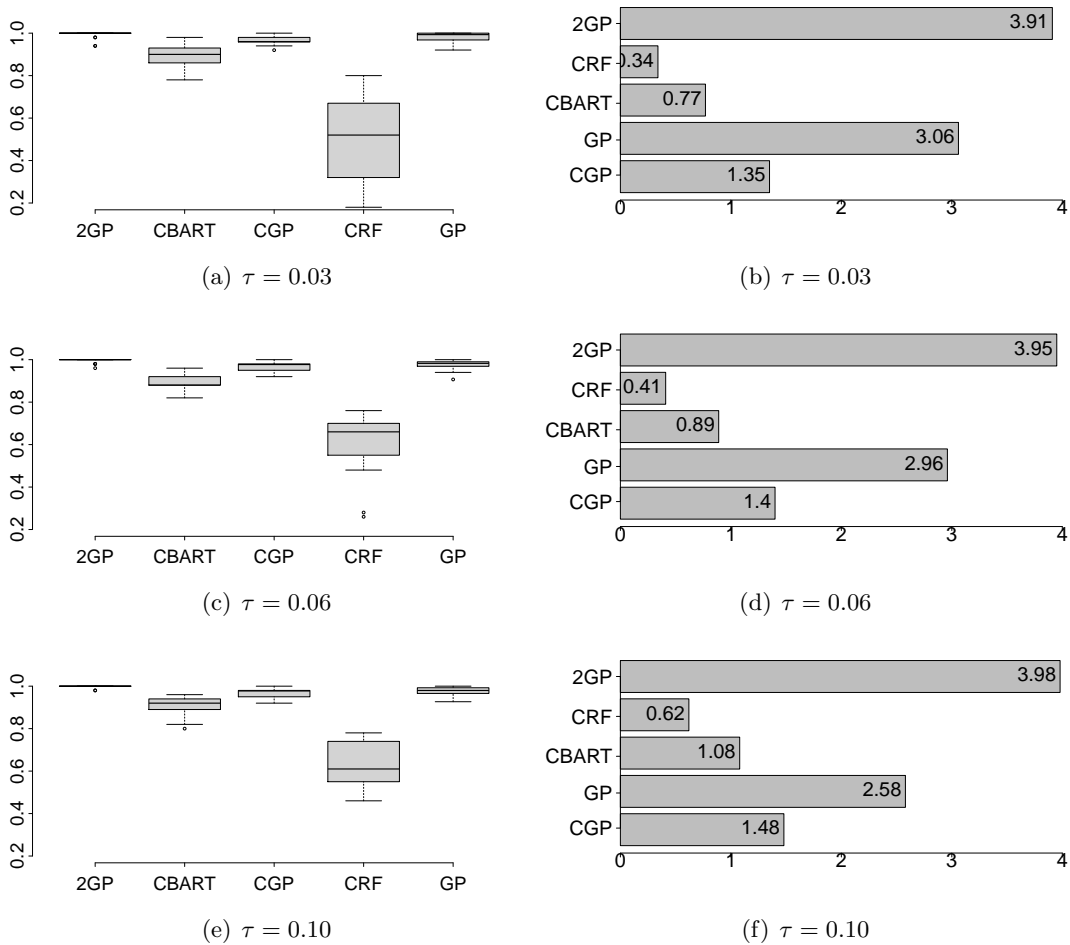


Figure 8: Left panel: Boxplot for coverage of 95% predictive intervals over 50 replications; Right panel: Boxplot for length of 95% predictive intervals over 50 replications for CGP, GP, 2GP, CBART, CRF. In the left and right panels y-axis corresponds to the coverage and length respectively.

7. Discussion

The overarching goal of this article is to develop nonparametric regression methods that scale to large/very large p and/or moderately large $n \sim 5000$ when features lie on a *noise corrupted* manifold. The statistical and machine learning literature is somewhat limited in robust and flexible methods that can accurately provide predictive inference for large p with moderately large sample size, while taking into account the geometric structure. We develop a method based on nonparametric *low-rank* Gaussian process methods combined with random feature compression to accurately characterize predictive uncertainties quickly, bypassing the need to estimate the underlying manifold. The computational template exploits model averaging to limit sensitivity of the inference to the specific choices of the random

projection matrix Ψ . The proposed method is also guaranteed to yield minimax optimal convergence rates.

There are many future directions motivated by our work. For example, the present work uses Banerjee et al. (2013) that is less suitable for massive n . It is quite straightforward to extend CGP to massive n by directly applying recently developed approaches for distributed computation in GP models (Deisenroth and Ng, 2015). Also the present work is not able to estimate the true dimensionality of the noise corrupted manifold. Arguably, a nonparametric method that can simultaneously estimate the intrinsic dimensionality of the manifold in the ambient space would improve performance both theoretically and practically. One possibility is to simultaneously learn the marginal distribution of the features, accounting for the low-dimensional structure. Other possible directions include adapting to massive streaming data where inference is to be made online. Although random compression both in n and p provides substantial benefit in terms of computation and inference, it might be worthwhile to learn the matrices Ψ , Φ while attempting to limit the associated computational burden.

References

- Anil Aswani, Peter Bickel, and Claire Tomlin. Regression on manifolds: Estimation of the exterior derivative. *The Annals of Statistics*, 39(1):48–81, 2011.
- Anjishnu Banerjee, David B Dunson, and Surya T Tokdar. Efficient Gaussian process regression for large datasets. *Biometrika*, 100(1):75–89, 2013.
- Sudipto Banerjee, Alan E Gelfand, Andrew O Finley, and Huiyan Sang. Gaussian predictive process models for large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(4):825–848, 2008.
- Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396, 2003.
- Peter J Bickel and Bo Li. Local polynomial regression on unknown manifolds. *Lecture Notes-Monograph Series*, 54:177–186, 2007.
- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.
- Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- Leo Breiman, Jerome Friedman, Charles J Stone, and Richard A Olshen. *Classification and regression trees*. CRC press, 1984.
- Roberto Calandra, Jan Peters, Carl Edward Rasmussen, and Marc Peter Deisenroth. Manifold Gaussian processes for regression. *arXiv preprint arXiv:1402.5876*, 2014.
- Minhua Chen, Jorge Silva, John Paisley, Chunping Wang, David Dunson, and Lawrence Carin. Compressive sensing on manifolds using a nonparametric mixture of factor ana-

- lyzers: Algorithm and performance bounds. *Signal Processing, IEEE Transactions on*, 58(12):6140–6155, 2010.
- Hugh Chipman, Robert McCulloch, and Maintainer Robert McCulloch. Package bayestree. 2009.
- Hugh A Chipman, Edward I George, and Robert E McCulloch. Bart: Bayesian additive regression trees. *The Annals of Applied Statistics*, 4(1):266–298, 2010.
- Abhirup Datta, Sudipto Banerjee, Andrew O Finley, and Alan E Gelfand. Hierarchical nearest-neighbor Gaussian process models for large geostatistical datasets. *arXiv preprint arXiv:1406.7343*, 2014.
- Marc Peter Deisenroth and Jun Wei Ng. Distributed Gaussian processes. *arXiv preprint arXiv:1502.02843*, 2015.
- Xavier Emery. The kriging update equations and their application to the selection of neighboring data. *Computational Geosciences*, 13(3):269–280, 2009.
- Mahdi Milani Fard, Yuri Grinberg, Joelle Pineau, and Doina Precup. Compressed least-squares regression on sparse spaces. In *AAAI*, 2012.
- Andrew O Finley, Sudipto Banerjee, Patrik Waldmann, and Tore Ericsson. Hierarchical spatial modeling of additive and dominance genetic variance for large spatial trial datasets. *Biometrics*, 65(2):441–451, 2009.
- Jerome Friedman, Trevor Hastie, and Rob Tibshirani. glmnet: Lasso and elastic-net regularized generalized linear models. *R package version*, 1, 2009.
- Robert B Gramacy. tgp: an r package for Bayesian nonstationary, semiparametric nonlinear regression and design by treed Gaussian process models. *Journal of Statistical Software*, 19(9):6, 2007.
- Robert B Gramacy and Daniel W Apley. Local Gaussian process approximation for large computer experiments. *Journal of Computational and Graphical Statistics*, 24(2):561–578, 2015.
- Robert B Gramacy and Herbert KH Lee. Bayesian treed Gaussian process models with an application to computer modeling. *Journal of the American Statistical Association*, 103(483):1119–1130, 2008.
- Ricardo Guerrero, Robin Wolz, and Daniel Rueckert. Laplacian eigenmaps manifold learning for landmark localization in brain mr images. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2011*, pages 566–573. Springer, 2011.
- Rajarshi Guhaniyogi and David B Dunson. Bayesian compressed regression. *arXiv preprint arXiv:1303.0642*, 2013.
- Dave Higdon. Space and space-time modeling using process convolutions. In *Quantitative methods for current environmental issues*, pages 37–56. Springer, 2002.

- Cari G Kaufman, Mark J Schervish, and Douglas W Nychka. Covariance tapering for likelihood-based estimation in large spatial data sets. *Journal of the American Statistical Association*, 103(484):1545–1555, 2008.
- Neil Lawrence. Probabilistic non-linear principal component analysis with Gaussian process latent variable models. *The Journal of Machine Learning Research*, 6:1783–1816, 2005.
- Elizaveta Levina and Peter J Bickel. Maximum likelihood estimation of intrinsic dimension. *Ann Arbor MI*, 48109:1092, 2004.
- Andy Liaw and Matthew Wiener. Classification and regression by randomforest. *R news*, 2(3):18–22, 2002.
- Odalric-Ambrym Maillard and Rémi Munos. Compressed least-squares regression. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1213–1221, 2009.
- Andrew Y Ng, Michael I Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. *Proceedings of Advances in Neural Information Processing Systems. Cambridge, MA: MIT Press*, 14:849–856, 2001.
- Garritt Page, Abhishek Bhattacharya, and David Dunson. Classification via Bayesian non-parametric learning of affine subspaces. *Journal of the American Statistical Association*, 108(501):187–201, 2013.
- Neal Parikh and Stephen Boyd. Block splitting for large-scale distributed learning. In *Neural Information Processing Systems (NIPS), Workshop on Big Learning*. Citeseer, 2011.
- Brian J Reich, Howard D Bondell, and Lexin Li. Sufficient dimension reduction via Bayesian mixture modeling. *Biometrics*, 67(3):886–895, 2011.
- Alex J Smola and Bernhard Schölkopf. Sparse greedy matrix approximation for machine learning. pages 911–918, 2000.
- Edward Snelson and Zoubin Ghahramani. Sparse Gaussian processes using pseudo-inputs. *Advances in neural information processing systems*, 18:1257, 2006.
- Edward Snelson and Zoubin Ghahramani. Variable noise and dimensionality reduction for sparse gaussian processes. *arXiv preprint arXiv:1206.6873*, 2012.
- Michael L Stein, Zhiyi Chi, and Leah J Welty. Approximating likelihoods for large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(2):275–296, 2004.
- Jonathan R Stroud, Michael L Stein, and Shaun Lysen. Bayesian and maximum likelihood estimation for Gaussian processes on an incomplete lattice. *arXiv preprint arXiv:1402.4281*, 2014.
- Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.

- Surya T Tokdar, Yu M Zhu, and Jayanta K Ghosh. Bayesian density regression with logistic gaussian process and subspace projection. *Bayesian analysis*, 5(2):319–344, 2010.
- Aad W van der Vaart and J Harry van Zanten. Adaptive Bayesian estimation using a gaussian random field with inverse gamma bandwidth. *The Annals of Statistics*, 37(5B):2655–2675, 2009.
- AW Van der Vaart and JA Wellner. *Weak Convergence and Empirical Processes*. Springer, New York, 1996.
- Aldo V Vecchia. Estimation and model identification for continuous spatial processes. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 297–312, 1988.
- Christopher K Wikle and Noel Cressie. A dimension-reduced approach to space-time kalman filtering. *Biometrika*, 86(4):815–829, 1999.
- Yun Yang and David B Dunson. Bayesian manifold regression. *arXiv preprint arXiv:1305.0617*, 2013.
- Yuchen Zhang, John C Duchi, and Martin J Wainwright. Divide and conquer kernel ridge regression: A distributed algorithm with minimax optimal rates. *arXiv preprint arXiv:1305.5029*, 2013.