

Efficacy of data fusion using convolved multi-output Gaussian processes

Shrihari Vasudevan, Arman Melkumyan and Steven Scheduling

Australian Centre for Field Robotics, The University of Sydney, NSW 2006, Australia
Email: shrihari.vasudevan@ieee.org | a.melkumyan@acfr.usyd.edu.au | s.scheduling@acfr.usyd.edu.au

Abstract

This paper evaluates the efficacy of a machine learning approach to data fusion using convolved multi-output Gaussian processes in the context of geological resource modeling. It empirically demonstrates that information integration across multiple information sources leads to superior estimates of all the quantities being modeled, compared to modeling them individually. Convolved multi-output Gaussian processes provide a powerful approach for simultaneous modeling of multiple quantities of interest while taking correlations between these quantities into consideration. Experiments are performed on large scale data taken from a mining context.

Keywords - Gaussian process, Machine learning, Data fusion, Geological resource modeling, Mining

1 Introduction

Gaussian processes (GPs) [Rasmussen and Williams, 2006] are powerful non-parametric Bayesian learning techniques that can handle correlated, uncertain and incomplete data. GPs yield a continuous domain representation of the data and hence can be sampled at any desired resolution (multi-scale model). They model and use the spatial correlation of the given data to estimate the values for unknown points of interest. GPs perform *Kriging* [Matheron, 1963, Cressie, 1993] interpolation. As a case in point, the work [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2009] modeled large scale terrain using GPs. It proposed the use of the non-stationary neural network kernel to model large scale discontinuous spatial data and empirically demonstrated that this kernel was superior (in modeling) to the stationary squared exponential kernel and at least as good as most standard interpolation techniques for a range of terrain (in terms of sparsity/complexity/discontinuities). Data fusion in the context of Gaussian processes is necessitated by the presence of multiple, multi-sensor, multi-attribute, incomplete and/or uncertain data sets of the entity being modeled. This paper presents an empirical evaluation of a machine learning approach to performing data fusion with Gaussian processes (based on convolved GPs) in the context of vector valued data. The objectives are to understand (1) if simultaneous modeling of multiple quantities of interest using GPs (i.e. modeling and using the correlations between them and

hence performing data fusion) is better than modeling these quantities independently and (2) if non-stationary kernels are more effective than stationary kernels for modeling geological data. Experiments are performed on large scale data obtained from a mining context.

2 Related work

This section first reviews recent works in the machine learning and related communities that addressed the problem of data fusion with Gaussian processes (GPs) and then positions the proposed approach within the wider geostatistics literature on which many of the techniques being presented are based.

Preliminary attempts at data fusion with GPs were shown in [El-Beltagy and Wright, 2001] and [Murray-Smith and Pearlmutter, 2005]. The former demonstrated how a GP could be used to model an expensive process by first modeling a GP on an approximate or cheap process and subsequently using the many input-output data from the approximate process together with the few samples available of the expensive process in order to learn a GP for the latter. The work [Murray-Smith and Pearlmutter, 2005] attempted to generalize arbitrary transformations on GP priors through linear transformations. It hinted at how this framework could be used to introduce heteroscedasticity (random variables with non-constant variance) and how information from different sources could be combined. However, specifics on how the fusion could actually be performed were beyond the scope of the work. Girolami in [Girolami, 2006] integrated heterogeneous feature types within a Gaussian process classification setting, in a protein fold recognition application domain. Each feature representation was represented by a separate GP. Fusion used the idea that individual feature representations were considered independent and hence a composite covariance function could be defined in terms of a linear sum of Gaussian process priors. Reece et al. in [Reece, Roberts, Nicholson, and Lloyd, 2011] integrated “hard” data obtained from sensors with “soft” information obtained from human sources within a Gaussian process classification framework. It used heterogeneous information-types as mutually independent sources of information that were transformed into the kernel representation (a kernel for each kind of information) and combined using a product rule. While [Girolami, 2006] and [Reece, Roberts, Nicholson, and Lloyd, 2011] demonstrated how multiple information-types may be combined using kernel methods, the approach presented in this paper addresses the problem of improving the predictions of several different (heterogeneous) quantities being simultaneously modeled by explicitly modeling correlations between them.

Recent approaches demonstrating data fusion with Gaussian processes in the context of large scale

terrain modeling were based on heteroscedastic GPs [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2010a] and dependent GPs [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2010b, 2011]. These addressed the problem of fusing multiple, multi-sensor data sets of a single quantity of interest. This paper describes the framework for extending this concept to multiple heterogeneous quantities of interest. The work [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2010a] treated the data-fusion problem as one of combining different noisy samples of a common entity (terrain) being modeled. The GP representing the fused output was assumed to have non-constant noise variance; the noise variance at any point was dependent on the input data from the different data sets. In the Machine Learning community, this idea is referred to as heteroscedastic GPs [Goldberg, Williams, and Bishop, 1998, Kersting, Plagemann, Pfaff, and Burgard, 2007]. Vasudevan et al in [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2010b, 2011] treated the data fusion problem as one of improving GP regression through modeling the spatial correlations (auto and cross covariances) between several dependent GPs representing the respective data sets. This idea was inspired by recent machine learning contributions in multi-task or multi-output or dependent GP modeling including [Bonilla, Chai, and Williams, 2007] and [Boyle and Frean, 2004], the latter being based on [Higdon, 2002]. In Kriging terminology, this idea is akin to Co-kriging [Wackernagel, 2003, Goovaerts, 1997]. The work [Vasudevan, 2012] performed a model complexity analysis of multiple approaches to data fusion using GPs, applied in the context of large scale terrain modeling. Melkumyan et al, in [Melkumyan and Ramos, 2011], compared various combinations of stationary kernel including the squared exponential (SQEXP), Matern 3/2 and a sparse covariance function [Melkumyan and Ramos, 2009] in the context of geological resource modeling. This paper evaluates a machine learning approach to data fusion using a convolved GP model applied to a multi-output (vector valued output) problem of geological resource modeling; multiple kernel combinations including the best results from past works mentioned above are compared.

The paper [Álvarez, Rosasco, and Lawrence, 2012] presents an extensive review of kernels for vector valued functions. Drawing parallels between regularization and Bayesian perspectives, the linear model of corregionalization (LMC) [Wackernagel, 2003, Goovaerts, 1997] is presented as a generalized model on which the geostatistical approaches to multivariate modeling and subsequently machine learning / Gaussian process modeling of multivariate data are based. The LMC incorporates a nested framework wherein outputs are first expressed as a linear combination of spatially uncorrelated processes, each of these processes is a linear sum of uncorrelated functions. Within this framework, kernels are broadly

classified as being separable (a Sum of Separable kernels) or non-separable (e.g. divergence-free and curl-free kernels). Separable kernels expresses the covariance between outputs as a product of two terms - one that captures the covariance between outputs, not considering the inputs and the other that captures the covariance between inputs, not considering the outputs. Such separation is not obtained in a nonseparable kernel formulation. Among approaches to developing nonseparable kernels, the process convolution (PC) approach provides a different and flexible way of modeling vector valued (multi-output or multi-task) functions. In the PC approach, each component of a vector output is modeled as a base process (e.g. a GP) convolved with a smoothing kernel and treated with a noise process (e.g. a GP). This transformation of base processes (or latent functions) contrasts with the LMC formulation. The PC approach captures the trend and complexity of the data through the parameters of the smoothing kernel (and additional parameters as required); the LMC relies on a repertoire of latent functions, a linear sum of which will adequately describe the data. Under assumption of Dirac delta smoothing kernels for each output, the covariance between outputs as obtained by the PC approach reduces to that of the LMC approach [Álvarez, Rosasco, and Lawrence, 2012]. This paper adopts the process convolution approach to modeling vector valued data. The proposed approach augments the PC based covariance between input data with a symmetric matrix of free parameters that is intended to capture the signal variance and similarities between outputs being modeled.

This paper reports a detailed multi-metric performance comparison experiment in a geological resource characterization context, performed between a convolved multi-output GP, an equivalent set of GPs (derived from the multi-output GP parameters) and a set of independently optimized GPs, to provide for an exact and an independent comparison between them. The objective is to quantify the benefit (if any) of simultaneous modeling of the multiple quantities by modeling and using the correlations between them as against modeling each of these quantities separately. This paper also compares data fusion using multiple stationary kernels and a nonstationary kernel in the context of modeling geological data. Further discussion on the broader issues relating to the application of the approach in practical problems and the tying together of different prior works that have studied this approach [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2010b, 2011, Melkumyan and Ramos, 2011] is presented in the technical report version of this paper [Vasudevan, Melkumyan, and Scheduling, 2012].

3 Approach

3.1 Gaussian processes for scalar and vector valued functions

Gaussian processes [Rasmussen and Williams, 2006] (GPs) are stochastic processes wherein any finite subset of random variables is jointly Gaussian distributed. They may be thought of as a Gaussian probability distribution in function space. They are characterized by a mean function $m(\mathbf{x})$ and the covariance function $k(\mathbf{x}, \mathbf{x}')$ that together specify a distribution over functions. In the context of geological resource modeling, each $\mathbf{x} \equiv (\textit{east}, \textit{north}, \textit{depth})$ (3D coordinates) and $f(\mathbf{x}) \equiv z$, the concentration of the quantity being modeled. Although not necessary, the mean function $m(\mathbf{x})$ may be assumed to be zero by scaling/shifting the data appropriately such that it has an empirical mean of zero.

The covariance function or kernel models the relationship between the random variables corresponding to the given data. It can take numerous forms [Rasmussen and Williams, 2006, chap. 4]. The stationary squared exponential (or Gaussian) kernel (SQEXP) is given by

$$k_{SQEXP}(\mathbf{x}, \mathbf{x}', \Sigma) = \sigma_f^2 \cdot \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T \Sigma (\mathbf{x} - \mathbf{x}')\right), \quad (1)$$

where k is the covariance function or kernel; $\Sigma = \textit{diag}[l_{\textit{east}}, l_{\textit{north}}, l_{\textit{depth}}]^{-2}$ is a $d \times d$ diagonal length-scale matrix ($d = \textit{dimensionality of input} = 3$ in this case), a measure of how quickly the modeled function changes in the east, north and depth directions; σ_f^2 is the signal variance. The set of parameters $\{l_{\textit{east}}, l_{\textit{north}}, l_{\textit{depth}}, \sigma_f\}$ are referred to as the kernel hyperparameters.

The non-stationary neural network (NN) kernel [Neal, 1996, Williams, 1998a,b] takes the form

$$k_{NN}(\mathbf{x}, \mathbf{x}', \Sigma) = \sigma_f^2 \cdot \frac{2}{\pi} \arcsin\left(\frac{2\tilde{\mathbf{x}}^T \Sigma \tilde{\mathbf{x}'}}{\sqrt{(1 + 2\tilde{\mathbf{x}}^T \Sigma \tilde{\mathbf{x}})(1 + 2\tilde{\mathbf{x}'^T \Sigma \tilde{\mathbf{x}'})}}\right), \quad (2)$$

where $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{x}'}$ are augmented input vectors (each point is augmented with a 1), Σ is a $(d + 1) \times (d + 1)$ diagonal length-scale matrix given by $\Sigma = \textit{diag}[\beta, l_{\textit{east}}, l_{\textit{north}}, l_{\textit{depth}}]^{-2}$, β being a bias factor and d being the dimensionality of the input data. The variables $\{\beta, l_{\textit{east}}, l_{\textit{north}}, l_{\textit{depth}}, \sigma_f\}$ constitute the kernel hyperparameters. The NN kernel represents the covariance function of a neural network with a single hidden layer between the input and output, infinitely many hidden nodes and using a Sigmoidal transfer function [Williams, 1998a] for the hidden nodes. Hornik, in [Hornik, 1993], showed that such neural networks are universal approximators and Neal, in [Neal, 1996], observed that the functions produced by such a network would tend to a Gaussian process. Prior work in [Vasudevan,

Ramos, Nettleton, and Durrant-Whyte, 2009] found the NN kernel to be more effective than the SQEXP kernel at modeling discontinuous data.

The Matern 3/2 kernel is another stationary kernel differing from the SQEXP kernel in that the latter is infinitely differentiable and consequently tends to have a strong smoothing nature, which is argued as being detrimental to modeling physical processes [Rasmussen and Williams, 2006]. It takes the form specified in Equation 3.

$$k_{MATERN3}(x, x', \Sigma) = \sigma_f^2 \cdot \prod_{1 \leq k \leq d} \left(1 + \frac{\sqrt{3}r_k}{l_k}\right) \exp\left(-\frac{\sqrt{3}r_k}{l_k}\right) \quad (3)$$

where $k \in 1 \dots d$ is the dimension of the input data ($d = \text{dimensionality of input} = 3$ in this case), $\Sigma = [l_{east}, l_{north}, l_{depth}]$ is a $1 \times d$ length-scale matrix and σ_f^2 is the signal variance. The set of parameters $\{l_{east}, l_{north}, l_{depth}, \sigma_f\}$ is referred to as the kernel hyperparameters.

Regression using GPs uses the fact that any finite set of training (evaluation) data and test data of a GP are jointly Gaussian distributed. Assuming noise free data, this idea is shown in Expression 4 (hereafter referred to as Equation 4). This leads to the standard GP regression equations yielding an estimate (the mean value, given by Equation 5) and its uncertainty (Equation 6).

$$\begin{bmatrix} \mathbf{z} \\ f_* \end{bmatrix} \sim N\left(0, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right) \quad (4)$$

$$\bar{f}_* = K(X_*, X) K(X, X)^{-1} \mathbf{z} \quad (5)$$

$$\text{cov}(f_*) = K(X_*, X_*) - K(X_*, X) K(X, X)^{-1} K(X, X_*) \quad (6)$$

For n training points $(X, \mathbf{z}) = (\mathbf{x}_i, z_i)_{i=1 \dots n}$ and n_* test points (X_*, f_*) , $K(X, X_*)$ denotes the $n \times n_*$ matrix of covariances evaluated at all pairs of training and test points. The terms $K(X, X)$, $K(X_*, X_*)$ and $K(X_*, X)$ are defined likewise. In the event that the data being modeled is noisy, a noise hyperparameter (σ) is also learnt with the other GP hyperparameters and the covariance matrix of the training data $K(X, X)$ is replaced by $[K(X, X) + \sigma^2 I]$ in Equations 4, 5 and 6. GP hyperparameters may be learnt using various techniques such as cross validation [Rasmussen and Williams, 2006], maximum-a-posteriori estimation using Markov Chain Monte Carlo techniques [Rasmussen and Williams, 2006, Williams, 1998b] and maximizing the marginal likelihood of the observed training data [Rasmussen and

Williams, 2006, Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2009]. This paper adopts the latter most approach based on the intuition that it may be more suited for large data sets. The marginal likelihood to be maximized is described in Equation 7.

$$\log p(\mathbf{z}|X, \theta) = -\frac{1}{2}\mathbf{z}^T K(X, X)^{-1}\mathbf{z} - \frac{1}{2} \log |K(X, X)| - \frac{n}{2} \log(2\pi) \quad (7)$$

In the geological resource modeling problem considered in this paper, the objective is to model concentrations of multiple elements across the field of interest. The data fusion aspect of this problem is the improved prediction of each one of these concentrations by integration or use of all other concentrations available. If the concentration of each element is modeled using a separate GP, the objective is to improve one GPs prediction estimates given all other GP models. Simultaneous modeling (and consequently data fusion) of multiple concentrations can be achieved by the use of multi-output (MOGP) or dependent GPs [Vasudevan, 2012], also known as multi-task GP if the input points for each element are different. One way of developing such models (specifically, kernels for vector valued functions) is through the process convolution approach. GP regression for vector valued output is presented in the following paragraphs; the next subsection discusses the process convolution approach and specifies the kernel functions used.

Let the number of outputs (element concentrations) that need to be simultaneously modeled be denoted by nt . Equations 4, 5 and 6 represent respectively the MOGP data fusion model, the regression estimates and their uncertainties, subject to the following modifications to the basic notation. The set $\mathbf{z} = [\mathbf{z}_1 , \mathbf{z}_2 , \mathbf{z}_3 , \dots , \mathbf{z}_{nt}]'$ represents the output values of the selected training data (element concentrations) from the individual nt outputs (elements) that need to be simultaneously modeled. The term $X = [X_1 , X_2 , X_3 , \dots , X_{nt}]$ denotes the input location values (east, north, depth) of the selected training data for each of the outputs. Any kernel [Rasmussen and Williams, 2006] may be used and even different kernel could be used for different data sets using the technique demonstrated in [Melkumyan and Ramos, 2011] (for stationary kernel) or the convolution process technique demonstrated in [Higdon, 2002, Boyle and Frean, 2004, Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2010b, 2011] and in this paper (for both stationary and nonstationary kernel). The covariance matrix of the training data

is given by

$$K(X, X) \equiv \begin{bmatrix} K_{11}^Y & K_{12}^Y & \dots & K_{1nt}^Y \\ K_{21}^Y & \dots & \dots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ K_{nt1}^Y & \dots & \dots & K_{ntnt}^Y \end{bmatrix}$$

where $K_{ii}^Y = K_{ii}^U(X_i, X_i) + \sigma_i^2 I$ and $K_{ij}^Y = K_{ij}^U(X_i, X_j)$ represent the auto-covariance of the i^{th} data set with itself and the cross covariance between the i^{th} and j^{th} data sets respectively. These terms model the covariance between the noisy observed data points (z values). Thus, they also take the noise components of the individual data sets / GPs into consideration. The corresponding noise free terms are respectively given by K_{ii}^U and K_{ij}^U . These are derived by using the process convolution approach to formulating Gaussian processes; details of this follow in the next subsection. The covariance matrix between the test points and training points is given by

$$K(X_*, X) = [K_{i1}^U(X_*, X_1), K_{i2}^U(X_*, X_2), \dots, K_{int}^U(X_*, X_{nt})],$$

where $i \in \{1 \dots nt\}$ is the GP that is being evaluated given all other GPs. The matrix $K(X, X_*)$ is defined likewise. Finally, the covariance of the test points is given by

$$K(X_*, X_*) = K_{ii}^U(X_*, X_*) + \sigma_i^2 I,$$

assuming the i^{th} GP needs to be evaluated for the particular test point. The mean and variance of the concentration estimate can thus be obtained by applying Equations 5 and 6, after incorporating multiple outputs, multiple GP/noise hyperparameters and deriving appropriate auto and cross covariances functions that model the spatial correlation between the individual data sets. Data fusion is thus achieved in the MOGP approach by correlating individual outputs and using this correlation information to improve the prediction estimates of each of them.

3.2 Convolved GPs, auto and cross covariance functions for vector valued functions

The process convolution approach to modeling GPs, proposed in [Higdon, 2002] models a GP as the convolution of a smoothing kernel and a Gaussian white noise process. Thus, the trend and complexity of the data is captured by the parameters of the smoothing kernel. In the case of vector valued outputs, each output (element concentration) can be assumed to be modeled by a separate convolved GP; the

covariance function between two outputs can be derived as a function of the respective smoothing kernel. The work [Higdon, 2002] expressed a relationship between the smoothing kernel and the corresponding covariance function through the Fourier transform and noted that for stationary isotropic kernels, there existed a one-to-one relationship between the covariance function and its smoothing kernel and that for non-isotropic and/or non-stationary kernels, there was no unique solution to the smoothing kernel. The paper suggested that the smoothing kernel for a covariance function could be obtained as the Inverse Fourier Transform of the square root of the spectrum (Fourier transform) of the covariance function. The process convolution approach to multi-output GPs has been used with the stationary SQEXP kernel in [Boyle and Frean, 2004, Álvarez and Lawrence, 2009, Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2010b, Melkumyan and Ramos, 2011] and the nonstationary NN kernel in [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2011, Vasudevan, 2012]. Given the smoothing kernel of the covariance functions in consideration, the cross-covariance terms can be derived as a kernel correlation between the respective smoothing kernels as demonstrated in [Higdon, 2002, Boyle and Frean, 2004, Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2011, Vasudevan, 2012, Melkumyan and Ramos, 2011]. A more detailed discussion of the approach can be obtained from the technical report version of this paper [Vasudevan, Melkumyan, and Scheduling, 2012] and the aforementioned references.

Assume two GPs $N(0, k_i)$ and $N(0, k_j)$, with length scale matrices Σ_i and Σ_j . Based on [Boyle and Frean, 2004], the cross and auto covariances for the stationary SQEXP kernel are given by Equations 8 and 9 respectively. The corresponding expressions for the nonstationary NN kernel are derived in [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2011, Vasudevan, 2012] and given in Equations 10 and 11 respectively. For the Matern 3/2 kernel, the expressions for the cross covariance and auto covariance are derived in [Melkumyan and Ramos, 2011] and given in Equations 12 and 13 respectively. Based on [Melkumyan and Ramos, 2011], the cross covariance function between a SQEXP and Matern 3/2 kernel is given by Equation 14.

$$K_{ij}^U(x, x') = K_f(i, j) \frac{(2\pi)^{\frac{d}{2}}}{|\Sigma_i + \Sigma_j|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - x')^T \Sigma_{ij} (x - x')\right) \quad (8)$$

where

$$\Sigma_{ij} = \Sigma_i(\Sigma_i + \Sigma_j)^{-1}\Sigma_j = \Sigma_j(\Sigma_i + \Sigma_j)^{-1}\Sigma_i$$

$$K_{ii}^U(x, x') = K_f(i, i) \frac{(\pi)^{\frac{d}{2}}}{|\Sigma_i|^{\frac{1}{2}}} \exp\left(-\frac{1}{4}(x - x')^T \Sigma_i (x - x')\right) \quad (9)$$

$$K_{ij}^U(x, x') = K_f(i, j) 2^{\frac{d+1}{2}} \frac{|\Sigma_i|^{\frac{1}{4}} |\Sigma_j|^{\frac{1}{4}}}{|\Sigma_i + \Sigma_j|^{\frac{1}{2}}} k_{NN}(\mathbf{x}, \mathbf{x}', \Sigma_{ij}) \quad (10)$$

where

$$\Sigma_{ij} = 2 \Sigma_i (\Sigma_i + \Sigma_j)^{-1} \Sigma_j$$

$$K_{ii}^U(x, x') = K_f(i, i) k_{NN}(\mathbf{x}, \mathbf{x}', \Sigma_i) \quad (11)$$

$$K_{ij}^U(x, x') = K_f(i, j) \prod_{1 \leq k \leq d} \frac{2l_{ik}^{\frac{1}{2}} l_{jk}^{\frac{1}{2}}}{l_{ik}^2 - l_{jk}^2} \left(l_{ik} e^{-\sqrt{3} \frac{r_k}{l_{ik}}} - l_{jk} e^{-\sqrt{3} \frac{r_k}{l_{jk}}} \right) \quad (12)$$

where $k \in 1 \dots d$ is the dimension of the input data, l_i and l_j are the length scales for the two Matern 3/2 kernel based GPs i and j , l_{ik} and l_{jk} are the k^{th} length scales (corresponding to the k^{th} dimensions) of these GPs and $r_k = |x_k - x'_k|$ is the distance in the k^{th} dimension between the input data.

$$K_{ii}^U(x, x') = K_f(i, i) k_{MATERN3}(\mathbf{x}, \mathbf{x}', \Sigma_i) \quad (13)$$

$$K_{ij}^U(x, x') = K_f(i, j) \prod_{1 \leq k \leq d} \sqrt{\lambda_k} \left(\frac{\pi}{2} \right)^{1/4} e^{\lambda_k^2} \left[2 \cosh \left(\frac{\sqrt{3} r_k}{l_{Mk}} \right) - e^{\frac{\sqrt{3} r_k}{l_{Mk}}} \operatorname{erf} \left(\lambda_k + \frac{r_k}{l_{SEk}} \right) - e^{-\frac{\sqrt{3} r_k}{l_{Mk}}} \operatorname{erf} \left(\lambda_k - \frac{r_k}{l_{SEk}} \right) \right] \quad (14)$$

where $\lambda_k = \frac{\sqrt{3}}{2} \frac{l_{SEk}}{l_{Mk}}$, $\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$, $k \in 1 \dots d$ is the dimension of the input data, l_{SE} and l_M are the respective length scales for the SQEXP and Matern 3/2 kernel based GPs i and j , l_{SEk} and l_{Mk} are the k^{th} length scales (corresponding to the k^{th} dimensions) of these GPs and $r_k = |x_k - x'_k|$ is the distance in the k^{th} dimension between the input data.

In Equations 10 and 11, the term, $k_{NN}(\mathbf{x}, \mathbf{x}', \Sigma_{ij})$, is the NN kernel for two data \mathbf{x} , \mathbf{x}' and length scale matrix Σ_{ij} . It is given by Equation 2, excluding the signal variance term (σ_f^2). Likewise, in Equation 13, $k_{MATERN}(\mathbf{x}, \mathbf{x}', \Sigma_i)$ refers to the Matern 3/2 kernel for two data \mathbf{x} , \mathbf{x}' and length scale matrix Σ_{ij} , given by Equation 3 (excluding the σ_f^2 term). The K_f terms in Equations 8, 9, 10 and 11 are inspired by [Bonilla, Chai, and Williams, 2007]. This term models the similarity between individual outputs and also incorporates the signal variance. It is a symmetric matrix of size $nt \times nt$ and is learnt along with the other GP hyperparameters. Thus, the hyperparameters of the system that need to be learnt include $(nt.(nt + 1))/2$ output similarity values, $nt \cdot 2$ or $nt \cdot 3$ length scale values respectively for the

individual SQEXP/MATERN3 or NN kernels and nt noise values corresponding to the noise in the observed data sets. Learning these hyperparameters by adapting the GP learning procedure described before (Equation 7) for multiple outputs [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2010b, 2011].

4 Experiments

4.1 Data set

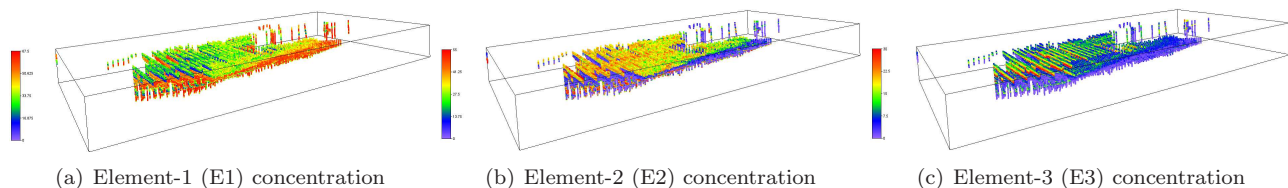


Figure 1: The geological resource data set. Figures 1(a), 1(b) and 1(c) respectively show the concentrations of three elements over the region of interest. The central region of points is surrounded by sparse sets of points which are not pre-filtered when applying the proposed algorithm.

Experiments were conducted on a large scale geological resource data set made up of real sensor data obtained from a mine site. The data consists of 63,667 measurements from a 3478.4 m x 1764.6 m x 345.9 m region in Australia that has undergone drilling and chemical assays to determine its composition. The holes are generally 25-100m apart and tens to hundreds of meters deep. Within each hole, data is collected at an interval of 2m. The measurements include the (east, north, depth) position data along with the concentrations of three elements, Element-1, Element-2 and Element-3, hereafter denoted as E1, E2 and E3 respectively. The names of these elements have been withheld at the request of the sponsor of this work. These three elements are known to be correlated and hence the objective is to use each of their GP models to improve the others' prediction estimates by capturing the correlation between these quantities. The data set is shown in Figure 1.

4.2 Testing procedure

The objective of the experiment was to compare the convolved multi-output GP approach with a conventional GP approach and quantify if the data fusion in the MOGP actually improves estimation. A second objective of the experiments was to compare the nonstationary NN kernel with the stationary SQEXP kernel, the Matern 3/2 kernel and a combination of them (Matern 3/2 - Matern 3/2 - SQEXP)

that proved effective in prior testing [Melkumyan and Ramos, 2011]. Towards these aims, a ten fold cross validation experiment was performed on the data set, with each of the kernels. This was motivated by the work [Kohavi, 1995], which suggests a ten fold stratified (similar number of samples in each fold) cross validation as the best way of testing the estimation accuracy of machine learning methods on real world data sets.

The MOGP and simple GP approaches each require an optimization step for model learning. The optimization step in each method can result in different local minima in each trial (and with each kernel). Thus, to do a one-on-one comparison between the two approaches and quantify their relative performances, an exact comparison is required. The performance comparison experiment presented in this paper provides an *exact* comparison between the MOGP and GP approaches. To do this,

- The best available MOGP parameters were found for each kernel. From this, appropriate subsets of the parameters were chosen for the GP approach. The idea of the exact comparison is to use separate GPs with parameters derived from the MOGP parameters and test the effect of discarding the other correlated data sources.
- The approaches were compared on identical test points and identical training/evaluation points selected for each of the test points.
- It is also necessary that the covariance function for the simple GP approach *must* be identical to the auto-covariance function of the DGP approach. For this reason, the auto-covariance function (for both kernels) is used as the covariance function for the GP approach to data fusion.

In addition to this, three independently optimized GPs (denoted as GPI here after) were optimized for E1, E2 and E3 and their estimates for the same set of test points were also compared. Thus the effect of information integration in the context of the geological resource modeling can be seen in terms of both an exact comparison (MOGP vs GP) and an independent comparison (MOGP vs GPI).

For the ten fold cross validation, a “block” sampling technique (see Figure 2) was used, a 3D version of the “patch” sampling method used in [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2009]. The idea was that rather than selecting test points uniformly, blocks of data test the robustness of the approach better as the support points to the query point are situated farther away (outside the block) than in uniform point selection. The GP models are used to predict (interpolate) the concentrations at all points within a test block; the larger the block size, the farther the support data to make predictions for points within the block. The data set is gridded into blocks of different sizes. Collections of blocks

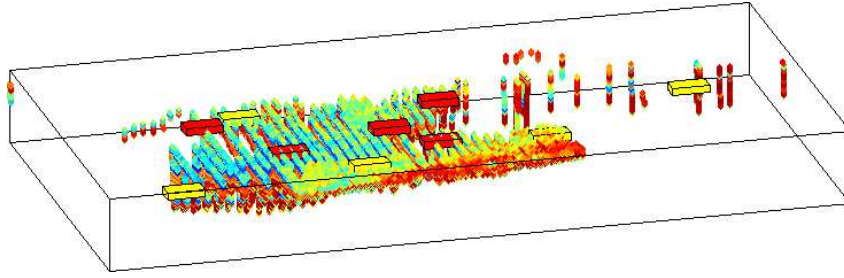


Figure 2: Example of 3D block sampling of a geological resource data set. Blocks may be sampled of different sizes. The red and yellow blocks represent blocks from two of the ten folds used in cross validation testing. Test points within these blocks have “support” data away from them, outside the blocks. This sampling method is therefore a stronger test of the robustness of an approach to estimating the quantity of interest, as compared to uniformly sampling test points. The estimation errors however, will be higher than that obtained for a uniformly sampled set of points.

represent individual folds. In each cross validation test, one fold was designated as a test fold and points from it were used exclusively for testing. All other folds together constituted the evaluation data, a small subset of which were labeled as the training data. Note that this technique of testing will naturally lead to larger errors. For the test fold, the E1, E2 and E3 concentrations (and error metrics defined in the following section) are estimated first using the MOGP approach, then with the GP approach for each of the three elements using parameters derived from the MOGP parameters and finally, with an independently optimized GP for each of the three quantities.

Table 1: Block sizes tested and implications on results - 10 fold cross validation with block sampling, 63667 points spread over 3478.4 m x 1764.6 m x 345.9 m

Block size (m)	Number of points in fold with MIN test points	Number of points in fold with MAX test points	Comments on cross validation test
22 x 11 x 2	6209	6454	Most stratified cross validation Least prediction error
44 x 22 x 4	6183	6456	stratification ↓ prediction error ↑
87 x 45 x 9	5807	6739	stratification ↓ prediction error ↑
174 x 89 x 18	5133	7549	stratification ↓ prediction error ↑
348 x 177 x 35	4976	9662	stratification ↓ prediction error ↑
696 x 353 x 70	1204	10371	Least Stratified cross validation Highest prediction error

Block sizes were chosen empirically, in proportion (arbitrarily rounded up or down) to the dimensions of the whole data set and with a view of performing a stratified cross validation test. The block sizes chosen and the resulting implications on the cross validation testing are shown in Table 1. The smaller block size of 22m x 11m x 2m results in each fold having a similar number of points (i.e. numbers of points in folds with min/max test points are similar) and thus results in the most stratified cross validation test. With increasing block size, prediction error increases (support data is farther away), stratification is reduced and hence, variance in prediction error also increases. Uniform sampling of test

points may be considered as a limiting case of block sampling with the smallest block size possible.

4.3 Metrics

Multiple metrics have been used to understand the various methods being tested. They are briefly described below. These are evaluated for each test point in each fold of the cross validation test. The result would then be represented by the mean and standard deviations of all values across all folds.

1. *Squared Error (SE)*: This represents the squared difference between the predicted concentration and the known concentrations for the set of test points. The mean over the set of all test points (Mean Squared Error or MSE) is the most popular metric for the context of this paper. Referring Equations 5 and 6, for the i^{th} test point,

$$SE(i) = (\bar{f}_*(i) - z_i)^2$$

2. *Variance (VAR)*: This represents the variance (uncertainty) in the predicted concentrations for the set of test points. a lower VAR is a good outcome, only if the SE is also low. A model that has high SE and low VAR would be a poor model as this result would suggest that the model is confident of its inaccurate estimates. A better outcome would be a model with high SE and correspondingly high VAR i.e. a model that has inaccurate predictions but is also uncertain about these predictions.
3. *Negative log probability / Log loss (NLP)*: Inspired by [Rasmussen and Williams, 2006, page 23], this is a measure of the extent to which the model (including the GP model, kernel, parameters and evaluation data) explain the current test point. The lower the value of this metric, the better the model. For the i^{th} test point,

$$NLP(i) = \frac{1}{2} \log(2\pi\sigma_*^2) + \frac{(\bar{f}_*(i) - z_i)^2}{2\sigma_*(i)^2}$$

4.4 Results

Figure 3 shows the predicted concentrations of E1, E2 and E3 over the entire region of interest as well as 2D section views of this output and the uncertainty of the predictions that constitute it; these were produced using multi-output GPs using the Neural Network kernel. Tables 2, 3 and 4 show the results of the cross validation testing on the geological resource data set with the Neural Network (NN), Matern

Table 2: E1 concentration estimation; 10 fold cross validation results using block sampling of various block sizes; Multi-output GP (MOGP) vs GP derived from MOGP (GP) vs Independently optimized GP (GPI) using Neural Network (NN), Matern 3/2 (MM), Squared exponential (SQEXP) and a Matern 3/2 - Matern 3/2 - Squared Exponential (MS) kernel combination on identical test data. The error metrics are expressed in squared units (squ).

Block size (m)	Method	NN kernel			MM kernel			SQEXP kernel			MS kernel		
		SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)	SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)	SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)	SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)
22 x 11 x 2	MOGP	1.59 (7.25)	1.47 (0.28)	1.68 (2.75)	2.66 (8.66)	0.96 (0.15)	2.35 (5.01)	23.14 (55.62)	0.39 (0.10)	31.63 (76.20)	34.08 (76.16)	33.48 (1.48)	3.18 (1.14)
	GP	36.52 (86.50)	14.93 (8.23)	3.56 (3.25)	41.42 (96.79)	5.96 (5.50)	5.77 (9.97)	44.43 (122.12)	0.51 (0.74)	45.37 (103.28)	43.49 (95.17)	36.98 (5.59)	3.30 (1.27)
	GPI	41.28 (90.35)	73.23 (7.30)	3.34 (0.60)	45.26 (96.79)	76.44 (5.23)	3.38 (0.62)	52.96 (107.62)	86.80 (4.74)	3.45 (0.61)	45.26 (96.79)	76.44 (5.23)	3.38 (0.62)
44 x 22 x 4	MOGP	1.86 (9.49)	1.81 (0.59)	1.74 (2.66)	2.79 (10.51)	1.17 (0.24)	2.19 (4.53)	24.29 (57.54)	0.48 (0.20)	28.69 (68.94)	39.49 (87.09)	34.57 (1.71)	3.26 (1.27)
	GP	52.75 (124.18)	26.55 (17.09)	3.60 (2.72)	65.00 (149.17)	13.91 (14.52)	5.05 (7.27)	76.89 (235.51)	1.07 (3.92)	50.71 (113.15)	57.55 (124.37)	40.75 (10.77)	3.46 (1.48)
	GPI	55.81 (119.69)	81.28 (11.10)	3.45 (0.71)	58.74 (124.07)	81.80 (8.02)	3.47 (0.74)	65.30 (132.02)	89.95 (6.61)	3.53 (0.72)	58.74 (124.07)	81.80 (8.02)	3.47 (0.74)
84 x 45 x 9	MOGP	3.38 (22.92)	3.24 (5.06)	1.91 (2.56)	7.50 (28.75)	1.65 (0.47)	3.04 (6.07)	30.08 (70.98)	0.85 (0.53)	21.35 (49.50)	52.13 (109.11)	36.90 (2.07)	3.42 (1.47)
	GP	85.88 (187.61)	58.60 (39.70)	3.71 (2.00)	114.05 (242.32)	44.66 (37.46)	4.44 (4.35)	265.25 (853.00)	8.06 (25.34)	45.35 (95.72)	91.75 (190.74)	57.17 (33.92)	3.71 (1.63)
	GPI	85.53 (171.85)	100.04 (21.34)	3.63 (0.81)	86.81 (175.15)	97.42 (15.97)	3.64 (0.85)	91.63 (179.65)	101.22 (13.42)	3.67 (0.85)	86.81 (175.15)	97.42 (15.97)	3.64 (0.85)
174 x 89 x 18	MOGP	14.60 (88.59)	10.97 (27.66)	2.24 (2.70)	32.05 (96.26)	2.24 (0.61)	7.29 (15.78)	52.96 (122.09)	1.54 (0.78)	19.27 (41.43)	83.79 (166.79)	39.05 (2.41)	3.81 (2.09)
	GP	128.39 (261.10)	113.03 (88.51)	3.84 (1.60)	156.56 (306.22)	95.10 (56.46)	4.23 (3.05)	701.62 (1787.63)	34.69 (59.68)	28.36 (59.58)	154.34 (319.96)	104.53 (106.83)	3.94 (1.58)
	GPI	124.52 (235.22)	129.20 (45.85)	3.79 (0.88)	124.93 (240.35)	121.09 (26.43)	3.80 (0.95)	128.86 (244.59)	122.29 (23.44)	3.82 (0.96)	124.93 (240.35)	121.09 (26.43)	3.80 (0.95)
348 x 177 x 35	MOGP	73.42 (213.89)	64.36 (86.94)	3.00 (2.09)	112.99 (206.86)	2.74 (0.64)	19.76 (32.27)	114.47 (214.94)	2.45 (0.98)	24.18 (42.85)	155.36 (257.28)	40.86 (2.50)	4.64 (3.08)
	GP	204.57 (387.09)	249.34 (232.37)	4.06 (1.39)	215.40 (373.05)	153.09 (71.73)	4.25 (2.51)	1091.93 (2801.29)	124.02 (94.50)	16.84 (43.56)	290.14 (541.45)	329.34 (376.16)	4.23 (1.21)
	GPI	189.14 (335.76)	199.86 (120.66)	3.98 (0.90)	189.21 (327.56)	151.49 (45.31)	4.00 (1.04)	190.92 (331.84)	155.24 (44.34)	4.01 (1.03)	189.21 (327.56)	151.49 (45.31)	4.00 (1.04)
696 x 353 x 70	MOGP	180.64 (368.05)	173.67 (154.95)	3.61 (1.56)	206.18 (349.73)	2.97 (0.52)	33.83 (54.10)	214.11 (357.84)	2.89 (0.84)	37.05 (58.21)	243.80 (380.23)	41.95 (2.21)	5.66 (4.46)
	GP	325.98 (546.05)	562.94 (523.54)	4.31 (1.19)	301.60 (452.75)	192.63 (64.46)	4.43 (2.29)	871.72 (2481.90)	185.80 (93.03)	12.58 (39.08)	460.23 (775.98)	976.69 (915.85)	4.55 (1.01)
	GPI	291.72 (465.37)	362.43 (271.48)	4.19 (0.80)	282.28 (428.51)	180.05 (49.34)	4.23 (1.12)	283.43 (430.84)	183.92 (49.45)	4.23 (1.11)	282.28 (428.51)	180.05 (49.34)	4.23 (1.12)

Table 3: E2 concentration estimation; 10 fold cross validation results using block sampling of various block sizes; Multi-output GP (MOGP) vs GP derived from MOGP (GP) vs Independently optimized GP (GPI) using Neural Network (NN), Matern 3/2 (MM), Squared exponential (SQEXP) and a Matern 3/2 - Matern 3/2 - Squared Exponential (MS) kernel combination on identical test data. The error metrics are expressed in squared units (squ).

Block size (m)	Method	NN kernel			MM kernel			SQEXP kernel			MS kernel		
		SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)	SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)	SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)	SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)
22x11x2	MOGP	3.85 (17.59)	3.13 (0.53)	2.10 (2.87)	3.88 (18.47)	2.02 (0.36)	2.23 (4.74)	52.32 (231.51)	0.04 (0.10)	650.87 (1848.02)	36.09 (102.25)	39.09 (2.48)	3.20 (1.27)
	GP	36.86 (118.99)	19.44 (8.89)	3.34 (3.10)	40.79 (129.77)	8.45 (6.63)	4.37 (7.41)	75.67 (869.23)	0.05 (0.30)	668.93 (1880.52)	46.60 (134.94)	43.34 (6.65)	3.31 (1.36)
	GPI	49.79 (143.90)	93.47 (7.40)	3.44 (0.71)	53.61 (151.69)	97.26 (5.97)	3.47 (0.72)	60.25 (161.56)	85.95 (5.68)	3.48 (0.86)	53.61 (151.69)	97.26 (5.97)	3.47 (0.72)
44x22x4	MOGP	4.79 (22.21)	3.80 (1.06)	2.20 (2.87)	4.72 (22.80)	2.53 (0.58)	2.29 (4.43)	88.25 (352.70)	0.10 (0.30)	657.09 (1840.03)	42.82 (117.40)	40.31 (2.88)	3.28 (1.40)
	GP	55.51 (174.24)	32.14 (18.07)	3.49 (2.70)	64.20 (192.62)	18.19 (17.10)	4.19 (5.58)	181.27 (1385.73)	0.26 (2.48)	694.20 (1901.88)	64.50 (186.52)	47.76 (12.97)	3.46 (1.58)
	GPI	68.69 (195.85)	101.11 (10.86)	3.55 (0.86)	71.21 (201.30)	102.80 (8.82)	3.56 (0.88)	77.86 (211.96)	89.38 (9.01)	3.57 (1.04)	71.21 (201.30)	102.80 (8.82)	3.56 (0.88)
84x45x9	MOGP	8.04 (40.48)	6.22 (6.00)	2.37 (2.71)	10.97 (44.78)	3.69 (1.12)	2.81 (5.02)	211.05 (753.57)	0.56 (1.15)	461.55 (1290.67)	56.72 (146.64)	42.90 (3.41)	3.44 (1.64)
	GP	95.98 (274.95)	66.23 (41.41)	3.71 (2.27)	116.98 (318.33)	54.56 (43.24)	4.05 (3.56)	1140.18 (8156.99)	4.72 (22.50)	532.52 (1420.66)	105.69 (288.37)	67.37 (42.25)	3.69 (1.70)
	GPI	105.20 (277.85)	119.02 (20.10)	3.72 (1.03)	105.27 (279.95)	119.32 (17.19)	3.72 (1.06)	115.96 (301.71)	103.63 (23.63)	3.74 (1.22)	105.27 (279.95)	119.32 (17.19)	3.72 (1.06)
174x89x18	MOGP	21.49 (102.60)	15.88 (29.16)	2.66 (2.61)	37.85 (117.51)	5.09 (1.44)	4.85 (8.77)	402.24 (1144.48)	1.92 (2.06)	228.66 (752.97)	90.32 (211.07)	45.49 (3.73)	3.79 (2.23)
	GP	142.62 (356.16)	123.42 (91.33)	3.88 (1.88)	165.30 (394.82)	112.93 (64.77)	4.07 (2.70)	3510.60 (14425.39)	25.86 (59.06)	312.01 (914.26)	170.93 (420.58)	125.06 (134.64)	3.91 (1.59)
	GPI	148.71 (347.79)	146.71 (41.34)	3.88 (1.09)	147.83 (351.45)	145.39 (28.92)	3.88 (1.12)	164.86 (376.69)	133.30 (41.47)	3.92 (1.25)	147.83 (351.45)	145.39 (28.92)	3.88 (1.12)
348x177x35	MOGP	82.02 (233.18)	72.71 (90.66)	3.23 (2.04)	119.86 (236.67)	6.26 (1.52)	10.28 (16.06)	419.72 (1196.06)	4.61 (2.90)	134.53 (599.11)	167.10 (320.09)	48.08 (4.15)	4.54 (3.23)
	GP	219.37 (484.76)	265.84 (239.30)	4.09 (1.71)	232.89 (475.45)	178.16 (81.64)	4.17 (2.19)	8397.84 (28045.42)	114.38 (102.39)	164.10 (633.41)	314.92 (689.64)	414.05 (500.88)	4.23 (1.24)
	GPI	213.44 (442.90)	213.16 (107.76)	4.05 (1.09)	216.33 (443.91)	182.88 (51.20)	4.06 (1.18)	234.24 (459.31)	194.54 (59.86)	4.11 (1.17)	216.33 (443.91)	182.88 (51.20)	4.06 (1.18)
696x353x70	MOGP	196.72 (379.37)	189.18 (162.78)	3.75 (1.48)	227.42 (370.75)	6.86 (1.26)	17.23 (24.68)	420.95 (1016.76)	6.14 (2.62)	94.13 (458.35)	273.23 (440.21)	49.71 (4.07)	5.58 (4.37)
	GP	340.01 (603.14)	594.95 (544.29)	4.34 (1.33)	331.10 (534.42)	223.12 (73.21)	4.37 (1.77)	5910.35 (23857.18)	189.13 (106.28)	105.40 (514.02)	493.28 (924.16)	1309.98 (1309.73)	4.59 (0.96)
	GPI	314.51 (529.84)	365.16 (250.80)	4.24 (0.91)	317.37 (519.71)	218.58 (58.44)	4.27 (1.15)	331.32 (531.79)	236.80 (60.61)	4.32 (1.13)	317.37 (519.71)	218.58 (58.44)	4.27 (1.15)

Table 4: E3 concentration estimation; 10 fold cross validation results using block sampling of various block sizes; Multi-output GP (MOGP) vs GP derived from MOGP (GP) vs Independently optimized GP (GPI) using Neural Network (NN), Matern 3/2 (MM), Squared exponential (SQEXP) and a Matern 3/2 - Matern 3/2 - Squared Exponential (MS) kernel combination on identical test data. The error metrics are expressed in squared units (squ).

Block size (m)	Method	NN kernel			MM kernel			SQEXP kernel			MS kernel		
		SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)	SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)	SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)	SE (squ) mean (std)	VAR (squ) mean (std)	NLP mean (std)
22 x 11 x 2	MOGP	1.60 (6.82)	1.29 (0.35)	1.64 (2.49)	1.29 (5.42)	0.98 (0.14)	1.57 (2.84)	7.30 (20.47)	0.64 (0.10)	6.50 (16.50)	19.15 (45.80)	36.32 (0.17)	2.98 (0.63)
	GP	9.09 (26.21)	3.84 (1.76)	2.83 (3.67)	9.65 (27.41)	2.36 (1.44)	3.57 (6.46)	10.47 (29.99)	0.75 (0.40)	8.21 (21.19)	19.18 (45.75)	36.47 (0.19)	2.98 (0.63)
	GPI	9.69 (27.69)	14.64 (3.24)	2.58 (0.90)	11.19 (31.47)	16.90 (1.16)	2.66 (0.94)	12.06 (33.46)	18.46 (1.15)	2.70 (0.91)	12.06 (33.46)	18.46 (1.15)	2.70 (0.91)
44 x 22 x 4	MOGP	1.90 (8.29)	1.68 (0.66)	1.71 (2.17)	1.50 (7.41)	1.19 (0.23)	1.62 (2.92)	7.99 (21.94)	0.73 (0.19)	6.49 (16.04)	22.53 (52.33)	36.35 (0.20)	3.03 (0.72)
	GP	12.80 (37.03)	6.38 (3.66)	2.91 (3.17)	14.47 (40.66)	4.52 (3.53)	3.50 (5.34)	15.59 (45.59)	1.07 (1.45)	9.56 (23.82)	22.56 (52.27)	36.51 (0.23)	3.03 (0.72)
	GPI	12.88 (36.55)	16.79 (4.65)	2.69 (1.00)	14.16 (39.23)	18.14 (1.79)	2.76 (1.08)	14.92 (40.78)	19.41 (1.73)	2.79 (1.05)	14.92 (40.78)	19.41 (1.73)	2.79 (1.05)
84 x 45 x 9	MOGP	2.87 (13.66)	2.79 (1.97)	1.88 (2.06)	3.08 (12.59)	1.65 (0.42)	1.99 (3.53)	9.19 (24.11)	1.09 (0.48)	5.46 (12.14)	29.82 (63.83)	36.43 (0.31)	3.13 (0.88)
	GP	20.35 (55.84)	13.22 (8.54)	3.01 (2.45)	24.26 (64.08)	12.01 (8.52)	3.33 (3.63)	35.53 (99.80)	3.73 (7.16)	9.73 (21.78)	29.82 (63.64)	36.61 (0.36)	3.13 (0.87)
	GPI	19.56 (51.95)	21.69 (8.26)	2.86 (1.08)	20.50 (53.15)	21.61 (3.60)	2.92 (1.19)	21.21 (54.63)	22.73 (3.76)	2.94 (1.17)	21.21 (54.63)	22.73 (3.76)	2.94 (1.17)
174 x 89 x 18	MOGP	6.63 (35.07)	6.00 (8.73)	2.17 (2.24)	8.67 (24.92)	2.19 (0.54)	3.07 (4.93)	14.56 (34.46)	1.73 (0.70)	5.51 (10.69)	39.96 (80.47)	36.66 (0.54)	3.27 (1.10)
	GP	29.91 (81.55)	25.03 (19.65)	3.13 (2.01)	32.57 (79.36)	23.31 (12.64)	3.28 (2.58)	73.85 (189.74)	12.06 (14.90)	7.34 (14.06)	39.81 (79.82)	36.89 (0.66)	3.26 (1.08)
	GPI	28.27 (69.54)	29.28 (16.01)	3.01 (1.10)	28.87 (69.91)	26.65 (6.01)	3.07 (1.27)	29.50 (72.11)	28.35 (6.17)	3.09 (1.24)	29.50 (72.11)	28.35 (6.17)	3.09 (1.24)
348 x 177 x 35	MOGP	22.76 (76.19)	23.57 (36.17)	2.64 (1.91)	25.66 (50.80)	2.62 (0.57)	5.81 (8.42)	27.36 (55.29)	2.58 (0.89)	6.50 (10.74)	54.46 (99.06)	37.62 (1.59)	3.46 (1.32)
	GP	48.12 (113.43)	55.21 (53.40)	3.32 (1.82)	47.08 (100.71)	34.95 (15.85)	3.39 (2.25)	96.50 (228.76)	33.92 (22.17)	5.49 (10.17)	54.31 (98.51)	38.13 (2.22)	3.45 (1.30)
	GPI	43.19 (97.09)	46.90 (36.64)	3.20 (1.15)	42.96 (92.63)	32.36 (10.28)	3.25 (1.40)	43.24 (93.69)	34.15 (10.70)	3.25 (1.35)	43.24 (93.69)	34.15 (10.70)	3.25 (1.35)
696 x 353 x 70	MOGP	50.93 (123.14)	81.32 (99.31)	3.10 (1.43)	43.57 (77.06)	2.86 (0.48)	8.55 (12.36)	46.94 (81.94)	3.00 (0.77)	8.98 (13.24)	66.54 (109.25)	40.95 (5.93)	3.58 (1.32)
	GP	73.71 (146.92)	129.87 (126.46)	3.57 (1.65)	64.20 (112.93)	43.42 (14.04)	3.58 (2.05)	93.45 (205.91)	47.21 (21.02)	4.85 (8.96)	67.03 (109.94)	43.80 (11.83)	3.57 (1.27)
	GPI	65.75 (122.28)	83.71 (72.39)	3.41 (1.03)	61.34 (108.31)	38.04 (10.97)	3.48 (1.39)	61.26 (108.52)	39.36 (10.98)	3.47 (1.34)	61.26 (108.52)	39.36 (10.98)	3.47 (1.34)

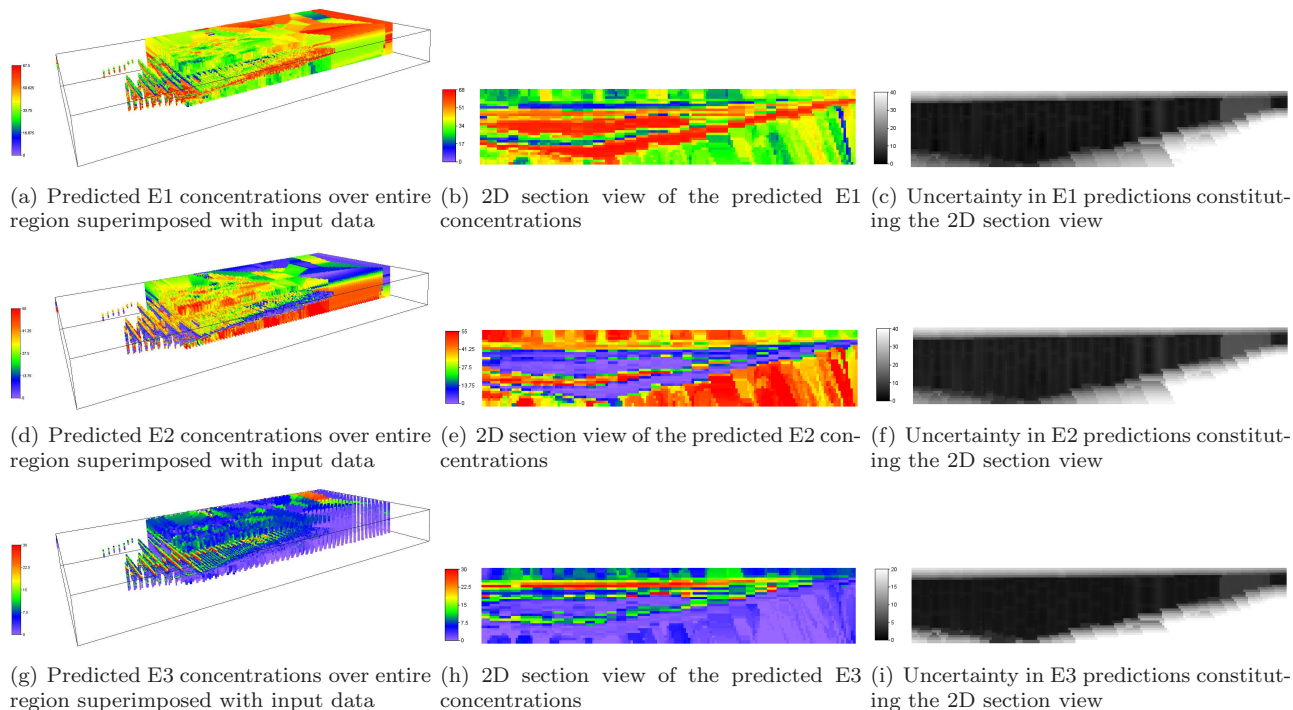


Figure 3: Figures 3(a), 3(b) and 3(c) respectively show the predicted E1 concentrations (over the entire region) superimposed with the input data, a 2D section-view of the output data and the uncertainty in the predicted concentrations for the 2D view. The corresponding figures for E2 and E3 are 3(d), 3(e), 3(f) and 3(g), 3(h), 3(i) respectively. Expectedly, the uncertainty is low around regions where input/given data exist and rapidly rises for predictions away from such areas - typically, the fringe areas. The 2D section view (Figure 3(b)) shows two red regions corresponding two regions of high E1 concentration. The corresponding regions in Figures 3(e) and 3(h) show low E2 and E3 concentrations respectively.

3/2 (MM), Squared Exponential (SQEXP) and Matern 3/2 - Matern 3/2 - Squared exponential (MS) kernels. Trends from the three tables are visualized through numerous graphs located in the appendix and summarized below. For a more complete set and large sized graphs, see [Vasudevan, Melkumyan, and Scheduling, 2012].

1. Prediction error (SE) increases with increase in test block size.

- See Figures 4(a), 4(d), 4(f), 4(h) for E1, 5(a), 5(d), 5(f), 5(h) for E2 and 6(a), 6(d), 6(f), 6(h) for E3.
- This behavior is expected. It happens because the support training data required for regressing at a test point is situated farther away. Increasing the test block size also results in reduced stratification as one fold of the cross validation may have e.g. 10,000 test points whereas another may have only 1000 points. This results in increased standard deviation of prediction error. A ten fold stratified cross validation is generally considered to be the most representative of performance measure [Kohavi, 1995], however testing multiple larger block sizes provides a better understanding of the model's behavior and robustness.

2. NN kernel based MOGP/GP models trained faster than other kernels

- Further optimization of each of the MOGP/GP models could yield better results. The results shown are the result of a reasonable amount of optimization applied to each kernel and GP model. Typically, multiple attempts were performed and the best results obtained were pursued/used. One iteration consisted of a stochastic optimization step (simulated annealing) and/or a gradient based optimization step (Quasi Newton optimization with BFGS Hessian update) with 10,000 training data chosen uniformly from the data. This work uses a “block-learning” approximation [Vasudevan, Ramos, Nettleton, and Durrant-Whyte, 2010b] which approximates the total marginal likelihood as a sum of a sequence of marginal likelihoods computed over blocks of points comprising the training data. The size of the block is defined by the computational resources available. The stochastic optimization step was the most time consuming part; each attempt was started with completely random parameters. The code was unoptimized MATLAB code running typically on an 8-core processor based machine. Most times, not all the cores were used for the same process; multiple processes also shared the same system. Note that the experiments in this paper do not use analytical gradients for the optimization of the hyperparameters; this was a design choice made in the interest of stability and comparability of the optimization results across kernels. The use of analytical gradients can significantly reduce the total training time. Training time may also be reduced significantly by various other ways including other approximations, intelligently setting initial parameters, scaling the data etc.

Model	Kernel	Number of training attempts, iterations Total training time for successful attempt
MOGP	NN	2 attempts, 3 iterations, total training time = 78.89 hours
	MM	3 attempts, 2 iterations, total training time = 222.15 hours
	SQEXP	4 attempts, 1 iteration, total training time = 92.52 hours
	MS	2 attempts, 3.5 iterations, 3 iterations took 113.69 hours
GPI	NN	3 attempts, 2 iterations, total training time = 41.07 hours
	MM	2 attempts, 1 iteration, total training time = 48.91 hours
	SQEXP	2 attempts, 1 iteration, total training time = 47.67 hours

Rather than the individual training times, the relative amount of training (under similar con-

ditions, with different kernel) required to produce a reasonable set of parameters is of more interest. Experience suggests that the NN kernel based MOGP/GP models converged faster and better as compared to other kernels.

3. MOGP models based on the NN kernel outperform other kernels tested.

- See Figures 4(a), 4(b), 4(c) for E1, 5(a), 5(b), 5(c) for E2 and 6(a), 6(b), 6(c) for E3.
- The NN kernel is the best performing kernel of the four tested, across all block sizes tested. The MOGP based on the NN kernel produces lower SE (better estimate) and reduced NLP (better model) estimates than other kernels tested.
- For small block sizes, both the NN and MM kernel are competitive; in case of E3, the MM even marginally outperforms the NN kernel for the two smallest block sizes tested. Note however that considering all test sizes and all three elements, the observation is that the MM kernel produces lower VAR for a higher SE, meaning that it is more confident of its SE values which are worse/higher than those of the NN kernel. This makes its NLP higher and the model poorer than an MOGP based on the NN kernel. Note also that as the test block size increases, the advantage in performance of the MOGP based on the NN kernel over that based on the MM kernel becomes more distinctive. Not only are the SE values smaller for the NN kernel, the NLP values remain in the same range whereas those of the MM kernel rise significantly. This proves that the MOGP-NN is better performing and more robust than the MOGP-MM. The latter property suggests that the MOGP-NN will be able to cope better with incomplete data sets.
- Both the MS and SQEXP kernels are not competitive with respect to the NN or MM kernels considering both the SE and NLP metrics. These kernels are discussed individually in the following paragraphs.

4. MOGP models perform significantly better than three separate GPs (using the MOGP parameters) or three independently optimized GPs as information fusion improves estimation.

- See Figures 4(d) and 4(e) for E1, 5(d) and 5(e) for E2 and 6(d) and 6(e) for E3.
- For the NN kernel, the MOGP metrics are always lower than the corresponding derived GP (GP) or independent GP (GPI) metrics - lower SE (better estimate) with lower NLP (better

model). This clearly demonstrates the benefits of information fusion across heterogeneous information sources so as to improve individual predictions using the MOGP model.

- From Tables 2, 3 and 4, the average reduction in error (i.e. improvement in performance) of MOGP models over GP/GPI models for the smallest, intermediate and largest test block sizes are -

– E1

- * 22 x 11 x 2 - 95.6% over GP, 96.2% over GPI
- * 84 x 45 x 9 - 96.1% over GP, 96.0% over GPI
- * 696 x 353 x 70 - 44.6% over GP, 38.1% over GPI

– E2

- * 22 x 11 x 2 - 89.6% over GP, 92.3% over GPI
- * 84 x 45 x 9 - 91.6% over GP, 92.4% over GPI
- * 696 x 353 x 70 - 42.1% over GP, 37.5% over GPI

– E3

- * 22 x 11 x 2 - 82.4% over GP, 83.5% over GPI
- * 84 x 45 x 9 - 85.9% over GP, 85.3% over GPI
- * 696 x 353 x 70 - 30.9% over GP, 22.5% over GPI

These numbers demonstrate significant improvements in performance, even in very large test block sizes, when using the MOGP-NN model for correlated data.

5. The MS kernel was uncompetitive

- See Figures 4(h) and 4(i) for E1, 5(h) and 5(i) for E2 and 6(h) and 6(i) for E3.
- The MS kernel is not competitive with respect to the NN and MM kernels as discussed earlier. However, the MOGP using this kernel combination proves to be better than a derived GP and an independently optimized GP with respect to the SE metric. From the NLP perspective, the MOGP-MS model is more competitive than the other GP models for small block sizes. For larger block sizes, using an independently optimized GP proves to be a more trust worthy modeling option as the increase in error is met with a corresponding increase in uncertainty (hence low NLP) for the independent GP models. The exception to this behavior is seen in

the results for E3, the MOGP model is poor in this case. This is attributed to do with inferior parameters relevant to the element E3 obtained from the optimization process.

- The MS kernel performs better than the SQEXP with respect to the NLP metric and hence can be trusted more (prediction error compensated by prediction uncertainty), but in two of the three elements (E1 and E3), its SE was inferior to that of the SQEXP.

6. The SQEXP kernel was uncompetitive and unreliable

- See Tables 2, 3 and 4; see Figures 4(a), 4(b) 4(c), 4(f), 4(g), 4(h) and 4(i) for E1, 5(a), 5(b), 5(c), 5(f), 5(g), 5(h) and 5(i) for E2 and 6(a), 6(b), 6(c), 6(f), 6(g), 6(h) and 6(i) for E3.
- The MOGP-SQEXP model performs poorly in comparison with the equivalent models using the NN/MM kernels, with respect to both SE and NLP.
- For elements E1 and E3, the MOGP-SQEXP has a better SE than the corresponding model based on the MS kernel; it has an SE better than the corresponding derived/independent GP models but an inferior (overconfident or low uncertainty) VAR and a fluctuating NLP trend. For element E2, the MOGP-SQEXP is worse off than both the equivalent model based on the MS kernel as well as its corresponding GP models.
- Considering the results for E2, the NLP is directly proportional to the SE and inversely to the prediction variance. At the smallest block size, the MOGP-SQEXP produces relatively high SE (with respect to e.g. MOGP-NN) but very low prediction variance. This basically suggests that the model is confident of its poor estimates - a bad outcome. This results in a high NLP and poor model. As the block size increases, the prediction variance increases more relative to the prediction error resulting in the decreasing NLP trend. For elements E1 and E3, the largest block size results in a stronger increase in prediction error than the variance in the prediction resulting in an increase in NLP. Overall, the MOGP-SQEXP model is poor.
- The SQEXP kernel is a limiting case of the MM kernel; both are stationary kernels. Considering the behavior of the GPI model using the SQEXP kernel and its competitive results with respect to those of the GPI-MM kernel, it is possible that the poor performance of the MOGP-SQEXP (as compared to the MOGP-MM) is due to poor optimization output (a bad local minima).

7. In general, the stationary kernels tested seemed to have an inadequate increase in prediction uncertainty with increasing test block size and worsening predictions. This leads a higher NLP metric

and a poor model that is overly confident of its worsening predictions. This behavior can be attributed to the correlation profile of the stationary kernels tested - they all share the “correlation decreases with increasing distance of support data from point of interest” trend. This results in stationary kernels not being able to cope with large test block sizes as the support data is farther away (i.e. less correlated and not of much use). In contrast, the nonstationary NN kernel has a sigmoidal profile that can handle this issue across a range of test block sizes.

8. The SE metric taken alone can be misleading. The experiments have reinforced the need for a multi-metric analysis. The SE metric only provides information on the prediction error but it does not describe the prediction uncertainty which is very important in understanding if a model is reliable or otherwise. The VAR and NLP metrics provided key insights on the difference in performance between different models and kernels. A model that is very confident of its poor predictions is unreliable (as was the case for the SQEXP kernel). Worsening predictions (due to increasing test block size) is itself not a bad outcome, provided it is met with an equivalent increase in prediction uncertainty.
9. Further discussion of the results in the context of broader practical issues like determining if a MOGP model is indeed good or identifying the kind of GP model suited for a particular task, are available in the technical report version of this paper [Vasudevan, Melkumyan, and Scheduling, 2012].

5 Conclusion

This paper empirically studied the problem of geological resource modeling using a machine learning approach to convolved multi-output Gaussian processes (MOGPs). The concentrations of three elements were modeled and predicted over a region of interest using multi-output Gaussian processes (MOGPs; joint modeling of multiple outputs) as well as individual GPs for each of these quantities separately. The paper demonstrates that MOGPs perform significantly better than individual GPs at the modeling problem as they effectively integrate heterogeneous sources of information (concentrations of individual elements) to improve the predictions of each of them. The benefits of information integration using the MOGP as against independent GPs for the task of geological resource modeling have been quantified by a multi-metric, multi-kernel and multi-test-size cross validation study that performed both an exact and an independent comparison between MOGPs and GPs. Multi-output Gaussian process models based

on the Neural Network kernel was shown to be a competitive and robust modeling option across a range of test block sizes.

Acknowledgements

This work has been funded by the Rio Tinto Centre for Mine Automation.

References

- M. A. Álvarez and N. D. Lawrence. Sparse convolved gaussian processes for multi-output regression. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems 21*, pages 57–64. 2009.
- M. A. Álvarez, L. Rosasco, and N. D. Lawrence. Kernels for Vector-Valued Functions: A Review. *Foundations and Trends in Machine Learning*, 4:195–266, 2012.
- E. Bonilla, K. M. Chai, and C. Williams. Multi-task Gaussian process prediction. In J.C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 153–160. MIT Press, Cambridge, MA, 2007.
- P. Boyle and M. Frean. Dependent Gaussian processes. In Lawrence K. Saul, Yair Weiss, and Léon Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 217–224. MIT Press, Cambridge, MA, 2004.
- N. Cressie. *Statistics for Spatial Data*. Wiley-Interscience, 1993.
- M. A. El-Beltagy and W. A. Wright. Gaussian processes for model fusion. In *International Conference on Artificial Neural Networks (ICANN)*, 2001.
- M. Girolami. Bayesian data fusion with gaussian process priors: An application to protein fold recognition. In *Workshop on Probabilistic Modeling and Machine Learning in Structural and Systems Biology (PMSB)*, 2006.
- P. W. Goldberg, C. K. I. Williams, and C. M. Bishop. Regression with Input-dependent Noise: A Gaussian Process Treatment. In M. I. Jordan, M. J. Kearns, S. A. Solla, and L. Erlbaum, editors, *Advances in Neural Information Processing Systems (NIPS) 10*. MIT Press, Cambridge, MA, 1998.

- P. Goovaerts. *Geostatistics for Natural Resources Evaluation*. Oxford University Press, 1997.
- D. Higdon. *Quantitative Methods for Current Environmental Issues*, chapter Space and Space-Time Modeling Using Process Convolutions, pages 37–54. Springer, 2002.
- K. Hornik. Some new results on neural network approximation. *Neural Networks*, 6(8):1069–1072, 1993.
- K. Kersting, C. Plagemann, P. Pfaff, and W. Burgard. Most Likely Heteroscedastic Gaussian Process Regression. In *International Conference on Machine Learning (ICML)*, 2007.
- R. Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the Fourteenth International Joint Conferences on Artificial Intelligence (IJCAI)*, volume 14, pages 1137–1145, 1995.
- G. Matheron. Principles of Geostatistics. *Economic Geology*, 58:1246–1266, 1963.
- A. Melkumyan and F. Ramos. A sparse covariance function for exact gaussian process inference in large data sets. In *International Joint Conferences on Artificial Intelligence*, volume 21, pages 1936–1942, 2009.
- A. Melkumyan and F. Ramos. Multi-Kernel Gaussian Processes. In *in the proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 2011.
- R. Murray-Smith and B. A. Pearlmutter. *Deterministic and Statistical Methods in Machine Learning, LNAI 3635*, chapter Transformations of Gaussian Process priors, pages 110–123. Springer-Verlag, 2005.
- R. M. Neal. *Bayesian Learning for Neural Networks*. Lecture Notes in Statistics 118. Springer, New York, 1996.
- C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- S. Reece, S. Roberts, D. Nicholson, and C. Lloyd. Determining intent using hard/soft data and gaussian process classifiers. In *Proceedings of the 14th International Conference on Information Fusion (FUSION)*, 2011.
- S. Vasudevan. Data fusion using gaussian processes. *Elsevier Journal of Robotics and Autonomous Systems*, 2012. URL <http://dx.doi.org/10.1016/j.robot.2012.08.006>. Available online 25 August 2012.

- S. Vasudevan, F. Ramos, E. Nettleton, and H. Durrant-Whyte. Gaussian Process Modeling of Large Scale Terrain. *Journal of Field Robotics*, 26(10):812–840, 2009.
- S. Vasudevan, F. Ramos, E. Nettleton, and H. Durrant-Whyte. Heteroscedastic Gaussian processes for data fusion in large scale terrain modeling. In *the International Conference for Robotics and Automation (ICRA)*, 2010a.
- S. Vasudevan, F. Ramos, E. Nettleton, and H. Durrant-Whyte. Large-scale terrain modeling from multiple sensors with dependent Gaussian processes. In *in the proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Taipei, October 2010b.
- S. Vasudevan, F. Ramos, E. Nettleton, and H. Durrant-Whyte. Non-stationary dependent Gaussian processes for data fusion in large scale terrain modeling. In *in the proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, Shanghai, China, 2011.
- S. Vasudevan, A. Melkumyan, and S. Scheduling. Information fusion in multi-task Gaussian process models. Technical report, Australian Centre for Field Robotics, The University of Sydney, 2012. arXiv report 1210.1928, available online at <http://arxiv.org/abs/1210.1928>.
- H. Wackernagel. *Multivariate geostatistics: an introduction with applications*. Springer, 2003.
- C. K. I. Williams. Computation with infinite neural networks. *Neural Computation*, 10(5):1203–1216, 1998a.
- C. K. I. Williams. Prediction with Gaussian processes: From linear regression to linear prediction and beyond. In M. I. Jordan, editor, *Learning in Graphical Models*, pages 599–622. Springer, 1998b.

Appendix: Graphs of results obtained in Tables 2, 3 and 4

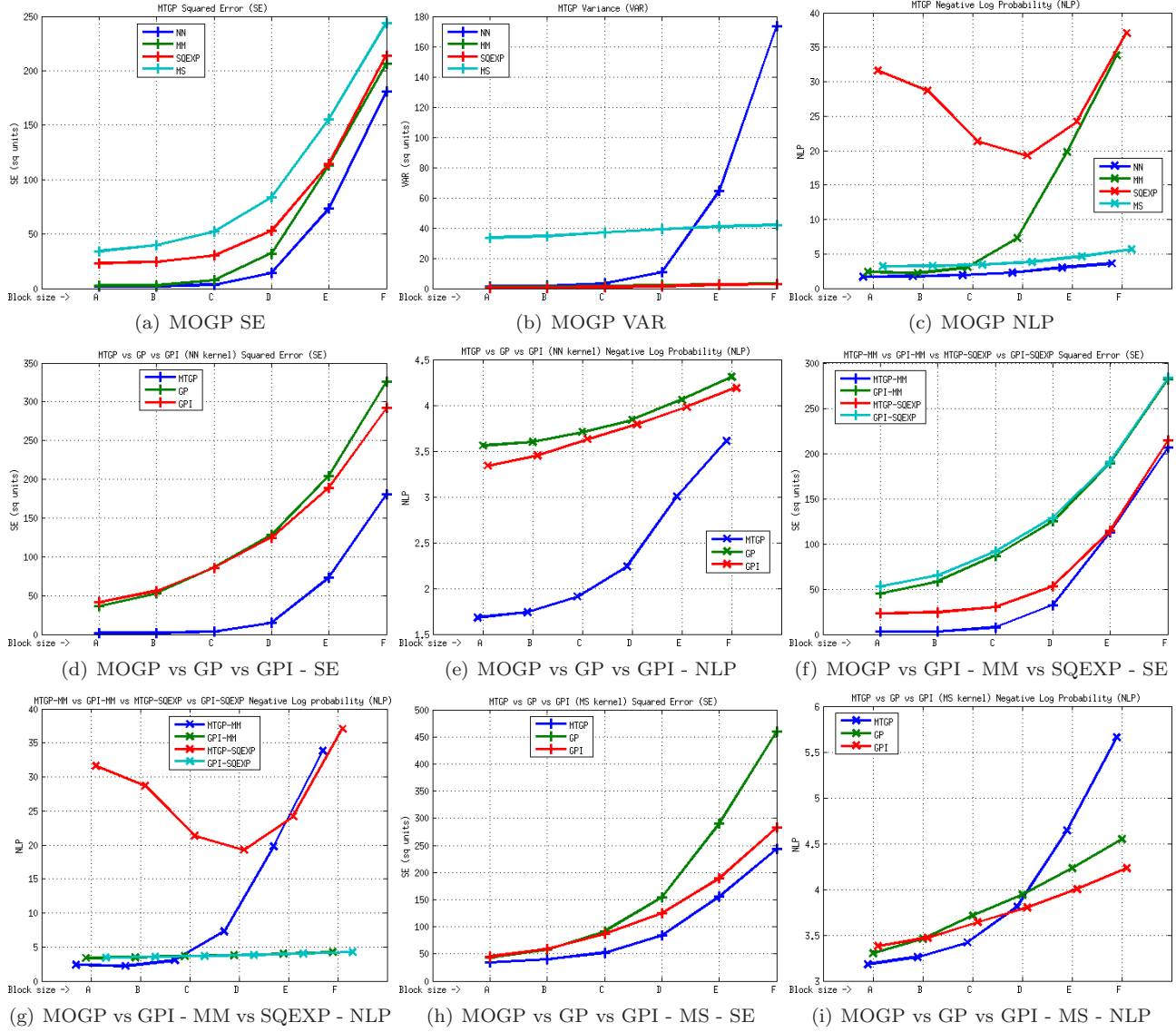


Figure 4: Element E1 - key trends. Test block sizes (m) - A (22 x 11 x 2), B (44 x 22 x 4), C (84 x 45 x 9), D (174 x 89 x 18), E (348 x 177 x 35) and F (696 x 353 x 70).

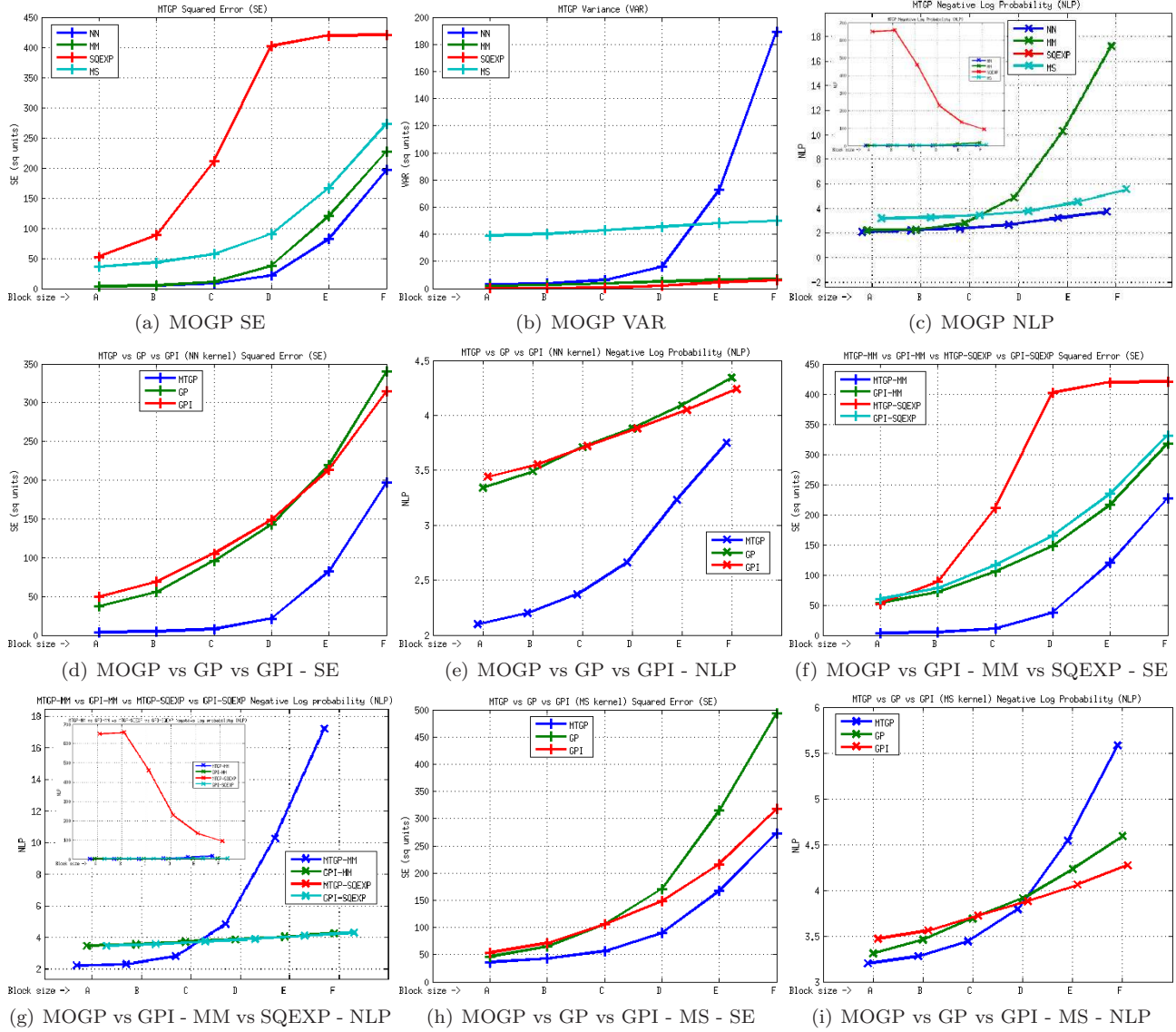


Figure 5: Element E2 - key trends. Test block sizes (m) - A (22 x 11 x 2), B (44 x 22 x 4), C (84 x 45 x 9), D (174 x 89 x 18), E (348 x 177 x 35) and F (696 x 353 x 70).

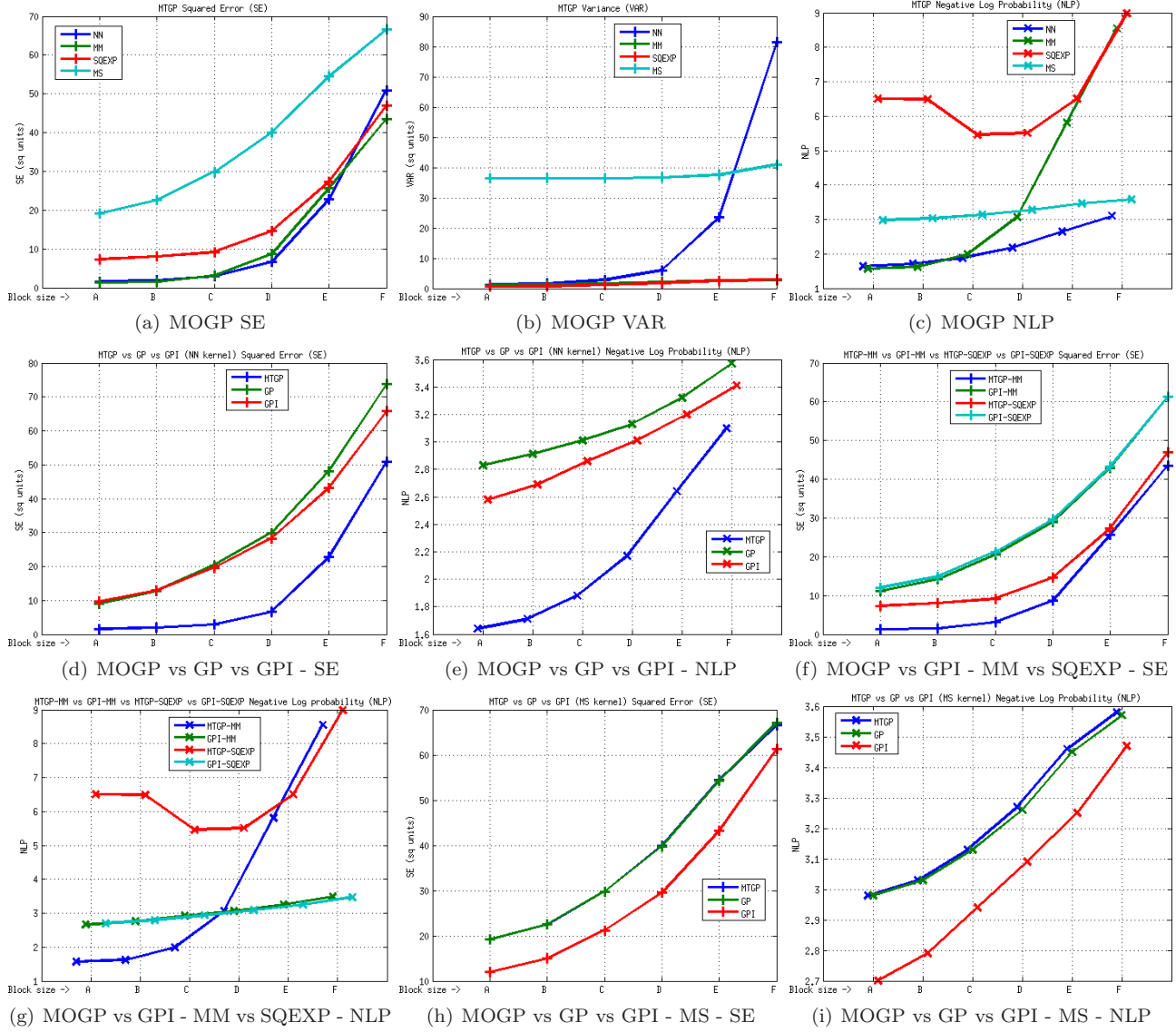


Figure 6: Element E3 - key trends. Test block sizes (m) - A (22 x 11 x 2), B (44 x 22 x 4), C (84 x 45 x 9), D (174 x 89 x 18), E (348 x 177 x 35) and F (696 x 353 x 70).