



Contents lists available at ScienceDirect

# Computational Statistics and Data Analysis

journal homepage: [www.elsevier.com/locate/csda](http://www.elsevier.com/locate/csda)



## Efficient optimization of the likelihood function in Gaussian process modelling<sup>☆</sup>



A. Butler<sup>a</sup>, R.D. Haynes<sup>a,\*</sup>, T.D. Humphries<sup>a</sup>, P. Ranjan<sup>b</sup>

<sup>a</sup> Department of Mathematics and Statistics, Memorial University, St. John's, NL, Canada, A1C 5S7

<sup>b</sup> Department of Mathematics and Statistics, Acadia University, Wolfville, NS, Canada, B4P 2R6

### ARTICLE INFO

#### Article history:

Received 1 August 2013

Received in revised form 22 November 2013

Accepted 24 November 2013

Available online 4 December 2013

#### Keywords:

BFGS

Clustering

Computer simulators

Dividing rectangles

Implicit filtering

Ill-conditioning

Nugget

### ABSTRACT

Gaussian Process (GP) models are popular statistical surrogates used for emulating computationally expensive computer simulators. The quality of a GP model fit can be assessed by a goodness of fit measure based on optimized likelihood. Finding the global maximum of the likelihood function for a GP model is typically challenging, as the likelihood surface often has multiple local optima, and an explicit expression for the gradient of the likelihood function may not be available. Previous methods for optimizing the likelihood function have proven to be robust and accurate, though relatively inefficient. Several likelihood optimization techniques are proposed, including two modified multi-start local search techniques, that are equally as reliable, and significantly more efficient than existing methods. A hybridization of the global search algorithm Dividing Rectangles (DIRECT) with the local optimization algorithm BFGS provides a comparable GP model quality for a fraction of the computational cost, and is the preferred optimization technique when computational resources are limited. Several test functions and an application motivated by oil reservoir development are used to test and compare the performance of the proposed methods with the implementation provided in the R library **GPMfit**. The proposed method is implemented in a Matlab package, **GPMfit**.

© 2013 Elsevier B.V. All rights reserved.

## 1. Introduction

Computer simulators are useful tools for modelling complex real world systems that are either impractical, expensive, or time consuming to physically observe. For example, the energy generated by the tides of large ocean basins (Greenberg, 1979), the estimation of the magnetic field generated near the Milky Way (Short et al., 2007), and the analysis of the flow of oil in a reservoir (Aziz and Settari, 1979) — the latter of which motivated this research — can be achieved through the use of computer simulators. That being said, realistic computer simulators can be computationally expensive to run, and as a result are often emulated using statistical models, such as Gaussian Process (GP) models (Sacks et al., 1989).

The maximum likelihood approach for fitting a GP model to deterministic simulator output requires minimizing the negative log-likelihood, or deviance. Rasmussen and Williams (2006) proposed the use of either a randomized multi-start conjugate gradient method or Newton's method for this problem. Explicit information about the gradient of deviance cannot be easily obtained, however, and the deviance function surface often has many local optima, making the optimization

<sup>☆</sup> Supplementary Material: The open source Matlab package **GPMfit** is available for download on SourceForge.net. See **Readme.txt** for detailed instruction. The main functions are **model1\_fit.m** and **predictor\_iterative.m**.

\* Corresponding author. Tel.: +1 7098648825.

E-mail addresses: [adsb85@mun.ca](mailto:adsb85@mun.ca) (A. Butler), [rhaynes@mun.ca](mailto:rhaynes@mun.ca), [rhaynes74@gmail.com](mailto:rhaynes74@gmail.com) (R.D. Haynes), [thumphries@mun.ca](mailto:thumphries@mun.ca) (T.D. Humphries), [pritam.ranjan@acadiau.ca](mailto:pritam.ranjan@acadiau.ca) (P. Ranjan).

problem challenging (MacDonald et al., 2013). Derivative-free optimization techniques, such as the genetic algorithm used by Ranjan et al. (2011), or the differential evolution algorithm used by Petelin et al. (2011), are robust, but can be computationally inefficient. Gradient approximation methods, such as the Broyden–Fletcher–Goldfarb–Shanno method (BFGS) (Broyden, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970), are generally faster, but have the potential to converge only locally if poorly initialized. MacDonald et al. (2013) proposed a clustering-based multi-start BFGS algorithm, which allows for a more global search to be performed. Nonetheless, this method requires multiple executions of BFGS, which is also computationally expensive.

In this paper we investigate several optimization techniques in order to improve the efficiency of the likelihood optimization process. Each technique is a combination of global and local search strategies. At the global level, we propose using the Dividing Rectangles algorithm (DIRECT) (Finkel, 2003) as an alternative to the clustering-based approach for choosing the starting point(s) of the local search. In terms of the local search, we compare the performance of BFGS with that of Implicit Filtering (IF), a sophisticated pattern search algorithm developed by Kelley (2011) for multimodal noisy functions. We use several test functions and an application motivated by real-world oil reservoir development to compare the performance of different optimization techniques, as measured by the prediction accuracy (optimized deviance and root mean squared prediction error) and number of deviance function evaluations (FEs) required to optimize the deviance. After an extensive case study we find that a hybrid approach of DIRECT and BFGS is the most efficient optimization technique for fitting such GP models.

The remainder of the paper is outlined as follows. Section 2 describes the GP model and the main components of the newly developed Matlab package **GPMfit**. In Section 3 we briefly outline the optimization techniques used for minimizing the deviance. Section 4 provides the results and analysis for several test functions, followed by an example in Section 5 where the GP model is fit to an oil reservoir simulator using our proposed method. Concluding remarks are provided in Section 6.

## 2. The Gaussian process model

The GP model requires as input a set of design points,  $x_i = (x_{i1}, \dots, x_{id})'$ , and the corresponding simulator outputs,  $y_i = y(x_i)$ , where  $i = 1, \dots, n$ , and  $n$  is the number of user supplied design points. Here, the prime symbol,  $'$ , denotes the transpose of vectors or matrices. We assume that the simulator provides a scalar valued output,  $y_i$ , for each  $d$ -dimensional design point  $x_i$ , and we use  $Y = (y_1, \dots, y_n)'$  to denote the  $n \times 1$  vector of simulator outputs. The simulator output is modelled as

$$y_i = \mu + z(x_i),$$

where  $\mu$  is the overall mean, and  $z(x_i)$  is a GP with  $E[z(x_i)] = 0$ ,  $\text{Var}[z(x_i)] = \sigma^2$ , and  $\text{Cov}[(z(x_i), z(x_j))] = \sigma^2 R_{ij}$ .

The  $n \times n$  spatial correlation matrix  $R$  defines the degree of dependency between design points, based on their observed simulator value. Following MacDonald et al. (2013), we use the Gaussian correlation matrix,  $R$ ; a special case of the power exponential correlation family defined as

$$R_{ij} = \prod_{k=1}^d \exp\{-10^{\beta_k} |x_{ik} - x_{jk}|^{p_k}\} \quad \text{for all } i, j. \quad (1)$$

Here  $p_k = 2$  is the smoothness parameter, and  $\beta = (\beta_1, \dots, \beta_d)$  is a  $1 \times d$  vector of correlation hyper-parameters which measures the sensitivity of the response to the spatial distribution of  $|x_{ik} - x_{jk}|^2$  for all  $i, j \in \{1, \dots, n\}$  and  $k \in \{1, \dots, d\}$  (Loepky et al., 2009).

The formulation of the correlation function in Eq. (1) is slightly different than the popular form of Gaussian correlation, which replaces  $10^{\beta_k}$  with  $\theta_k$  (e.g., in Ranjan et al., 2011). MacDonald et al. (2013) demonstrate that the deviance surface with  $\beta$ -parametrization shown in Eq. (1) is much easier to optimize as compared to the commonly used  $\theta$ -parametrization.

Sacks et al. (1989) show that the best linear unbiased predictor (BLUP) at a given point  $x^*$  in the input space (typically normalized to  $[0, 1]^d$ ) is

$$\begin{aligned} \hat{y}(x^*) &= \hat{\mu} + r'R^{-1}(Y - \mathbf{1}_n \hat{\mu}) \\ &= \left[ \frac{(1 - r'R^{-1}\mathbf{1}_n)}{\mathbf{1}_n'R^{-1}\mathbf{1}_n} \mathbf{1}_n' + r' \right] R^{-1}Y \\ &\equiv C'Y, \end{aligned}$$

where  $r = [r_1(x^*), \dots, r_n(x^*)]'$ , and  $r_i(x^*) = \text{corr}[z(x^*), z(x_i)]$  is the correlation between  $z(x^*)$  and  $z(x_i)$ . The GP model also returns the associated uncertainty estimate,  $s^2(x^*)$ , as measured by the mean squared error (MSE),

$$s^2(x^*) = E[(\hat{y}(x^*) - y(x^*))^2] = \hat{\sigma}^2(1 - 2C'r + C'RC). \quad (2)$$

The model fitting process requires the estimation of  $\mu$ ,  $\sigma^2$  and  $\beta$ . The closed form estimators of the mean and variance are given by

$$\hat{\mu}(\beta) = (\mathbf{1}_n'R^{-1}\mathbf{1}_n)^{-1}(\mathbf{1}_n'R^{-1}Y)$$

and

$$\hat{\sigma}^2(\beta) = \frac{(Y - \mathbf{1}_n \hat{\mu}(\beta))' R^{-1} (Y - \mathbf{1}_n \hat{\mu}(\beta))}{n},$$

respectively, and are used to obtain the profiled negative log-likelihood or deviance (ignoring the unimportant terms like  $\log(\sqrt{2\pi})$  and  $\log(n)$ ):

$$\mathcal{L}_\beta = \log(|R|) + n \log[(Y - \mathbf{1}_n \hat{\mu}(\beta))' R^{-1} (Y - \mathbf{1}_n \hat{\mu}(\beta))], \quad (3)$$

where  $\mathbf{1}_n$  is an  $n \times 1$  vector of all ones. The most difficult part of fitting the GP model is to find

$$\beta^* = \arg \min_{\beta \in \mathbb{R}^d} (\mathcal{L}_\beta).$$

Eq. (3) shows that evaluating  $\mathcal{L}_\beta$  requires computing both the action of the inverse,  $R^{-1}$ , and determinant,  $|R|$ , of the correlation matrix. If any pair of design points are sufficiently close to one another, the matrix  $R$  can become near-singular, resulting in unstable computation of  $R^{-1}$  and  $|R|$ . This can lead to an unreliable model fit.

To overcome this instability, we follow a popular technique developed by Sacks et al. (1989), Neal (1997), and Booker et al. (1999), which adds a small “nugget” parameter,  $\delta$ , to the model fitting procedure. The inclusion of  $\delta$  smoothes the model predictions, and consequently the GP model fit will no longer exactly interpolate the design points (O’Hagan and Kingman, 1978; Wahba, 1978). To avoid over-smoothing, Ranjan et al. (2011) introduce a lower bound on the nugget parameter,

$$\delta_{lb} = \max \left\{ \frac{\lambda_n(\kappa(R) - e^a)}{\kappa(R)(e^a - 1)}, 0 \right\},$$

where  $\kappa(R) = \lambda_n/\lambda_1$  is the 2-norm condition number and  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  are the eigenvalues of  $R$ . That is,  $R$  is simply replaced by  $R_\delta = R + \delta_{lb}I$  in Eq. (3). The over-smoothness can be reduced further by using the iterative regularization technique proposed by Ranjan et al. (2011). In our simulations, we scale the input-space domain to  $[0, 1]^d$ , and the design points are generated via a space-filling maximin Latin hypercube design (LHD). The desired threshold for  $\kappa(R)$  is  $e^a$ , where  $a = 25$  is recommended for a space-filling LHD (McKay et al., 1979). Under this configuration,  $\delta_{lb}$  values remain relatively small, if not zero, and therefore the iterative approach is not needed.

Although there are several choices for the correlation function, we focus on the Gaussian correlation, with  $p_k = 2$ , because of its smoothness property and popularity in other areas such as machine learning (radial basis kernels) and geostatistics (kriging). In practice, however, we can increase the stability of inverse and determinant computation by slightly lowering the smoothness parameter,  $p_k$ , of Eq. (1), such that  $p_k \lesssim 2$  (e.g.,  $p_k = 1.99$ ). By setting  $p_k = 1.99$ , the smoothness of the fitted surface does not appear to be affected significantly and the occurrence of near-singularity is substantially reduced; though not completely resolved as instability may still occur if the design points are extremely close to each other in input space. In this paper, we used  $p_k = 2$  with  $\delta_{lb}$  as given above for all simulated examples in Section 4, and  $p_k = 1.99$  with the same formula for  $\delta_{lb}$  in the oil reservoir application in Section 5.

### 3. Optimization methodology

Our objective function,  $\mathcal{L}_\beta$ , has a complicated dependency on  $\beta$ , namely in the form of the mean estimator,  $\hat{\mu}(\beta)$ , the correlation matrix,  $R$ , and the nugget parameter,  $\delta$ . With the inclusion of  $\delta$  it is difficult to extract an explicit gradient,  $\nabla \mathcal{L}_\beta$ , and therefore optimization methods that require the user to provide an expression for  $\nabla \mathcal{L}_\beta$  are not applicable here. We can, however, compute numerical approximations to  $\nabla \mathcal{L}_\beta$ , which is implicitly performed in both the BFGS and IF algorithms. That said,  $\mathcal{L}_\beta$  may contain several local optima and flat regions (see Fig. 1). Thus, a strictly descent-based optimization approach may not be desirable.

In the next section, we present a brief description of the local optimization algorithms used herein, namely BFGS and IF. In Sections 3.2 and 3.3 we describe two global optimization techniques: multi-start clustering and DIRECT. We conclude with a brief discussion on the bound constraints imposed on the  $\beta$  search space.

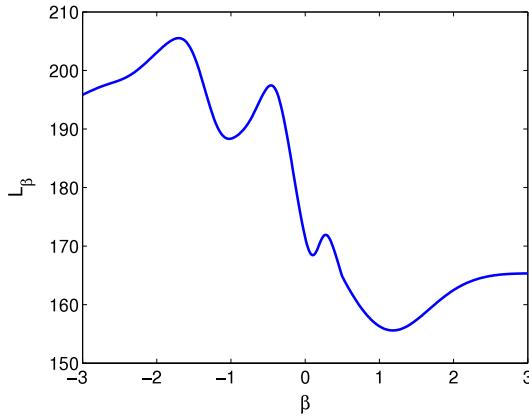
#### 3.1. BFGS and implicit filtering

The BFGS algorithm is a quasi-Newton optimization technique that uses a rank-two Hessian update formula to locate an optimal  $\beta$  (Coleman et al., 1999). At iteration  $k$ , the algorithm obtains a descent direction,  $q^{(k)}$ , by approximating the Hessian matrix,  $H^{(k)}$ , and the gradient,  $\nabla \mathcal{L}_\beta^{(k)}$ , at the current point,  $\beta^{(k)}$ , and solving

$$H^{(k)} q^{(k)} = -\nabla \mathcal{L}_\beta^{(k)} \cdot \beta^{(k)}.$$

A line search procedure then determines a suitable step-size,  $\alpha^{(k)}$ , along  $q^{(k)}$ , in order to obtain an updated solution,  $\beta^{(k+1)}$ , given by

$$\beta^{(k+1)} = \beta^{(k)} + \alpha^{(k)} q^{(k)}.$$



**Fig. 1.** 1-D plot of  $\mathcal{L}_\beta$  surface for the 1-D Hump function and a given set of design points. The Hump function is described in the Appendix.

We use Matlab's built-in unconstrained optimization routine `fminunc` for an implementation of BFGS. We have chosen to use the *medium-scale* implementation of `fminunc`, where the user must supply an initial value,  $\beta^{(0)}$ . Matlab is chosen for all experiments in this paper, due to the ready availability of Matlab implementations of the optimization approaches we consider.

Implicit Filtering (IF) is a sophisticated, deterministic pattern search algorithm designed by Kelley (2011) for bound constrained optimization. Specifically, IF is a local optimization algorithm that hybridizes a general pattern search algorithm with a BFGS derivative-approximation algorithm. The pattern search is arranged on a stencil, which, given an incumbent solution  $\beta^{(k)}$ , will evaluate  $\mathcal{L}_\beta$  at  $\beta^{(k)} \pm hv_j$ , in all coordinate directions  $j = \{1, \dots, d\}$ . Here,  $v_j = (U_j - L_j)e_j$ , where  $L_j$  and  $U_j$  represent the components of the lower and upper bounds of the search space, respectively, and  $e_j$  is the unit coordinate vector. The scale,  $h$ , varies as the optimization progresses, according to the sequence  $h = \{2^{-m}\}_{m=1}^7$ , where  $m$  is incremented each time the current stencil fails to find a more optimal position than the incumbent point. As the pattern search progresses, IF constructs a linear least squares interpolant from previously sampled points. After each pattern search phase, the linear interpolant surface is optimized locally using the BFGS algorithm. This process repeats until an optimal  $\beta$ -parameter is located. The idea is that the pattern search phase, with a suitable step-size, could step over local minima, while the quasi-Newton phase of IF will give efficient convergence in regions near the global optimum.

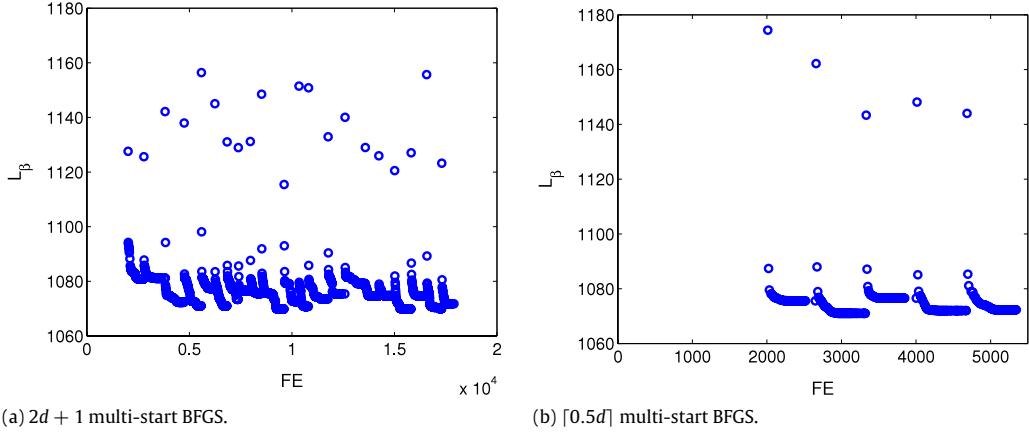
The BFGS algorithm does not require user specified bound constraints and works efficiently for smooth objective functions. With appropriate Hessian and gradient approximations and an efficient line search, the BFGS algorithm converges superlinearly to a locally optimal  $\beta$ -parameter near the initial solution  $\beta^{(0)}$  (Griva et al., 2009). Conversely, IF does require user specified bound constraints and is designed for functions that are noisy, with many local optima. The rate of convergence of IF is somewhat slower than superlinear (Kelley, 2011). However, IF is, in principle, more likely to find the global optimum within the bound constraints, due to its pattern search component. The quasi-Newton phase of both algorithms requires computation and storage of an approximate Hessian matrix and the solution of a linear system in order to obtain a suitable descent direction, which are relatively expensive operations.

### 3.2. Clustering-based multistart technique

MacDonald et al. (2013) proposed using a clustering-based multistart BFGS to replace the computationally expensive genetic algorithm (GA) used by Ranjan et al. (2011). The BFGS algorithm converges more rapidly than the GA, but lacks robustness, in the sense that it has the potential to get stuck in a local minimum, depending on the starting position,  $\beta^{(0)}$ . MacDonald et al. (2013) therefore proposed using  $2d + 1$  starting points of BFGS to improve the chances of global convergence, where  $d$  is the dimension of the simulator input (and therefore of  $\beta$  as well). These points are determined through sampling and a  $k$ -means clustering method (MacQueen, 1967), as described below.

1. Generate  $200d$   $\beta$ -vectors within the search space,  $S_\beta \subset (-\infty, \infty)^d$  (defined in Section 3.4), using a random maximin LHD, and evaluate  $\mathcal{L}_\beta$  for each  $\beta$ .
2. From the  $200d$  evaluations of  $\mathcal{L}_\beta$  obtained from Step 1, select the  $80d$   $\beta$ -vectors with the smallest  $\mathcal{L}_\beta$  values.
3. Cluster these  $80d$  points from Step 2 into  $2d$  groups, using the best of 5 random restarts of  $k$ -means clustering method. The  $2d$  cluster centres serve as the  $2d$  starting points of BFGS.
4. Evaluate  $\mathcal{L}_\beta$  at three equidistant points along the main diagonal of the search space,  $S_\beta$ . The  $\beta$ -vector with the lowest  $\mathcal{L}_\beta$  value is chosen as the  $(2d + 1)$ -th starting point.
5. Begin a run of BFGS from each of the  $2d + 1$  starting points.

From thorough experimentation on several test functions, we have observed that executing  $2d + 1$  starts of BFGS is excessive, and often results in several runs converging to comparable optima. This is evident from the results shown in



**Fig. 2.** Plot showing the convergence of each BFGS start versus the cumulative number of FEs (includes the initial  $200d = 2000$  FEs used for clustering), for fitting the 10-D Rastrigin function. The starting points are generated using a random maximin LHD.

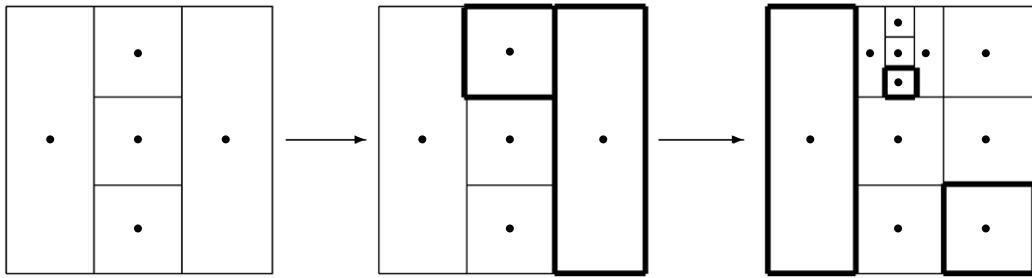
**Table 1**

Performance comparison of different likelihood optimization techniques on six test functions. The  $\% \Delta$  notation denotes the percent relative difference between the value of the performance measure returned by a given technique and the best value found among all techniques. Dashed values in the  $\% \Delta \mathcal{L}_\beta$  and  $\% \Delta \text{RMSPE}$  columns indicate that the best overall value was found by this algorithm. Underlined values in the FE column indicate the smallest number of FEs required by any algorithm.

Algorithm	Goldstein-Price (2-D)			Schwefel (5-D)		
	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE
$\lceil 0.5d \rceil$ BFGS	0.017	0.201	<u>439</u>	0.100	2.885	1859
$2d + 1$ BFGS	–	3.213	653	–	2.885	4277
$\lceil 0.5d \rceil$ IF	0.007	–	518	0.206	0.962	2381
$2d + 1$ IF	–	3.213	995	0.074	2.404	5735
$\lceil 0.5d \rceil$ IF-2	0.007	–	546	0.272	1.442	1677
DIRECT-BFGS	–	3.213	449	0.232	2.885	<u>1296</u>
DIRECT-IF	–	3.213	498	0.243	–	1304
Hartmann (6-D)			Rastrigin (10-D)			
	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE
	0.034	–	2068	0.071	3.865	5553
$\lceil 0.5d \rceil$ BFGS	0.102	1.064	5526	–	3.865	17853
$\lceil 0.5d \rceil$ IF	1.703	6.991	2293	0.624	–	4346
$2d + 1$ IF	–	0.304	7497	0.134	2.899	17284
$\lceil 0.5d \rceil$ IF-2	1.223	8.511	1866	0.545	0.483	4011
DIRECT-BFGS	0.207	1.216	<u>1526</u>	0.255	1.932	<u>2682</u>
DIRECT-IF	0.207	1.216	1533	0.287	1.932	2774
Rosenbrock (10-D)			Perm (12-D)			
	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE
	–	–	4562	0.003	1.370	8170
$\lceil 0.5d \rceil$ BFGS	–	–	16245	–	1.643	27622
$2d + 1$ BFGS	–	–	6654	0.010	–	12411
$\lceil 0.5d \rceil$ IF	–	–	24725	0.001	2.055	44375
$\lceil 0.5d \rceil$ IF-2	0.031	3.196	3775	0.025	–	4935
DIRECT-BFGS	0.003	0.457	<u>2332</u>	0.011	0.545	<u>3338</u>
DIRECT-IF	0.003	0.457	2654	0.012	0.685	3459

**Table 1** and is discussed in further detail in Section 4.3. As an example, Fig. 2 shows the convergence of  $2d + 1$  multistarts of BFGS (left plot) and  $\lceil 0.5d \rceil$  multistarts of BFGS (right plot) as a function of the number of deviance function evaluations (FE), for the problem of fitting a GP model to the 10-D Rastrigin test function (described in the Appendix). It is apparent that three of the five BFGS runs shown in the right plot converge to a value that is comparable to the best value found from the twenty-one BFGS runs shown in the left plot. We therefore propose reducing the number of cluster centres to  $\lceil 0.5d \rceil$  and eliminating the additional starting point obtained from sampling along the main diagonal. We implement two clustering-based techniques: (i)  $\lceil 0.5d \rceil$  multistarts of BFGS (i.e. the method of MacDonald et al. (2013) with fewer starting points), and (ii)  $\lceil 0.5d \rceil$  multistarts of IF.

The simulation results (in Table 1 and Figs. 2 and 4) show that the  $\mathcal{L}_\beta$  values returned by BFGS and IF do not usually change significantly after the first few iterations. Therefore, performing a full search (until the stopping criterion is satisfied) from each starting point  $\beta^{(0)}$  may be excessive as well. We therefore investigate a two stage multistart IF method, denoted by



**Fig. 3.** A few iterations of the Dividing Rectangles (DIRECT) algorithm. Bold rectangles are identified as being potentially optimal and are divided in the following iteration. The •'s denote the sampled points.

IF-2. The process starts with a clustering-based,  $\lceil 0.5d \rceil$  multistart IF approach, where each run of IF is limited to a budget of  $20d$  FEs. In the second stage, the single run of IF that returns the lowest  $\mathcal{L}_\beta$  value will then run to completion.

### 3.3. DIRECT hybrid techniques

Dividing Rectangles (DIRECT) is a derivative-free, block partitioning algorithm that sequentially samples points in the search space and partitions the domain into hyper-rectangles based on the objective function value (here,  $\mathcal{L}_\beta$ ) at the sampled points (Finkel, 2003). Hyper-rectangles are then identified as being potentially optimal if they contain a sampled point whose function value is more optimal than the sampled points contained by all other hyper-rectangles of equal size. Each potentially optimal hyper-rectangle is then divided into thirds along its longest dimension, and the process repeats. Fig. 3 provides a 2-dimensional visualization of how DIRECT samples and divides the search space. The left panel in Fig. 3 shows the initial sampling phase and partitioning of the domain. The bold rectangles in the middle panel of Fig. 3 are identified as being potentially optimal in the following iteration. The rectangles are then partitioned in turn (rightmost figure), and three new rectangles are identified as potentially optimal. The alternating process of partitioning and then identifying potentially optimal rectangles continues, until a stopping criterion is met.

DIRECT is specifically designed as a global optimization approach, and since it provides a thorough exploration of search space, it can be slow to converge locally. We therefore use DIRECT to provide a single, somewhat optimized starting point,  $\beta^{(0)}$ , from which to begin a run of either BFGS or IF, thereby eliminating the need for a multi-start approach. For ease of comparison with the clustering-based approach, we provide the same budget of  $200d$  FEs to DIRECT (as in Step 1 of the clustering based approach).

### 3.4. Boundaries of the search space

Many optimization techniques, including IF and DIRECT, require user-supplied bound constraints on the optimization parameters,  $\beta$ . Although the exact position of the global optimum in  $\beta$ -space is unknown, we can determine a region,  $S_\beta$ , where the optimum is likely to be found. MacDonald et al. (2013) use the structural form of the correlation matrix  $R$  to provide an approximate bound for  $R_{ij}$ . Specifically

$$\exp(-5) \approx 0.0067 \leq R_{ij} \leq 0.9999 \approx \exp(-10^{-4}),$$

or equivalently,

$$10^{-4} \leq \sum_{k=1}^d 10^{\beta_k} |x_{ik} - x_{jk}|^2 \leq 5.$$

Since the assumed input-space is  $[0, 1]^d$ , and the design points are generated via a maximin LHD, the approximate spatial distribution of the  $10d$  design points that we use is at most  $|x_{ik} - x_{jk}| \approx 1/10$  in any dimension  $k$ . Moreover, if we assume that the simulator function is equally smooth in all coordinate directions, i.e.,  $\beta_1 \approx \beta_2 \approx \dots \approx \beta_d$ , then the set of  $\beta$  values that is likely to contain the global optimum is

$$S_\beta = \{(\beta_1, \dots, \beta_d) : -2 - \log_{10}(d) \leq \beta_k \leq \log_{10}(500) - \log_{10}(d), k = 1, \dots, d\}.$$

The starting points,  $\beta^{(0)}$ , determined using either clustering or DIRECT, will be confined to the region  $S_\beta$ .

BFGS is an unbounded optimization algorithm and does not require any user-specified bound constraints on the optimization variables, whereas IF does. The step-size calculated in the pattern search phase of IF is proportional to the physical size of the user supplied bound constraints, which we denote by  $S_\beta^{IF}$  (defined below). From our experimentation on the test functions, we have noticed that if the size of  $S_\beta^{IF}$  is too small, then the efficiency of the optimization is compromised because IF requires a large number of FEs to converge to an optimal  $\beta$ -parametrization. Conversely, if  $S_\beta^{IF}$  is too large, IF has the tendency to “jump” around the potentially optimal regions of  $\mathcal{L}_\beta$ , resulting in convergence to a suboptimal  $\beta$  value.

MacDonald et al. (2013) note that the optimal  $\beta$  values are rarely large and positive, and hence, we modify the bound constraints on IF in which the negative  $\beta$  region occupies a larger portion of the domain, i.e.,

$$S_{\beta}^{IF} = \{(\beta_1, \dots, \beta_d) : d(-2 - \log_{10}(d)) \leq \beta_k \leq \log_{10}(500), k = 1, \dots, d\}.$$

We acknowledge that the  $\beta$ -value that globally minimizes  $\mathcal{L}_{\beta}$  may occasionally be positioned outside the provided bounds,  $S_{\beta}$  and  $S_{\beta}^{IF}$ . Therefore, included in the **GPMfit** package is the option to multiplicatively expand or contract the default bound constraints.

#### 4. Simulation results

We use seven test functions, with input dimensions varying from  $d = 1$  to  $d = 12$  to compare the performance of the different optimization techniques discussed in Section 3. The formulae for all the test functions are provided in the [Appendix](#). The performance of each optimization technique is averaged over 25 simulations. For each simulation,  $10d$  training design points ( $x_i$ ) and  $100d$  validation prediction points ( $x_i^*$ ) are chosen in  $[0, 1]^d$  via a space-filling maximin LHD. The initial sample points in  $\beta$ -space for the clustering procedure are randomly generated in  $S_{\beta}$  using the LHD. All simulations were performed using 64-bit Matlab 2012(b) on a Gentoo Linux operating system with a Core 2 Quad Xeon processor.

##### 4.1. Optimization accuracy

Recall that our objective is to minimize  $\mathcal{L}_{\beta}$ . Typically, the parameter estimate that corresponds to the smallest  $\mathcal{L}_{\beta}$  will provide the most accurate model fit, as measured by the average relative root mean square prediction error (RMSPE) between the GP model fit and the true simulator (test function) response. That is,

$$\text{RMSPE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2} / \sqrt{\frac{1}{N} \sum_{i=1}^N y_i^2}, \quad \text{where } N = 100d. \quad (4)$$

We note that in most real applications, the true RMSPE values cannot be calculated, as the true simulator outputs are unknown at the validation points. One can, however, use the average RMSE estimates (see Eq. (2)) for performance comparison. The consistency of RMSPE values over 25 simulations is measured by the standard error in the RMSPE value, given by

$$\text{Std. Err.} = \sigma_{\text{RMSPE}} / \sqrt{25}, \quad (5)$$

where  $\sigma_{\text{RMSPE}}$  denotes the standard deviation of the RMSPE values. A standard error of one order of magnitude less than the corresponding RMSPE value indicates that our results are fairly consistent over all 25 simulations.

##### 4.2. Convergence efficiency

We measure the efficiency of an optimization technique by the number of likelihood function evaluations (FEs) required for optimization of  $\mathcal{L}_{\beta}$ . Using Matlab's profiler, we determined that evaluating  $\mathcal{L}_{\beta}$  constituted the bulk of the computational load for all optimization techniques considered. In particular, we determined that computation of the correlation matrix,  $R_{\delta}$ , demands anywhere from 60%–90% of the total computation time, depending on the simulator input dimension,  $d$ . Computation of  $R_{\delta}$  is nested within the calculation of  $\mathcal{L}_{\beta}$  and is evaluated once per FE. Therefore, the number of FEs, which is not affected by issues like server load, can be used in place of computation time as a fair measure of optimization efficiency.

##### 4.3. Discussion

**Table 1** summarizes the accuracy ( $\mathcal{L}_{\beta}$  and RMSPE) and efficiency (FE) of each optimization technique for six of the seven test functions, namely the 2-D Goldstein–Price function, the 5-D Schwefel function, the 6-D Hartmann function, the 10-D Rastrigin function, the 10-D Rosenbrock function and the 12-D Perm function (see [Appendix](#) for closed form expressions). Results for fitting the 1-D Hump function are not shown here, as the performance of all the techniques was essentially the same for this simple test function. The  $\% \Delta$  notation in **Table 1** denotes the percent relative difference between the value of the performance measure returned by a given technique and the best value found among all techniques. The standard errors are not included in this table; we found that for all cases the standard error was indeed at least one order of magnitude less than the corresponding RMSPE, suggesting that each technique is consistently able to provide the same GP model quality.

The results in **Table 1** show that the  $[0.5d]$  multi-start techniques and DIRECT-based techniques provide efficient and reliable alternatives to the  $2d + 1$  BFGS technique. We observe that the  $[0.5d]$  multi-start and DIRECT-based techniques require anywhere from 20% to 90% fewer FEs than the  $2d + 1$  BFGS technique (depending on the dimension of the test function), while maintaining a comparable level of optimization accuracy. The relative difference in the  $\mathcal{L}_{\beta}$  value returned by each technique is typically less than 1%. As a result, each optimization technique provides comparable GP model quality,

**Table 2**

Ranks of the optimization techniques based on their overall average performance for the seven test functions.

Rank	Algorithm	# of starts
1	DIRECT-BFGS	1
2	DIRECT-IF	1
3	BFGS	$\lceil 0.5d \rceil$
4	IF	$\lceil 0.5d \rceil$
5	BFGS	$2d + 1$
6	IF	$2d + 1$
7	IF-2	$\lceil 0.5d \rceil$

as measured by the RMSPE value. Table 1 shows, however, that for all test functions, the  $\lceil 0.5d \rceil$  multi-start IF-2 technique returns a larger average  $\mathcal{L}_\beta$  value and requires anywhere from 10% to 50% more FEs than both of the DIRECT-based methods. As a result, IF-2 is not included in the **GPMfit** package, as more accurate and efficient techniques are clearly available.

The results in Table 1 also indicate that a slightly suboptimal  $\mathcal{L}_\beta$  value can result in an equal, or even slightly better GP model quality, as measured by the RMSPE. For example, for fitting the 10-D Rastrigin function, DIRECT-BFGS provides an RMSPE value that is 1.93% smaller than the RMSPE value returned by  $2d + 1$  BFGS, despite converging to a  $\mathcal{L}_\beta$  value that is sub-optimal by 0.255%. The non-monotonic relationship between optimal  $\mathcal{L}_\beta$  and GP model quality is presented in a recent paper by Kalaitzis and Lawrence (2011), who argue that one can maintain the quality of the GP model even with a slightly suboptimal  $\mathcal{L}_\beta$  value. Furthermore, Nguyen et al. (2011) suggest that, due to the difficulties in finding optimal  $\beta$ -parameters, particularly when the training data  $(x_i, y_i)$  is sparse, GP models can be prone to overfitting, which can lead to larger than expected RMSPE values. The authors acknowledge that the degree to which one optimizes  $\mathcal{L}_\beta$ , in practice, is often situation specific, as GP model prediction accuracy can be largely affected by overfitting. Thus, motivated by this non-monotonic relationship, our goal is to efficiently determine a sufficiently optimal  $\mathcal{L}_\beta$  value, without compromising the resulting GP model quality.

Fig. 4 shows the convergence performance of both the BFGS and IF algorithms after the initial  $\beta^{(0)}$  point(s) are determined using either clustering or DIRECT. For the multi-start clustering-based techniques, the convergence plots are displayed as though each run of BFGS or IF was implemented in parallel, from which the best of the  $\lceil 0.5d \rceil \mathcal{L}_\beta$  values is plotted versus the cumulative number of FEs. The optimization performance for each technique has been averaged over all 25 simulations, and is plotted on a semi-log scale as the absolute difference between  $\mathcal{L}_\beta$  and the minimum  $\mathcal{L}_\beta$  value,  $\mathcal{L}_{\min}$ , (rounded down to the nearest  $10^{-2}$ ) determined by one of the four techniques.

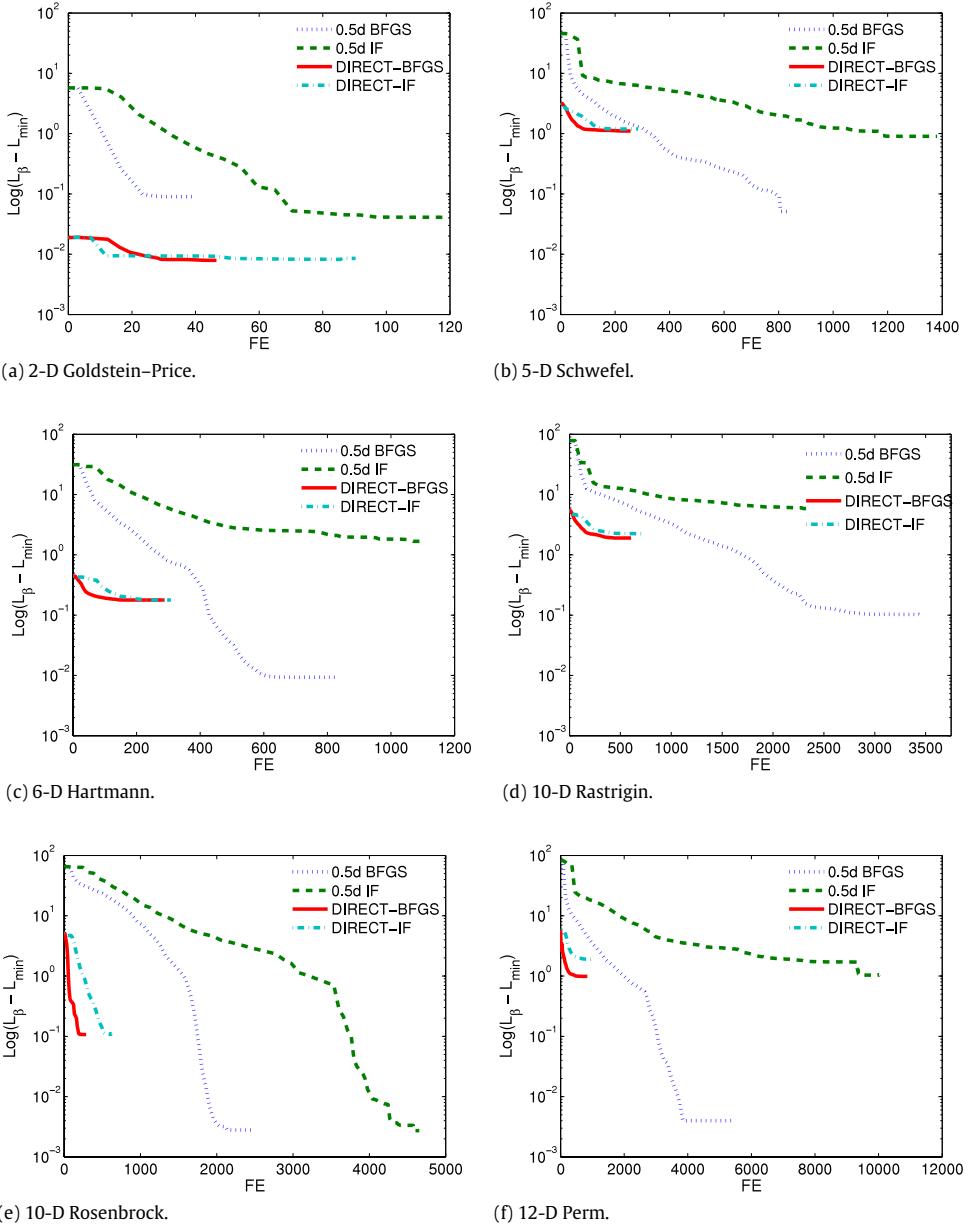
We first observe that, for both the DIRECT and multi-start approaches, BFGS is more efficient than IF, and converges to a more optimal value of  $\mathcal{L}_\beta$ . This suggests that BFGS is generally better suited to this optimization problem than IF. Secondly, with the exception of the 2-D test case (Fig. 4(a)), the  $\lceil 0.5d \rceil$  BFGS technique converges to the best  $\mathcal{L}_\beta$  value. In general, however, the difference between the  $\mathcal{L}_\beta$  value returned by the  $\lceil 0.5d \rceil$  BFGS technique and DIRECT-based methods is on the order of  $10^{-1}$  to  $10^0$ , which has only a small effect on the resulting quality of the GP model, as measured by the RMSPE. Moreover, Fig. 4 shows that once the starting  $\beta^{(0)}$  point(s) have been determined, the DIRECT-based methods require approximately  $\frac{1}{\lceil 0.5d \rceil}$  as many FEs as the multi-start clustering techniques. This represents a substantial increase in optimization efficiency, particularly as  $d$  increases. Thus, the DIRECT-based approaches are able to find an optimum that is only slightly worse than those found by the significantly more expensive clustering-based approaches.

We ran an additional experiment to determine whether providing additional FEs to the DIRECT-based methods would result in these methods converging to the optimal  $\mathcal{L}_\beta$  value found by  $\lceil 0.5d \rceil$  BFGS. Fig. 5 compares the optimization performance of the DIRECT-based and  $\lceil 0.5d \rceil$  multi-start clustering techniques for a single GP model realization of the 5-D Schwefel function. For this particular case, the local search techniques used after DIRECT have been allotted 900 FEs, which is the number of FEs that were required for  $\lceil 0.5d \rceil$  BFGS to converge to the optimal  $\mathcal{L}_\beta$ . From the logarithmic plot (Fig. 5(a)), we can see that there is no improvement in the solution found by the DIRECT-BFGS and DIRECT-IF methods after roughly 200 FEs, indicating that allotment of additional FEs provides no benefit to these approaches. We note again, however, that when plotted on regular axes, as in Fig. 5(b), the discrepancy in the  $\mathcal{L}_\beta$  values returned by the various methods is small.

Figs. 4 and 5 also show that the DIRECT algorithm is able to determine a more optimal starting position,  $\beta^{(0)}$ , for initialization of BFGS or IF. This enables us to implement a single run of BFGS or IF, with only a small loss of optimization accuracy. These results were observed for all test functions and support a fundamental conclusion; if the user has significant computational resources at their disposal and if a highly accurate optimal  $\mathcal{L}_\beta$  is desired, then the multi-start  $\lceil 0.5d \rceil$  BFGS technique is preferred. If computational resources are limited, however, then one can use DIRECT-BFGS to obtain comparable GP model fit for a fraction of the computational cost. From these results, we are able to establish an overall performance ranking of each optimization technique, shown in Table 2.

## 5. Oil reservoir simulator example

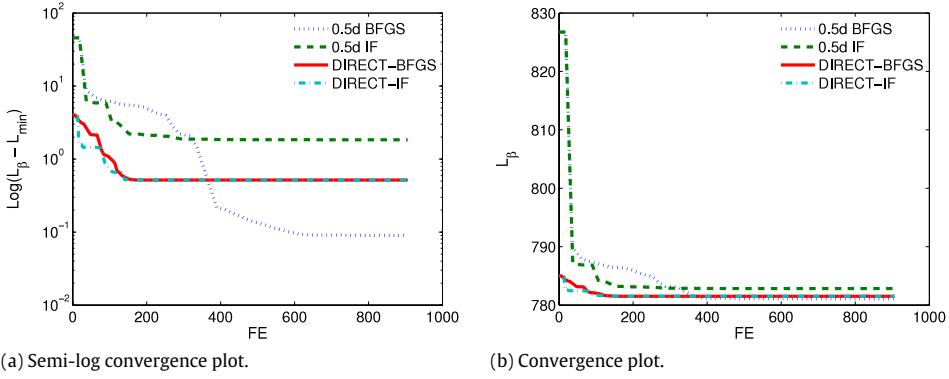
Determining optimal drilling locations for production and injection wells in an oil reservoir is a problem of considerable industrial interest (see for instance Yeten et al., 2003; Bangerth et al., 2006; Onwunalu and Durlofsky, 2010). The variables



**Fig. 4.** Semi-log plots comparing the performance of the single start DIRECT-BFGS and DIRECT-IF techniques and the multi-start [0.5d] BFGS and IF techniques, averaged over 25 simulations.

in this problem correspond to positional parameters for each well; in this example, we will consider only vertical wells, each of which can be parameterized by its  $(x_1, x_2)$  co-ordinates, representing a grid location in the discrete reservoir model. The well locations serve as input to a computationally expensive complex reservoir simulator – in our case, the Matlab Reservoir Simulator (MRST) (Lie et al., 2011; SINTEF Applied Mathematics, 2012). The simulator output, along with various economic parameters, are then usually combined to provide the net present value (NPV) of the produced oil. The goal is to determine the configuration of wells that yields the best NPV.

We consider two problems using a simple 2-D reservoir model based on a  $60 \times 50$  grid. For the first placement problem, we assume that two injection wells ( $\times$ ) and one production well ( $\circ$ ) have already been drilled at the positions shown in Fig. 7(a), and the goal is to find the optimal location for the second production well. The NPV surface corresponding to this problem is shown in Fig. 7(a). One could use an expected improvement based sequential design scheme (Jones et al., 1998) for finding this optimal location; the key component in such a sequential optimization is to efficiently emulate (i.e., fit a GP model to) the simulator response after every iteration of this sequential procedure. In this paper, we focus on this first step of fitting a GP model-based surrogate to the simulator output. For the second problem, we allow the positions of all four



**Fig. 5.** Semi-log convergence plot and convergence plot comparing the  $\mathcal{L}_\beta$  optimization performance of the DIRECT-based techniques and the multi-start clustering techniques.

**Table 3**

Performance comparison of different  $\mathcal{L}_\beta$  optimization methods for the 2-D reservoir simulator. Dashed values in the  $\% \Delta \mathcal{L}_\beta$  and  $\% \Delta \text{RMSPE}$  columns indicate that the best overall value was found by this algorithm. Underlined values in the FE column indicate the smallest number of FEs required by any algorithm.

Algorithm	$n = 20$			$n = 40$		
	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE
[0.5d] BFGS	0.110	2.461	<u>435</u>	2.153	41.602	452
2d + 1 BFGS	–	2.461	691	–	3.359	721
DIRECT-BFGS	0.029	–	<u>435</u>	0.004	–	<u>440</u>
$n = 80$				$n = 100$		
Algorithm	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE
	2.757	47.734	461	2.648	22.459	465
[0.5d] BFGS	–	–	756	–	3.476	822
DIRECT-BFGS	–	–	<u>443</u>	0.072	–	<u>459</u>

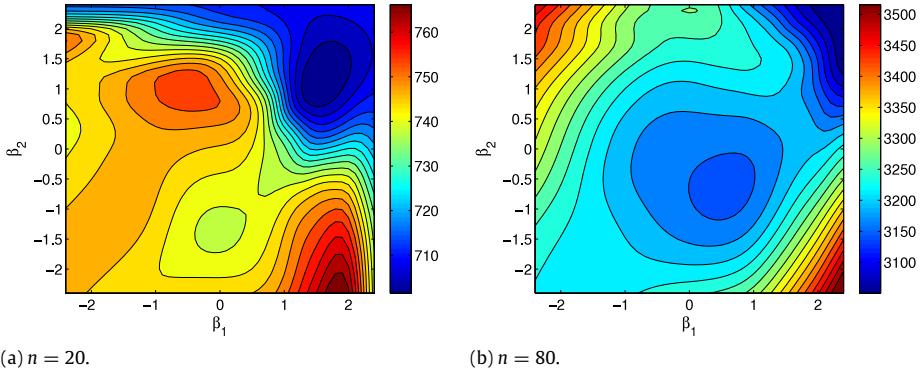
wells to be chosen freely, meaning that the NPV now depends on 8 variables. Both of these problems are variants of cases considered in Humphries et al. (in press).

### 5.1. 2-D reservoir simulator

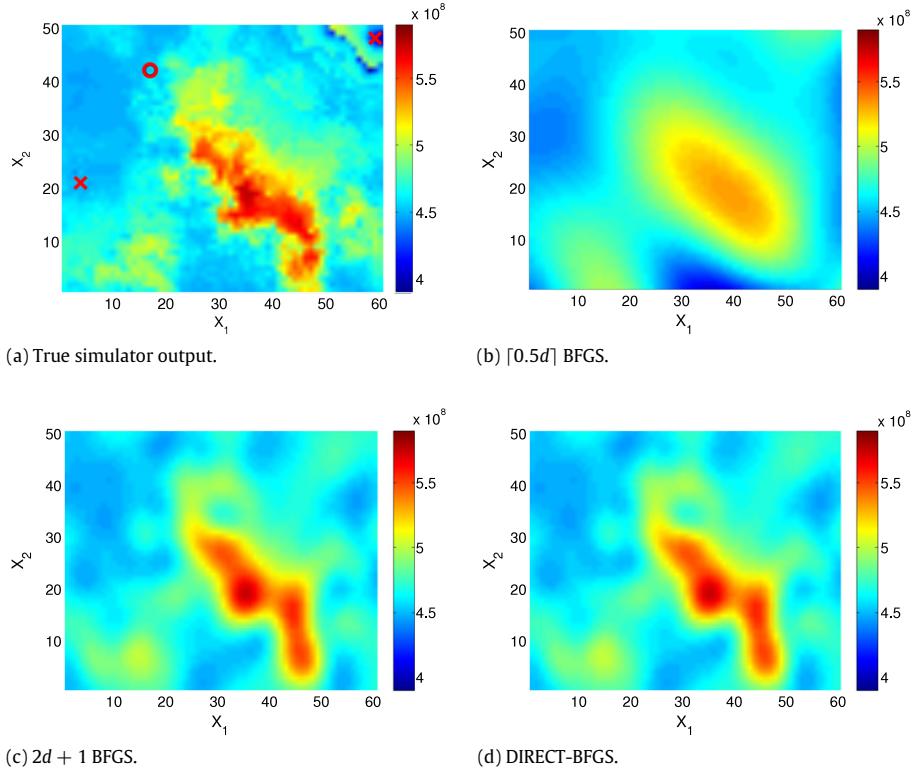
**Table 3** compares the performance of three methods for fitting the GP model to the 2-D reservoir simulator: (i) [0.5d] BFGS, (ii) 2d + 1 multi-start BFGS, and (iii) DIRECT-BFGS. The number of training design points used,  $n$ , ranges from 20 to 100, and are generated via a space-filling maximin LHD. The performance of each optimization technique is averaged over 25 simulations. For every value of  $n$ , the 2d + 1 BFGS technique returns the best  $\mathcal{L}_\beta$  value. The relative percent difference of the  $\mathcal{L}_\beta$  value returned by the DIRECT-BFGS technique, however, is always less than 0.1%. Moreover, in all cases, DIRECT-BFGS returns the smallest RMSPE value and requires anywhere from 37% to 45% fewer FEs than 2d + 1 BFGS, and is therefore the preferred technique.

A close look at one realization (see Fig. 6) reveals that the likelihood surface changes significantly as the number of design points increases. In particular, when  $n = 80$ , a local optimum with a large basin of attraction is created near the centre of the  $S^\beta$  domain. As a result, the [0.5d] multi-start BFGS technique, which is a single-start technique when  $d = 2$ , converges to this sub-optimal point when initialized in that region. DIRECT, on the other hand, is able to determine a starting point that is significantly closer to the global optimum (located in the top right corner) and thus avoids converging to the sub-optimal point. This explains why the  $\mathcal{L}_\beta$  and corresponding RMPSE values determined by [0.5d] BFGS are significantly higher than those values determined by 2d + 1 BFGS and DIRECT-BFGS, when  $n \geq 40$ .

Fig. 7 shows the true simulator output and an example of the GP models that were found using each of the three optimization techniques, with  $n = 100$  design points. The 2d + 1 BFGS and DIRECT-BFGS approaches provide near identical GP model predictions, with DIRECT-BFGS requiring a fraction of the computational cost. As mentioned, the [0.5d] BFGS technique converges to a sub-optimal  $\beta$ -parameter, and as a result the GP model quality is significantly worse. In general we observe that for small dimensional simulators, DIRECT-BFGS outperforms the [0.5d] multi-start BFGS technique in terms of both optimization accuracy and efficiency (**Tables 1** and **3**, **Fig. 7**). The 2d + 1 multi-start BFGS method often provides a more accurate GP model fit for any number of design points, but requires up to 80% more FEs than both DIRECT-BFGS and [0.5d] BFGS for a 2-D function. This example shows that DIRECT-BFGS provides an efficient alternative to the 2d + 1 multi-start BFGS approach, without sacrificing model accuracy.



**Fig. 6.** Comparing  $\mathcal{L}_\beta$  surfaces of the 2-D reservoir simulator for varying design points.



**Fig. 7.** True 2-D reservoir simulator output and GP fit surface obtained through  $\mathcal{L}_\beta$  optimization using the  $[0.5d]$  BFGS,  $2d + 1$  BFGS and DIRECT-BFGS techniques, with  $n = 100$ . The locations of the two existing injection ( $\times$ ) wells and single production ( $\circ$ ) well are shown in plot (a).

## 5.2. 8-D reservoir simulator

In our second example, we fit a GP model to the oil reservoir simulator with 8 variables (the positions of all four wells), again using  $[0.5d]$  BFGS,  $2d + 1$  BFGS and DIRECT-BFGS for  $\mathcal{L}_\beta$  optimization. The GP model is initialized using both  $10d$  and  $20d$  design points, and predicts for  $100d$  points of unknown function value. The performance of each of the three optimization techniques is averaged over 25 simulations.

Table 4 compares the  $\mathcal{L}_\beta$  optimization and resulting GP model fitting performance for each of the three techniques. Again, we observe that the  $2d + 1$  BFGS technique converges to the best  $\mathcal{L}_\beta$  value on average. Nonetheless, this method requires roughly 12000 FEs, which is more than 3 times the number of FEs required by  $[0.5d]$  BFGS and almost 6 times the number of FEs required by DIRECT-BFGS. Moreover, DIRECT-BFGS, on average, provides the highest quality GP model, as measured by the RMSPE value, despite converging to a slightly sub-optimal  $\mathcal{L}_\beta$  value. The results obtained in fitting the GP model to a true reservoir simulator provide evidence that, in practice, one can employ the single start DIRECT-BFGS optimization technique to greatly increase the efficiency of the GP model fitting procedure, without compromising the quality of the model.

**Table 4**

Performance comparison of different  $\mathcal{L}_\beta$  optimization methods for the 8-D reservoir simulator. Dashed values in the  $\% \Delta \mathcal{L}_\beta$  and  $\% \Delta \text{RMSPE}$  columns indicate that the best overall value was found by this algorithm. Underlined values in the FE column indicate the smallest number of FEs required by any algorithm.

Algorithm	$n = 80$			$n = 160$		
	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE	$\% \Delta \mathcal{L}_\beta$	$\% \Delta \text{RMSPE}$	FE
$\lceil 0.5d \rceil$ BFGS	0.013	0.711	3706	0.007	1.139	3567
$2d + 1$ BFGS	—	2.370	11896	—	1.635	11893
DIRECT-BFGS	0.032	—	<u>2070</u>	0.026	—	<u>2005</u>

**Table A.5**

Test functions and corresponding formula used for evaluating the performance of the  $\mathcal{L}_\beta$  optimization process in GP modelling.

Test function (d)	Formula $y = f(x)$
Hump (1)	$y = 1.0316285 + 4x^2 - 2.1x^4 + \frac{1}{3}x^6$
Goldstein-Price (2)	$y = [1 + (x_1 + x_2 + 1)^2(19 - 14x_1 + 3x_1^2 - 14x_2 + 6x_1x_2 + 3x_2^2)] \times [30 + (2x_1 - 3x_2)^2(18 - 32x_1 + 12x_1^2 + 48x_2 - 36x_1x_2 + 27x_2^2)]$
Schwefel (5)	$y = 2094.9 - \sum_{i=1}^5 x_i \sin(\sqrt{ x_i })$
Hartmann (6)	$y = -\sum_{i=1}^6 \alpha \cdot \exp[-\sum_{j=1}^6 B_{ij}(x_j - Q_{ij})^2]$ $\alpha = [1, 1.2, 3, 3.2], B = \begin{bmatrix} 10 & 3 & 17 & 3.5 & 1.7 & 8 \\ 0.02 & 10 & 17 & 0.1 & 8. & 14 \\ 3 & 3.5 & 1.7 & 10 & 17 & 8 \\ 17 & 8 & 0.05 & 10 & 0.1 & 14 \end{bmatrix},$ $Q = \begin{bmatrix} 0.1312 & 0.1696 & 0.5569 & 0.0124 & 0.8283 & 0.588 \\ 0.2329 & 0.4135 & 0.8307 & 0.3736 & 0.1004 & 0.9991 \\ 0.2348 & 0.1451 & 0.3522 & 0.2883 & 0.3047 & 0.6650 \\ 0.4047 & 0.8828 & 0.8732 & 0.5743 & 0.1091 & 0.0381 \end{bmatrix}$
Rastrigin (10)	$y = 10n + \sum_{i=1}^{10} (x_i^2 - 10 \cos(2\pi x_i))$
Rosenbrock (10)	$y = \sum_{i=1}^9 [100(x_i^2 - x_{i+1})^2 + (x_i - 1)^2]$
Perm (12)	$y = \sum_{i=1}^{12} \left[ \sum_{j=1}^{12} \left[ (j^i + 0.5) \left( \frac{x_j}{j} \right)^{i-1} \right]^2 \right]$

## 6. Conclusion

In this paper we have investigated several techniques for efficient optimization of the deviance function,  $\mathcal{L}_\beta$ , in GP modelling. These techniques provide the foundation for an improved Matlab package, **GPMfit**. The results obtained from simulated examples and real applications show that using  $2d + 1$  multi-starts of BFGS is computationally expensive, and that we are able to significantly improve the optimization efficiency by reducing the number of multi-starts to  $\lceil 0.5d \rceil$ , while maintaining the quality of the GP model in most cases.

Implicit filtering proves to be slightly less accurate and efficient than BFGS, and therefore is included in the **GPMfit** package solely as a secondary option. The modified multi-start technique IF-2 is generally unreliable; specifically, the slight reduction in computational cost that is gained does not outweigh the reduced accuracy, particularly when more efficient and robust optimization techniques exist.

Replacing the  $\lceil 0.5d \rceil$  multi-start technique with the DIRECT optimization algorithm enables us to further reduce the number of starts of BFGS or IF from  $\lceil 0.5d \rceil$  to 1. After an initial  $\beta_0$  value has been determined, the DIRECT-based hybrid techniques require approximately  $\frac{1}{\lceil 0.5d \rceil}$  as many function evaluations as the multi-start techniques and provide an almost equally accurate model fit. As a result, the DIRECT-BFGS hybrid technique is the default  $\mathcal{L}_\beta$  optimization algorithm used in the **GPMfit** package. The slightly less efficient DIRECT-IF hybrid technique is also included in the **GPMfit** package as an additional optimization option, along with the more computationally expensive  $2d + 1$  and  $\lceil 0.5d \rceil$  multi-start BFGS and IF techniques.

## Acknowledgements

The authors would like to thank the editor, AE and two referees for their valuable feedback. Haynes, Ranjan and Butler would like to acknowledge the support of the Natural Sciences and Engineering Research Council of Canada's Discovery Grant and USRA programs. Humphries was supported by the Research Development Corporation of Newfoundland and the Atlantic Canada Opportunities Agency.

## Appendix. Test functions

See Table A.5.

## References

- Aziz, K., Settari, A., 1979. *Petroleum Reservoir Simulation*. Applied Science Publishers, London.
- Bangerth, W., Klie, H., Wheeler, M., Stoffa, P., Sen, M., 2006. On optimization algorithms for the reservoir oil well placement problem. *Comput. Geosci.* 10, 303–319.
- Booker, A.J., Dennis Jr., J.E., Frank, P.D., Serafini, D.B., Torczon, V., Trosset, M.W., 1999. A rigorous framework for optimization of expensive functions by surrogates. *Struct. Optim.* 17 (17), 1–13.
- Broyden, C., 1970. The convergence of a class of double-rank minimization algorithms. *J. Appl. Math.* 6 (1), 76–90.
- Coleman, T., Branch, M.A., Grace, A., 1999. Optimization Toolbox: For Use with Matlab. The Math Works Inc.
- Finkel, D.E., 2003. DIRECT Optimization Algorithm User Guide. Center for Research in Scientific Computation.
- Fletcher, R., 1970. A new approach to variable metric algorithms. *Comput. J.* 13 (3), 317–322.
- Goldfarb, D., 1970. A family of variable metric updates derived by variational means. *Math. Comp.* 24 (1), 23–26.
- Greenberg, D., 1979. A numerical model investigation of tidal phenomena in the bay of fundy and gulf of maine. *Mar. Geod.* 2, 161–187.
- Griva, I., Nash, S., Sofer, A., 2009. *Linear and Nonlinear Optimization*. Society for Industrial and Applied Mathematics, pp. 355–448.
- Humphries, T., Haynes, R., James, L., 2013. Simultaneous and sequential approaches to joint optimization of well placement and control. *Comput. Geosci.* <http://dx.doi.org/10.1007/s10596-013-9375-x>. (in press).
- Jones, D.R., Schonlau, M., Welch, W., 1998. Efficient global optimization of expensive black-box functions. *J. Global Optim.* 13 (4), 445–492.
- Kalaitzis, A.A., Lawrence, N.D., 2011. A simple approach to ranking differentially expressed gene expression time courses through Gaussian process regression. *BMC Bioinformatics* 12 (1).
- Kelley, C., 2011. *Implicit Filtering*. Society for Industrial and Applied Mathematics.
- Lie, K.-A., Krogstad, S., Ligaarden, I., Natvig, J., Nilsen, H., Skaflestad, B., 2011. Open-source MATLAB implementation of consistent discretisations on complex grids. *Comput. Geosci.* 1–26.
- Loepky, J.L., Sacks, J., Welch, W.J., 2009. Choosing the sample size of a computer experiment: a practical guide. *Technometrics* 51 (4), 366–376.
- MacDonald, B., Ranjan, P., Chipman, H., 2013. GPfit: an R package for Gaussian process model fitting using a new optimization algorithm.
- MacQueen, J.B., 1967. Some methods for classification and analysis of multivariate observations. In: Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability. University of California Press, pp. 281–297.
- McKay, M.D., Beckman, R.J., Conover, W.J., 1979. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics* 21 (2), 239–245.
- Neal, M.R., 1997. Monte Carlo implementation of Gaussian process models for bayesian regression and classification.
- Nguyen, H.M., Couckuyt, I., Knockaert, L., Dhaene, T., Gorissen, D., Saeyns, Y., 2011. An alternative approach to avoid overfitting for surrogate models. In: Winter Simulation Conference. IEEE, pp. 2760–2771.
- O'Hagan, A., Kingman, J.F.C., 1978. Curve fitting and optimal design for prediction. *J. Roy. Statist. Soc. Ser. B* 40 (1), 1–42.
- Onwunalu, J., Durlofsky, L., 2010. Application of a particle swarm optimization algorithm for determining optimum well location and type. *Comput. Geosci.* 14, 183–198.
- Petelin, D., Filipič, B., Kocijan, J., 2011. Optimization of Gaussian Process Models with Evolutionary Algorithms. In: Lecture Notes in Computer Science, vol. 6593. Springer, Berlin Heidelberg, pp. 420–429.
- Ranjan, P., Haynes, R., Karsten, R., 2011. A computationally stable approach to Gaussian process interpolation of deterministic computer simulation data. *Technometrics* 52 (4), 366–378.
- Rasmussen, C.E., Williams, C.K.I., 2006. Gaussian processes for machine learning.
- Sacks, J., Welch, W., Mitchell, T., Wynn, H., 1989. Design and analysis of computer experiments. *Statist. Sci.* 4 (4), 409–435.
- Shanno, D., 1970. Conditioning of quasi-Newton methods for function minimization. *Math. Comp.* 24 (1), 647–656.
- Short, B., Higdon, D.M., Kronberg, P.P., 2007. Estimation of Faraday rotation measures of the near galactic sky using Gaussian process models. *Int. Soc. Bayesian Anal.* 2 (4), 665–680.
- SINTEF Applied Mathematics, 2012. Matlab reservoir simulator toolbox v. 2012a. <http://www.sintef.no/Projectweb/MRST/>.
- Wahba, G., 1978. Improper priors, spline smoothing and the problem of guarding against model errors in regression. *J. Roy. Statist. Soc. Ser. B* 40 (1), 364–372.
- Yeten, B., Durlofsky, L., Aziz, K., 2003. Optimization of nonconventional well type, location and trajectory. *SPE J.* 8 (3), 200–210.