



# A Radial Basis Function Method for Global Optimization

H.-M. GUTMANN\*

*Department of Applied Mathematics and Theoretical Physics, University of Cambridge, Silver Street, Cambridge CB3 9EW, England, UK (e-mail: h.m.gutmann@damtp.cam.ac.uk)*

**Abstract.** We introduce a method that aims to find the global minimum of a continuous nonconvex function on a compact subset of  $\mathbb{R}^d$ . It is assumed that function evaluations are expensive and that no additional information is available. Radial basis function interpolation is used to define a utility function. The maximizer of this function is the next point where the objective function is evaluated. We show that, for most types of radial basis functions that are considered in this paper, convergence can be achieved without further assumptions on the objective function. Besides, it turns out that our method is closely related to a statistical global optimization method, the P-algorithm. A general framework for both methods is presented. Finally, a few numerical examples show that on the set of Dixon-Szegö test functions our method yields favourable results in comparison to other global optimization methods.

**Key words:** Global optimization, radial basis functions, interpolation, P-algorithm.

## 1. Introduction

Global optimization has attracted a lot of attention in the last 20 years. In many applications, the objective function is nonlinear and nonconvex. Often, the number of local minima is large. Therefore standard nonlinear programming methods may fail to locate the global minimum.

In the most general way, the Global Optimization Problem can be stated as

$$\text{(GOP)} \quad \text{find } x^* \in \mathcal{D} \text{ such that } f(x^*) \leq f(x), \quad x \in \mathcal{D},$$

where  $\mathcal{D} \subset \mathbb{R}^d$  is compact, and  $f : \mathcal{D} \rightarrow \mathbb{R}$  is a continuous function defined on  $\mathcal{D}$ . Under these assumptions, (GOP) is solvable, because  $f$  attains its minimum on  $\mathcal{D}$ .

Numerous methods to solve (GOP) have been developed (see e.g. Horst and Pardalos [4] and Törn and Žilinskas [19]). Stochastic methods like simulated annealing and genetic algorithms which use only function values are very popular among users, although their rate of convergence is usually rather slow. Deterministic methods like Branch-and-Bound, however, assume that one can compute a lower bound of  $f$  on a subset of  $\mathcal{D}$ . This can be done, for example, when the

---

\* Supported by an Engineering and Physical Sciences Research Council grant and a doctoral scholarship (HSP III) of the German Academic Exchange Service.

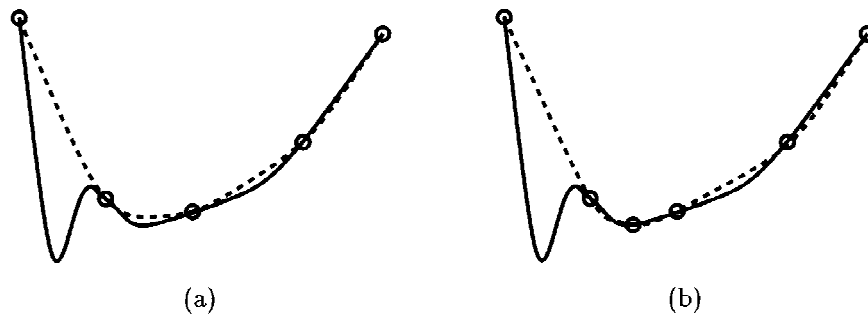


Figure 1. The function whose graph is the solid line is to be minimized. The dots in (a) indicate the points where the function values are known. The dotted line in (a) is the graph of the response surface. Sampling the function at the global minimizer of this surface gives the new response surface in (b). A better estimate of a local minimum has been found, but the global minimum is missed.

Lipschitz constant on  $f$  is available. The further assumptions make these methods very powerful, but often they are not satisfied or it is too expensive to provide the necessary information.

For the method investigated in this paper, we have in mind problems when the only information available is the possibility to evaluate the objective function, and each evaluation is very expensive. This may mean that it takes several hours to calculate a function value. For example, a function evaluation at a point may be done by building an experiment, by running a long computer simulation or by using a finite element method. Therefore, the duration of an optimization process is dominated by the function evaluations. As it can take very long to compute a global minimum in such a case, users often are satisfied when an adequate estimate of the global minimum is obtained. Thus, our goal is to require as few function evaluations as possible to find such an estimate.

Response surface methods have been developed to solve this kind of problem. Given points and their function values, a response surface can be computed that interpolates the objective function at these points. For many smooth objective functions such a response surface can identify the region of a global minimum after only a few function evaluations.

After having found a response surface, a naive idea would be to choose the global minimizer of the surface and evaluate the objective function there. However, if this process is iterated, the global minimum might be missed (see Figure 1). This happens because one trusts the surface model without taking into account possible errors.

To avoid this problem the decision on where to evaluate the objective function next must be based on the response surface model and a measure of the error in this model. If one knew the global minimum value, one could choose any point  $x$  and assume that it is a global minimizer. Then a response surface can be fitted through this point and the existing points. Intuitively, if this surface is very ‘bumpy’, it is

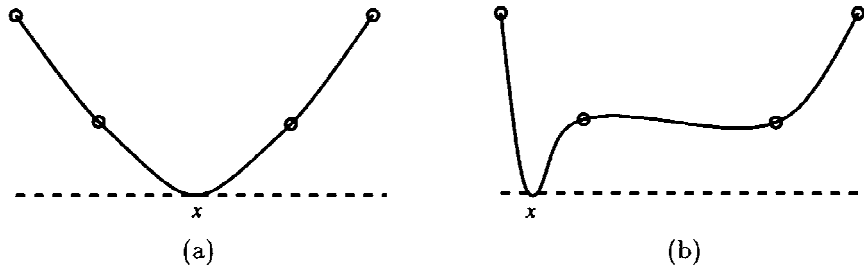


Figure 2. An example of the measure of ‘bumpiness’, where the dotted line indicates the global minimum. The response surface in (a) is less ‘bumpy’ than the one in (b).

unreasonable to expect that  $x$  is a global minimizer. So one would choose the next point to be the one that yields the ‘least bumpy’ of all these response surfaces. Normally, of course, the optimal value is not known. Then one can choose an estimate instead of the true value and follow the idea above. An example of two different levels of ‘bumpiness’ is given in Figure 2.

A general response surface technique has been proposed by Jones [8]. Let  $\mathcal{A}$  be a linear space of functions, and assume that, for  $s \in \mathcal{A}$ ,  $\sigma(s)$  is a measure of the ‘bumpiness’ of  $s$ . Now assume that we have calculated  $x_1, \dots, x_n$  and the function values  $f(x_1), \dots, f(x_n)$ . A target value  $f^*$  is chosen that can be regarded as an estimate of the optimal value, but it might be very crude. For each  $y \notin \{x_1, \dots, x_n\}$ , let  $s_y \in \mathcal{A}$  be defined by the interpolation conditions

$$\begin{aligned} s_y(x_i) &= f(x_i), \quad i = 1, \dots, n, \\ s_y(y) &= f^*. \end{aligned} \tag{1.1}$$

The new point  $x_{n+1}$  is chosen to be the value of  $y$  that minimizes  $\sigma(s_y)$ ,  $y \notin \{x_1, \dots, x_n\}$ .

Our method is based on this technique where we use radial basis functions as interpolants. Their interpolation properties are very suitable. Specifically, the uniqueness of an interpolant is achieved under very mild conditions on the location of the interpolation points, and a measure of bumpiness is also available.

Close relations can be established between our method and one from statistical global optimization, namely the P-algorithm (Žilinskas [22]). Although being derived using a completely different approach, it is very similar to our method. One special case of a P-algorithm, developed by Kushner [12], is even equivalent to a special case of our radial basis function method.

Other global optimization methods based on radial basis functions have been developed. Alotto et al. [1] use interpolation by multiquadrics to accelerate a simulated annealing method. Ishikawa et al. [6, 7] employ radial basis functions to estimate the global minimizer and run an SQP algorithm to locate it.

The properties of radial basis functions that are necessary for the description of our method are introduced in the following section. In particular, we address

the question of interpolation and introduce a suitable measure of ‘bumpiness’. The global optimization method is described in detail in Section 3. Convergence of the method is the subject of Section 4. The proof of the main theorem can be found in the Appendix. The relation between our method and the P-algorithm is addressed in Section 5. The final section deals with search strategies, but a complete analysis is beyond the scope of this paper.

## 2. Interpolation by Radial Basis Functions and a Measure of Bumpiness

The radial basis function interpolation problem is as follows. Let  $n$  pairwise different points  $x_1, \dots, x_n \in \mathbb{R}^d$  and data  $f_1, \dots, f_n \in \mathbb{R}$  be given, where  $n$  and  $d$  are any positive integers. We seek a function  $s$  of the form

$$s(x) = \sum_{i=1}^n \lambda_i \phi(\|x - x_i\|) + p(x), \quad x \in \mathbb{R}^d, \quad (2.1)$$

that interpolates the data  $(x_1, f_1), \dots, (x_n, f_n)$ . The coefficients  $\lambda_i$ ,  $i = 1, \dots, n$ , are real numbers, and the norm  $\|\cdot\|$  is the Euclidean norm in  $\mathbb{R}^d$ .  $p$  is from  $\Pi_m$ , the space of polynomials of degree less than or equal to  $m$ , i.e. it can be expressed as a linear combination of functions  $x_1^{k_1} \dots x_d^{k_d}$ ,  $x \in \mathbb{R}^d$ , where  $k_1 + \dots + k_d \leq m$ . We let  $\Pi_{-1} := \{0\}$ . The following choices of  $\phi$  are considered:

$$\left. \begin{aligned} \phi(r) &= r && \text{(linear),} \\ \phi(r) &= r^3 && \text{(cubic),} \\ \phi(r) &= r^2 \log r && \text{(thin plate spline),} \\ \phi(r) &= \sqrt{r^2 + \gamma^2} && \text{(multiquadric),} \\ \phi(r) &= e^{-\gamma r^2} && \text{(Gaussian),} \end{aligned} \right\} r \geq 0, \quad (2.2)$$

where  $\gamma$  is a prescribed positive constant.

It would be obvious to set  $m = -1$  so that (2.1) is a linear combination of the basis functions  $\phi(\| \cdot - x_i \|)$ ,  $i = 1, \dots, n$ , only. However, the matrix  $\Phi \in \mathbb{R}^{n \times n}$  that is defined by

$$(\Phi)_{ij} := \phi(\|x_i - x_j\|), \quad i, j = 1, \dots, n, \quad (2.3)$$

might be singular. For example, if  $\phi(r) = r^2 \log r$ ,  $n = d + 1$  and the points  $x_1, \dots, x_{d+1}$  form a simplex where all the edges have length 1, then  $\Phi = 0$ . So for  $m = -1$  and nonzero data there is no interpolant (2.1). However, any data  $f_1, \dots, f_{d+1}$  can be interpolated by a linear polynomial. Thus the interpolant (2.1) exists if  $\lambda_i = 0$ ,  $i = 1, \dots, n$ , and if  $p$  is this interpolating polynomial. Further, the general form (2.1) allows more freedom in defining a suitable measure of bumpiness. In general, our task now is to find values of  $m$  that guarantee existence and uniqueness of an interpolant (2.1) and a measure of its bumpiness.

The key to this task is the concept of conditional definiteness. For each  $\phi$  from (2.2),  $\Phi$  is conditionally positive or negative definite. Specifically, let  $\mathcal{V}_m \subset \mathbb{R}^n$  be

the linear space of all  $\lambda \in \mathbb{R}^n$  that satisfy

$$\sum_{i=1}^n \lambda_i q(x_i) = 0 \quad \forall q \in \Pi_m. \tag{2.4}$$

Formally, we set  $\mathcal{V}_{-1} := \mathbb{R}^n$ . Obviously,  $\mathcal{V}_{m+1} \subset \mathcal{V}_m$  for all  $m \geq -1$ . Powell [15] shows that, in the cubic and thin plate spline cases

$$\lambda^T \Phi \lambda > 0 \quad \forall \lambda \in \mathcal{V}_1 \setminus \{0\}, \tag{2.5}$$

in the linear and multiquadric cases

$$\lambda^T \Phi \lambda < 0 \quad \forall \lambda \in \mathcal{V}_0 \setminus \{0\}, \tag{2.6}$$

and in the Gaussian case

$$\lambda^T \Phi \lambda > 0 \quad \forall \lambda \in \mathbb{R}^n \setminus \{0\}. \tag{2.7}$$

We let  $m_0$  be 1 in the cubic and thin plate spline cases, 0 in the linear and multiquadric cases and  $-1$  in the Gaussian case. Then the inequalities (2.5) – (2.7) can be merged into

$$(-1)^{m_0+1} \lambda^T \Phi \lambda > 0 \quad \forall \lambda \in \mathcal{V}_{m_0} \setminus \{0\}. \tag{2.8}$$

After choosing  $\phi$ , we let  $m$  be an integer that is not less than  $m_0$ , and  $\lambda$  is confined to  $\mathcal{V}_m$ .

Let  $\hat{m}$  be the dimension of  $\Pi_m$ , let  $p_1, \dots, p_{\hat{m}}$  be a basis of this linear space, and let  $P$  be the matrix

$$P := \begin{pmatrix} p_1(x_1) & \cdots & p_{\hat{m}}(x_1) \\ \vdots & & \vdots \\ p_1(x_n) & \cdots & p_{\hat{m}}(x_n) \end{pmatrix}. \tag{2.9}$$

Then  $\mathcal{V}_m$  is the space of all  $\lambda \in \mathbb{R}^n$  that satisfy  $P^T \lambda = 0$ . Further, it can be shown (see [15]) that the matrix

$$A = \begin{pmatrix} \Phi & P \\ P^T & 0 \end{pmatrix} \in \mathbb{R}^{(n+\hat{m}) \times (n+\hat{m})} \tag{2.10}$$

is nonsingular if and only if  $x_1, \dots, x_n$  satisfy

$$q \in \Pi_m \quad \text{and} \quad q(x_i) = 0, \quad i = 1, \dots, n, \quad \implies \quad q \equiv 0. \tag{2.11}$$

In the Gaussian case with  $m = -1$ ,  $P$  and condition (2.11) are omitted. Therefore the coefficients of the function  $s$  in (2.1) are defined uniquely by the system

$$\left. \begin{aligned} s(x_i) &= f_i, \quad i = 1, \dots, n \\ \sum_{i=1}^n \lambda_i p_j(x_i) &= 0, \quad j = 1, \dots, \hat{m} \end{aligned} \right\}. \tag{2.12}$$

Let  $F$  be the vector whose entries are the data values  $f_1, \dots, f_n$ . Then the system (2.12) becomes

$$\begin{pmatrix} \Phi & P \\ P^T & 0 \end{pmatrix} \begin{pmatrix} \lambda \\ c \end{pmatrix} = \begin{pmatrix} F \\ 0_{\hat{m}} \end{pmatrix}, \quad (2.13)$$

where  $\lambda = (\lambda_1, \dots, \lambda_n)^T \in \mathbb{R}^n$ ,  $c \in \mathbb{R}^{\hat{m}}$  and  $0_{\hat{m}}$  is the zero in  $\mathbb{R}^{\hat{m}}$ . The components of  $c$  are the coefficients of the polynomial  $p$  with respect to the basis  $p_1, \dots, p_{\hat{m}}$ .

The motivation for the measurement of the bumpiness of a radial basis function interpolant can be developed from the theory of natural cubic splines in one dimension. They can be written in the form (2.1), where  $\phi(r) = r^3$ ,  $\lambda \in \mathcal{V}_1$  and  $p \in \Pi_1$ . It is well known (e.g. Powell [14]) that the interpolant  $s$  that is defined by the system (2.12) minimizes  $I(g) := \int_{\mathbb{R}} [g''(x)]^2 dx$  among all functions  $g : \mathbb{R} \rightarrow \mathbb{R}$  that satisfy the interpolation conditions  $g(x_i) = f_i$ ,  $i = 1, \dots, n$ , and for which  $I(g)$  exists and is finite. Therefore  $I(g)$  is a suitable measure of bumpiness. The second derivative  $s''$  is piecewise linear and vanishes outside a bounded interval. Thus one obtains by integration by parts

$$\begin{aligned} I(s) &= \int_{\mathbb{R}} [s''(x)]^2 dx = 12 \sum_{i=1}^n \lambda_i s(x_i) \\ &= 12 \sum_{i=1}^n \lambda_i \left( \sum_{j=1}^n \lambda_j |x_i - x_j|^3 + p(x_i) \right) = 12 \lambda^T \Phi \lambda, \end{aligned}$$

where the last equation follows from  $\lambda \in \mathcal{V}_1$ . This relation suggests that expression (2.8) can provide a semi-inner product and a semi-norm for each  $\phi$  in (2.2) and  $m \geq m_0$ . Also, the semi-norm will be the measure of bumpiness of a radial basis function (2.1). A semi-inner product  $\langle \cdot, \cdot \rangle$  satisfies the same properties as an inner product, except that  $\langle s, s \rangle = 0$  need not imply  $s = 0$ . Similarly, for a semi-norm  $\|\cdot\|$ ,  $\|s\| = 0$  does not imply  $s = 0$ .

We choose any radial basis function from (2.2) and  $m \geq m_0$ , and we define  $\mathcal{A}_{\phi, m}$  to be the linear space of all functions of the form

$$\sum_{i=1}^N \lambda_i \phi(\|x - y_i\|) + p(x), \quad x \in \mathbb{R}^d,$$

where  $N \in \mathbb{N}$ ,  $y_1, \dots, y_N \in \mathbb{R}^d$ ,  $p \in \Pi_m$ , and  $\lambda = (\lambda_1, \dots, \lambda_N)^T$  satisfies (2.4) for  $n = N$ . On this space, the semi-inner product and the semi-norm are defined as follows. Let  $s$  and  $u$  be any functions in  $\mathcal{A}_{\phi, m}$ , i.e.

$$s(x) = \sum_{i=1}^{N(s)} \lambda_i \phi(\|x - y_i\|) + p(x) \quad \text{and} \quad u(x) = \sum_{j=1}^{N(u)} \mu_j \phi(\|x - z_j\|) + q(x).$$

We let the semi-inner product be the expression

$$\langle s, u \rangle := (-1)^{m_0+1} \sum_{i=1}^{N(s)} \lambda_i u(y_i). \quad (2.14)$$

Clearly, it is bilinear. To show symmetry, we use

$$\sum_{i=1}^{N(s)} \lambda_i q(y_i) = 0 \quad \text{and} \quad \sum_{j=1}^{N(u)} \mu_j p(z_j) = 0,$$

to deduce

$$\begin{aligned} \langle s, u \rangle &= (-1)^{m_0+1} \sum_{i=1}^{N(s)} \lambda_i \left( \sum_{j=1}^{N(u)} \mu_j \phi(\|y_i - z_j\|) + q(y_i) \right) \\ &= (-1)^{m_0+1} \sum_{i=1}^{N(s)} \sum_{j=1}^{N(u)} \lambda_i \mu_j \phi(\|y_i - z_j\|) \\ &= (-1)^{m_0+1} \sum_{j=1}^{N(u)} \mu_j \left( \sum_{i=1}^{N(s)} \lambda_i \phi(\|z_j - y_i\|) + p(z_j) \right) \\ &= (-1)^{m_0+1} \sum_{j=1}^{N(u)} \mu_j s(z_j) = \langle u, s \rangle. \end{aligned}$$

By (2.8),

$$\langle s, s \rangle = (-1)^{m_0+1} \sum_{i=1}^{N(s)} \lambda_i s(y_i) = (-1)^{m_0+1} \sum_{i,j=1}^{N(s)} \lambda_i \lambda_j \phi(\|y_i - y_j\|) \quad (2.15)$$

is strictly positive, if  $\lambda \in \mathcal{V}_m \setminus \{0\}$  and  $m \geq m_0$ , i.e.  $s \in \mathcal{A}_{\phi,m}$  is not a polynomial in  $\Pi_m$ . Thus (2.14) is a semi-inner product on  $\mathcal{A}_{\phi,m}$  that induces the semi-norm  $\langle s, s \rangle$  with null space  $\Pi_m$  (for details see Powell [16] and Schaback [17]).

In analogy to the variational principle for cubic splines in one dimension, mentioned above, there is a theorem that states that the given interpolant is the solution to a minimization problem.

**THEOREM 1.** (Schaback [17]). *Let  $\phi$  be any radial basis function from (2.2), and let  $m$  be chosen such that  $m \geq m_0$ . Given are points  $x_1, \dots, x_n$  in  $\mathbb{R}^d$  having the property (2.11) and values  $f_1, \dots, f_n$  in  $\mathbb{R}$ . Let  $s$  be the radial function of the form (2.1) that solves the system (2.12). Then  $s$  minimizes the semi-norm  $\langle g, g \rangle^{1/2}$  on the set of functions  $g \in \mathcal{A}_{\phi,m}$  that satisfy*

$$g(x_i) = f_i, \quad i = 1, \dots, n. \quad (2.16)$$

### 3. A Radial Basis Function Method

It will be shown how radial basis functions can be used in the general method of Jones [8] to solve the problem (GOP) (cf. Powell [16]). As in Section 2, we pick  $\phi$  from (2.2) and  $m \geq m_0$ . Let  $p_1, \dots, p_{\hat{m}}$  be a basis of  $\Pi_m$ , where  $\hat{m} = \dim \Pi_m$ . Assume we have chosen  $x_1, \dots, x_n \in \mathcal{D}$  that satisfy (2.11), and we know the function values  $f(x_1), \dots, f(x_n)$ . Let the function

$$s_n(x) = \sum_{i=1}^n \lambda_i \phi(\|x - x_i\|) + p(x), \quad x \in \mathbb{R}^d,$$

interpolate  $(x_1, f(x_1)), \dots, (x_n, f(x_n))$ . Our task is to determine  $x_{n+1}$ . For a target value  $f_n^*$  and a point  $y \in \mathcal{D} \setminus \{x_1, \dots, x_n\}$  the radial basis function  $s_y$  that satisfies (1.1) can be written as

$$s_y(x) = s_n(x) + [f_n^* - s_n(y)] \ell_n(y, x), \quad x \in \mathbb{R}^d, \quad (3.1)$$

where  $\ell_n(y, x)$  is the radial basis function solution to the interpolation conditions

$$\begin{aligned} \ell_n(y, x_i) &= 0, \quad i = 1, \dots, n, \\ \ell_n(y, y) &= 1. \end{aligned} \quad (3.2)$$

Therefore  $\ell_n(y, \cdot)$  can be expressed as

$$\ell_n(y, x) = \sum_{i=1}^n \alpha_i(y) \phi(\|x - x_i\|) + \mu_n(y) \phi(\|x - y\|) + \sum_{i=1}^{\hat{m}} b_i(y) p_i(x), \quad x \in \mathbb{R}^d. \quad (3.3)$$

As in equation (2.10), let  $A(y)$  be the matrix

$$A(y) := \begin{pmatrix} \Phi & u(y) & P \\ u(y)^T & \phi(0) & \pi(y)^T \\ P^T & \pi(y) & 0_{\hat{m} \times \hat{m}} \end{pmatrix}, \quad (3.4)$$

where  $u(y)$  and  $\pi(y)$  are the vectors

$$u(y) := (\phi(\|y - x_1\|), \dots, \phi(\|y - x_n\|))^T$$

and

$$\pi(y) := (p_1(y), \dots, p_{\hat{m}}(y))^T,$$

respectively. Then the coefficients of  $\ell_n(y, \cdot)$  are defined by the equations

$$A(y) \begin{pmatrix} \alpha(y) \\ \mu_n(y) \\ b(y) \end{pmatrix} = \begin{pmatrix} 0_n \\ 1 \\ 0_{\hat{m}} \end{pmatrix}, \quad (3.5)$$



where  $\alpha(y) = (\alpha_1(y), \dots, \alpha_n(y))^T \in \mathbb{R}^n$ ,  $b(y) = (b_1(y), \dots, b_{\hat{m}}(y))^T \in \mathbb{R}^{\hat{m}}$ ,  $\mu_n(y) \in \mathbb{R}$ ,  $0_n$  and  $0_{\hat{m}}$  denote the zero in  $\mathbb{R}^n$  and  $\mathbb{R}^{\hat{m}}$ , respectively.

The square of the semi-norm  $\langle s_y, s_y \rangle$  of the new interpolant (3.1), as defined in the previous section, has the value

$$\begin{aligned} \langle s_y, s_y \rangle &= \langle s_n, s_n \rangle + 2[f_n^* - s_n(y)] \langle s_n, \ell_n(y, \cdot) \rangle \\ &\quad + [f_n^* - s_n(y)]^2 \langle \ell_n(y, \cdot), \ell_n(y, \cdot) \rangle. \end{aligned}$$

Equations (2.14) and (3.2) imply

$$\langle s_n, \ell_n(y, \cdot) \rangle = (-1)^{m_0+1} \sum_{i=1}^n \lambda_i \ell_n(y, x_i) = 0,$$

and, using expressions (3.2) and (3.3), we find the expression

$$\begin{aligned} \langle \ell_n(y, \cdot), \ell_n(y, \cdot) \rangle &= (-1)^{m_0+1} \left[ \sum_{i=1}^n \alpha_i(y) \ell_n(y, x_i) + \mu_n(y) \ell_n(y, y) \right] \\ &= (-1)^{m_0+1} \mu_n(y). \end{aligned} \quad (3.6)$$

Thus we deduce the formula

$$\langle s_y, s_y \rangle = \langle s_n, s_n \rangle + (-1)^{m_0+1} \mu_n(y) [f_n^* - s_n(y)]^2.$$

Further, we define the function  $g_n : \mathcal{D} \setminus \{x_1, \dots, x_n\} \rightarrow \mathbb{R}$  as the difference

$$g_n(y) := \langle s_y, s_y \rangle - \langle s_n, s_n \rangle = (-1)^{m_0+1} \mu_n(y) [f_n^* - s_n(y)]^2,$$

which is nonnegative. Since  $\langle s_n, s_n \rangle$  is independent of  $y$ , the required minimization of  $\langle s_y, s_y \rangle$  and the minimization of  $g_n(y)$  are equivalent.

The choice of  $f_n^*$  determines the location of  $x_{n+1}$ . If

$$\max_{y \in \mathcal{D}} s_n(y) \geq f_n^* \geq \min_{y \in \mathcal{D}} s_n(y),$$

then  $g_n(y) = 0$  can be achieved. However, if

$$f_n^* < \min_{y \in \mathcal{D}} s_n(y),$$

then  $x_{n+1}$  will be away from the  $x_i$ ,  $i = 1, \dots, n$ . In particular, for  $f_n^* \rightarrow -\infty$ , we make the following deduction.

**REMARK 2.** For  $f_n^* < \min_{y \in \mathcal{D}} s_n(y)$  let  $x(f_n^*)$  be the minimizer of  $g_n$ , i.e.

$$\begin{aligned} (-1)^{m_0+1} \mu_n(x(f_n^*)) [s_n(x(f_n^*)) - f_n^*]^2 &\leq (-1)^{m_0+1} \mu_n(y) [s_n(y) - f_n^*]^2 \\ &\quad \forall y \in \mathcal{D} \setminus \{x_1, \dots, x_n\}. \end{aligned}$$

This is equivalent to

$$(-1)^{m_0+1} \mu_n(x(f_n^*)) \leq (-1)^{m_0+1} \mu_n(y) \left[ 1 + \frac{s_n(y) - s_n(x(f_n^*))}{s_n(x(f_n^*)) - f_n^*} \right]^2.$$

As  $f_n^* \rightarrow -\infty$ , the boundedness of  $s_n$  on  $\mathcal{D}$  implies

$$(-1)^{m_0+1} \mu_n(x(-\infty)) \leq (-1)^{m_0+1} \mu_n(y) \quad \forall y \in \mathcal{D} \setminus \{x_1, \dots, x_n\}.$$

Therefore, the choice  $f_n^* = -\infty$  requires the minimization of the function  $(-1)^{m_0+1} \mu_n(y)$ . This process puts  $x_{n+1}$  in a large gap between  $x_i$ ,  $i = 1, \dots, n$ , a property that is of fundamental importance to global optimization.

The following basic algorithm employs the given method.

### ALGORITHM 3.

**Initial step:** Pick  $\phi$  from (2.2) and  $m \geq m_0$ .

Choose points  $x_1, \dots, x_{n_0} \in \mathcal{D}$  that satisfy (2.11). Compute the radial function  $s_n$  that minimizes  $\langle s, s \rangle$  on  $\mathcal{A}_{\phi, m}$ , subject to the interpolation conditions

$$s(x_i) = f(x_i), \quad i = 1, \dots, n.$$

**Iteration step:**  $x_1, \dots, x_n$  are the points in  $\mathcal{D}$  where the value of  $f$  is known, and  $s_n$  minimizes  $\langle s, s \rangle$ , subject to  $s(x_i) = f(x_i)$ ,  $i = 1, \dots, n$ .

Choose a target value  $f_n^* \in [-\infty, \min_{y \in \mathcal{D}} s_n(y)]$ . (The choice  $f_n^* = \min s_n(y)$  is admissible only if none of the  $x_i$  is a global minimizer of  $s_n$ ).

Calculate  $x_{n+1}$ , which is the value of  $y$  that minimizes the function

$$g_n(y) = (-1)^{m_0+1} \mu_n(y) [s_n(y) - f_n^*]^2, \quad y \in \mathcal{D} \setminus \{x_1, \dots, x_n\}. \quad (3.7)$$

Evaluate  $f$  at  $x_{n+1}$  and set  $n := n + 1$ .

Stop, if  $n$  is greater than a prescribed  $n_{max}$ .

The function  $g_n$  is infinitely differentiable on  $\mathcal{D} \setminus \{x_1, \dots, x_n\}$ , but is not defined at the interpolation points. If  $f_n^* = \min s_n(y)$ ,  $y \in \mathcal{D}$ , and if  $s_n(x_i) > f_n^*$ ,  $i = 1, \dots, n$ , then the global minimizers of  $g_n$  are the global minimizers of  $s_n$ . Thus one can minimize  $s_n$ , which is defined on the whole of  $\mathcal{D}$ , to obtain  $x_{n+1}$ . If  $f_n^* < \min s_n(y)$ , however, then  $g_n(x)$  tends to infinity as  $x$  tends to  $x_i$ ,  $i = 1, \dots, n$ . Let  $h_n : \mathcal{D} \rightarrow \mathbb{R}$  be defined as

$$h_n(y) := \begin{cases} \frac{1}{g_n(y)}, & y \notin \{x_1, \dots, x_n\} \\ 0, & y = x_i, \quad i = 1, \dots, n \end{cases}. \quad (3.8)$$

The maximization of  $h_n$  on  $\mathcal{D}$  is equivalent to the minimization of  $g_n$ . Further,  $h_n$  is infinitely differentiable on  $\mathcal{D} \setminus \{x_1, \dots, x_n\}$ . It can also be shown, using the system

(3.5), that it is in  $C(\mathcal{D})$  in the linear case, in  $C^1(\mathcal{D})$  in the thin plate spline case, in  $C^2(\mathcal{D})$  in the cubic case, and in  $C^\infty(\mathcal{D})$  in the multiquadric and Gaussian cases.

Under certain conditions on  $f$  and the values  $f_n^*, n \rightarrow \infty$ , it can be proved that a subsequence of the generated points  $(x_n)_{n \in \mathbb{N}}$  converges to a global minimum. This is the subject of the following section.

#### 4. Convergence of the Method

Our aim is to prove convergence of the method for any continuous function  $f$ . A theorem by Törn and Zilinskas [19] tells us that the sequence that is generated by Algorithm 3 should be dense. Applied to our method, it states

**THEOREM 4.** *The algorithm converges for every continuous function  $f$  if and only if it generates a sequence of points that is dense in  $\mathcal{D}$ .*

So our task is to establish the density of the sequence of generated points.

The convergence result does not allow a free choice of the target values  $f_n^*$ . Figure 1 shows that the global minimum might be missed, if  $f_n^*$  is set to  $\min_{y \in \mathcal{D}} s_n(y)$  on each iteration, provided this choice is admissible. Therefore, we have to assume that enough of the numbers  $\min s_n(y) - f_n^*$  are sufficiently large. Specifically, let  $\tau > 0$  and  $\rho \geq 0$  be constants, where additionally  $\rho < 1$  in the linear case and  $\rho < 2$  in the thin plate spline and cubic cases, and define

$$\Delta_n := \min_{1 \leq i \leq n-1} \|x_n - x_i\|. \tag{4.1}$$

Then the condition

$$\min_{y \in \mathcal{D}} s_n(y) - f_n^* > \tau \Delta_n^{\rho/2} \|s_n\|_\infty, \tag{4.2}$$

for infinitely many  $n \in \mathbb{N}$ , will lead to the required result. Here  $\|\cdot\|_\infty$  denotes the maximum norm of a function on  $\mathcal{D}$ , defined by

$$\|g\|_\infty := \max_{x \in \mathcal{D}} |g(x)|, \quad g \in C(\mathcal{D}).$$

We note that the norms  $\|s_n\|_\infty$  may diverge as  $n \rightarrow \infty$ .

Unfortunately, the choice of  $\phi$  and  $m$  is restricted. In the proof of Theorem 7 we need the result that, for any  $y \in \mathcal{D}$  and any neighbourhood  $U$  of  $y$ ,  $(-1)^{m_0+1} \mu_n(y)$  can be bounded above by a number that does not depend on  $n$ , if none of the points  $x_1, \dots, x_n$  is in  $U$ . This condition is achieved, if there is a function that takes the value 1 at  $y$ , that is identically zero outside  $U$ , and that is in the function space  $\mathcal{N}_{\phi,m}(\mathbb{R}^d)$  as defined below.

**DEFINITION 5.** *Let  $\phi$  from (2.2) and  $m \geq m_0$  be given. A continuous function  $F : D \rightarrow \mathbb{R}$ ,  $D \subset \mathbb{R}^d$ , belongs to the function space  $\mathcal{N}_{\phi,m}(D)$ , if there exists a positive constant  $C$  such that, for any choice of interpolation points  $x_1, \dots, x_n \in$*

$D$  for which (2.11) holds, the interpolant  $s_n \in \mathcal{A}_{\phi,m}$  to  $F$  at these points has the property

$$\langle s_n, s_n \rangle \leq C.$$

The characterization of  $\mathcal{N}_{\phi,m}(D)$  is rather abstract. In the linear, cubic and thin plate spline cases, the following proposition that is taken from Gutmann [3] provides a useful criterion to check whether it is satisfied. In the multiquadric and Gaussian cases, however, no such criterion is known.

**PROPOSITION 6.** *Let  $\phi(r) = r$ ,  $\phi(r) = r^2 \log r$  or  $\phi(r) = r^3$ . Further, let  $\kappa = 1$  in the linear case,  $\kappa = 2$  in the thin plate spline case and  $\kappa = 3$  in the cubic case, and choose the integer  $m$  such that  $0 \leq m \leq d$  in the linear case,  $1 \leq m \leq d + 1$  in the thin plate spline case and  $1 \leq m \leq d + 2$  in the cubic case. Define  $v := (d + \kappa)/2$  if  $d + \kappa$  is even, and  $v := (d + \kappa + 1)/2$  otherwise. If  $F \in C^v(\mathbb{R}^d)$  has bounded support, then  $F \in \mathcal{N}_{\phi,m}(\mathbb{R}^d)$ .*

Global convergence will be established only for the cases covered by this proposition. It remains an open problem whether a similar property is achieved in other cases. Thus we have the following theorem.

**THEOREM 7.** *Let  $\phi(r) = r$ ,  $\phi(r) = r^2 \log r$  or  $\phi(r) = r^3$ . Further, choose the integer  $m$  such that  $0 \leq m \leq d$  in the linear case,  $1 \leq m \leq d + 1$  in the thin plate spline case and  $1 \leq m \leq d + 2$  in the cubic case. Let  $(x_n)_{n \in \mathbb{N}}$  be the sequence generated by Algorithm 3, and  $s_n$  be the radial function that interpolates  $(x_i, f(x_i))$ ,  $i = 1, \dots, n$ . Assume that, for infinitely many  $n \in \mathbb{N}$ , the choice of  $f_n^*$  satisfies (4.2). Then the sequence  $(x_n)$  is dense in  $\mathcal{D}$ .*

The proof of Theorem 7 is given in the Appendix.

A particular convergence result follows immediately from Theorems 4 and 7, because the right hand side of (4.2) is some real number.

**COROLLARY 8.** *Let the assumptions of Theorem 7 on  $\phi$  and  $m$  hold. Further, let  $f$  be continuous, and, for infinitely many  $n \in \mathbb{N}$ , let  $f_n^* = -\infty$ . Then the method converges.*

An interesting question is to find conditions on  $f$  such that the maximum norm of an interpolant is uniformly bounded. If they hold, then the right-hand side of (4.2) can be replaced by  $\tau \Delta_n^{\rho/2}$ , so this constraint on  $f_n^*$  can be checked easily. We consider the special case of linear splines in one dimension, when  $d = 1$ ,  $\phi(r) = r$  and  $m = 0$ . For arbitrary points  $x_1, \dots, x_n$ , the piecewise linear interpolant  $s_n$  attains its maximum and minimum values at interpolation points. Thus  $\|s_n\|_\infty$  is bounded by  $\|f\|_\infty$ , a number that does not depend on the interpolation points. Therefore in this case the term  $\|s_n\|_\infty$  may be dropped from (4.2).

For other radial basis functions and other dimensions this simplification may fail. It is shown in the next lemma, however, that the uniform boundedness of

the semi-norm of an interpolant is sufficient for the uniform boundedness of the maximum norm. Thus, the second convergence result applies to functions  $f$  in  $\mathcal{N}_{\phi,m}(\mathcal{D})$ .

LEMMA 9. *Let  $f$  be in  $\mathcal{N}_{\phi,m}(\mathcal{D})$ . Further, let  $(x_n)_{n \in \mathbb{N}}$  be a sequence in  $\mathcal{D}$  with pairwise different elements, such that (2.11) holds for  $n = n_0$ . For  $n \geq n_0$ , denote the radial basis function interpolant to  $f$  at  $x_1, \dots, x_n$  by  $s_n$ . Then  $\|s_n\|_\infty$  is uniformly bounded by a number that only depends on  $x_1, \dots, x_{n_0}$ .*

*Proof.* We fix  $n$ , and we let  $y$  be any point of  $\mathcal{D} \setminus \{x_1, \dots, x_n\}$ . Let  $\tilde{s}_n$  be the radial function that interpolates  $(y, f(y))$  and  $(x_i, f(x_i))$ ,  $i = 1, \dots, n$ . By analogy with Equation (3.1), it can be written as

$$\tilde{s}_n(x) = s_n(x) + [f(y) - s_n(y)]\ell_n(y, x), \quad x \in \mathbb{R}^d,$$

where  $\ell_n(y, \cdot)$  is still the cardinal function that interpolates  $(x_i, 0)$ ,  $i = 1, \dots, n$ , and  $(y, 1)$ . Thus, as shown in Section 3,

$$\langle \tilde{s}_n, \tilde{s}_n \rangle = \langle s_n, s_n \rangle + [f(y) - s_n(y)]^2 (-1)^{m_0+1} \mu_n(y),$$

which gives the equation

$$[f(y) - s_n(y)]^2 = \frac{\langle \tilde{s}_n, \tilde{s}_n \rangle - \langle s_n, s_n \rangle}{(-1)^{m_0+1} \mu_n(y)}, \tag{4.3}$$

the value of  $(-1)^{m_0+1} \mu_n(y)$  being strictly positive.

Next we show that  $(-1)^{m_0+1} \mu_n(y)$  is bounded away from zero. Let  $\ell_{n_0}(y, \cdot)$  be the cardinal function that interpolates  $(x_1, 0), \dots, (x_{n_0}, 0)$  and  $(y, 1)$ . Then the semi-norm properties of  $\langle \cdot, \cdot \rangle$  and Theorem 1 imply

$$\begin{aligned} 0 < (-1)^{m_0+1} \mu_{n_0}(y) &= \langle \ell_{n_0}(y, \cdot), \ell_{n_0}(y, \cdot) \rangle \\ &\leq \langle \ell_n(y, \cdot), \ell_n(y, \cdot) \rangle = (-1)^{m_0+1} \mu_n(y). \end{aligned}$$

For  $n = n_0$ , let  $A$  and  $A(y)$  be the matrices (2.10) and (3.4), respectively. By using Cramer's Rule to solve (3.5), we find

$$\mu_{n_0}(y) = \frac{\det A}{\det A(y)}.$$

Now  $\det A$  is a nonzero constant, and  $\det A(y)$  is bounded on  $\mathcal{D}$ . It follows that  $(-1)^{m_0+1} \mu_{n_0}(y)$  is bounded away from zero. Therefore there exists a constant  $\alpha > 0$  such that

$$(-1)^{m_0+1} \mu_n(y) \geq \alpha \quad \forall y \in \mathcal{D} \setminus \{x_1, \dots, x_{n_0}\}, \quad n \geq n_0. \tag{4.4}$$

As  $f \in \mathcal{N}_{\phi,m}(\mathcal{D})$ ,  $\langle \tilde{s}_n, \tilde{s}_n \rangle$  is bounded above by a positive constant  $C$ . Further,  $\langle s_n, s_n \rangle$  is nonnegative. It follows from (4.3) and (4.4) that

$$|f(y) - s_n(y)| \leq \sqrt{\frac{C}{\alpha}}, \quad y \in \mathcal{D} \setminus \{x_1, \dots, x_n\}.$$

Further, because  $f$  is bounded on  $\mathcal{D}$ , we obtain

$$|s_n(y)| \leq \sqrt{\frac{C}{\alpha}} + \|f\|_\infty.$$

Note that the right-hand side is independent of  $n$  and  $y$ , as required. Alternatively, if  $y \in \{x_1, \dots, x_n\}$ , we have

$$|s_n(y)| = |f(y)| \leq \|f\|_\infty,$$

which completes the proof.  $\square$

Next, by applying Proposition 6, we obtain a criterion that ensures that  $f$  is in  $\mathcal{N}_{\phi,m}(\mathcal{D})$ .

**PROPOSITION 10.** *Let  $\phi$ ,  $m$  and  $\nu$  be defined as in Proposition 6, and let  $f \in C^\nu(\mathcal{D})$ , where  $\mathcal{D} \subset \mathbb{R}^d$  is compact. Then  $f \in \mathcal{N}_{\phi,m}(\mathcal{D})$ .*

*Proof.* By Whitney's theorem ([20]),  $f$  can be extended to a function  $F \in C^\nu(\mathbb{R}^d)$  that is equal to  $f$  on  $\mathcal{D}$ . Now  $\mathcal{D}$  is contained in a closed ball of radius  $\delta$ , say, and there is an infinitely differentiable function  $g$  with  $g(x) = 1$ ,  $\|x\| \leq \delta$ , and  $g(x) = 0$ ,  $\|x\| \geq 2\delta$ . Thus  $F \cdot g$  is in  $C^\nu(\mathbb{R}^d)$ , and by Proposition 6 it is in  $\mathcal{N}_{\phi,m}(\mathbb{R}^d)$ . Since  $F \cdot g$  is equal to  $f$  on  $\mathcal{D}$ , it follows from the definition of the semi-norm that  $f \in \mathcal{N}_{\phi,m}(\mathcal{D})$ .  $\square$

We complete this section by combining Theorem 7, Lemma 9 and Proposition 10.

**COROLLARY 11.** *Let the assumptions of Theorem 7 on  $\phi$  and  $m$  hold. Let  $\nu$  be as in Proposition 6, and let  $f \in C^\nu(\mathcal{D})$ . Further, assume that, for infinitely many  $n \in \mathbb{N}$ ,  $f_n^*$  has the property*

$$\min_{y \in \mathcal{D}} s_n(y) - f_n^* \geq \tau \Delta_n^{\rho/2},$$

where  $\tau > 0$  is a constant, and where  $\Delta_n$  and  $\rho$  are as in the beginning of this section. Then the method converges.

## 5. Relations to Statistical Global Optimization

In this section we consider the similarities between the given radial basis function method and the P-algorithm. The idea of that method is proposed by Kushner [11, 12] for one-dimensional problems. Here the objective function is regarded as a realization of a Brownian motion stochastic process. If real numbers  $x_1 < \dots < x_n$  are given, and their function values  $f(x_1), \dots, f(x_n)$  have been calculated, the model yields, for each  $x$  in the feasible set, a mean value  $\text{Mean}(x)$  and a variance  $\text{Var}(x)$ .  $\text{Mean}(x)$  serves as a prediction of the true function value at  $x$ , while  $\text{Var}(x)$  is a measure of uncertainty. It turns out that  $\text{Mean}$  is the piecewise linear interpolant

of the given data. The variance is piecewise quadratic on  $[x_1, x_n]$ , nonnegative and takes the value zero at  $x_1, \dots, x_n$ . For a real number  $x$ , let  $F_x$  be the normally distributed random variable  $F_x$  with mean  $\text{Mean}(x)$  and variance  $\text{Var}(x)$ . Now a nonnegative  $\epsilon_n$  is chosen, and the next point  $x_{n+1}$  will be the one that maximizes the utility function

$$P \left( F_x \leq \min_{i=1, \dots, n} f(x_i) - \epsilon_n \right), \quad x \in \mathcal{D}, \tag{5.1}$$

where  $P$  denotes probability. One can show that maximizing (5.1) is equivalent to maximizing

$$\frac{\text{Var}(x)}{[\text{Mean}(x) - \min\{f(x_1), \dots, f(x_n)\} + \epsilon_n]^2}, \quad x \in \mathcal{D}. \tag{5.2}$$

Compare our method in one dimension with the choice of linear splines, i.e.  $\phi(r) = r$  and  $m = 0$ , and with the target values

$$f_n^* = \min\{f(x_1), \dots, f(x_n)\} - \epsilon_n.$$

In this case, the interpolant  $s_n$  is identical to  $\text{Mean}$ . Further, except for a constant factor,

$$\text{Var}(x) = -\frac{1}{\mu_n(x)}.$$

Therefore, Kushner’s method and our method using linear splines are equivalent.

Žilinskas [21, 22] extends this approach to Gaussian random processes in several dimensions. He uses the selection rule (5.1) and introduces the name ‘P-algorithm’. In addition, he gives an axiomatic description of the terms involved in (5.1), namely the mean value function, the variance function and the utility function. We relate these results to our method.

Consider a symmetric function  $\sigma : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}, (x, z) \mapsto \sigma(x, z)$ , and assume that  $\sigma$  is conditionally positive or negative definite of order  $m$ . This means, there exists  $\alpha \in \{0, 1\}$  such that, given  $n$  different points  $x_1, \dots, x_n \in \mathbb{R}^d$  and multipliers  $\lambda_1, \dots, \lambda_n \in \mathbb{R}$ , we have

$$(-1)^\alpha \sum_{i,j=1}^n \lambda_i \lambda_j \sigma(x_i, x_j) > 0,$$

if the  $\lambda_i, i = 1, \dots, n$ , are not all zero and satisfy

$$\sum_{i=1}^n \lambda_i p(x_i) = 0, \quad p \in \Pi_m.$$

Denote the matrix with the elements  $\sigma(x_i, x_j), i, j = 1, \dots, n$ , by  $\Sigma$ , and the matrix with the elements  $p_j(x_i), i = 1, \dots, n, j = 1, \dots, \hat{m}$ , by  $P$ , where  $\{p_j :$

$j = 1, \dots, \hat{m}$  is a basis of  $\Pi_m$  and  $\hat{m}$  its dimension. The analogue of expression (2.10) is the matrix

$$A = \begin{pmatrix} \Sigma & P \\ P^T & 0 \end{pmatrix}. \quad (5.3)$$

We now let the interpolant to the components of  $F = (f(x_1), \dots, f(x_n))^T$  be the function

$$s(y) = \sum_{i=1}^n \lambda_i \sigma(y, x_i) + \sum_{j=1}^{\hat{m}} c_j p_j(y),$$

whose coefficients solve the system

$$A \begin{pmatrix} \lambda \\ c \end{pmatrix} = \begin{pmatrix} F \\ 0_{\hat{m}} \end{pmatrix},$$

It can be written as

$$s(y) = v_m(y)^T A^{-1} \begin{pmatrix} F \\ 0_{\hat{m}} \end{pmatrix}, \quad y \in \mathcal{D}, \quad (5.4)$$

where  $v_m(y)$  is the vector

$$v_m(y) := (\sigma(y, x_1), \dots, \sigma(y, x_n), p_1(y), \dots, p_{\hat{m}}(y))^T. \quad (5.5)$$

The nonnegative function

$$\text{Var}(y) = |\sigma(y, y) - v_m(y)^T A^{-1} v_m(y)|, \quad y \in \mathcal{D}, \quad (5.6)$$

is assumed to be a measure of uncertainty. Note that  $v_m(x_i)$  is the  $i$ -th column of  $A$ ,  $i = 1, \dots, n$ , so we obtain

$$\begin{aligned} \text{Var}(x_i) &= |\sigma(x_i, x_i) - v_m(x_i)^T A^{-1} v_m(x_i)| \\ &= |\sigma(x_i, x_i) - v_m(x_i)^T e_i| \\ &= 0. \end{aligned}$$

Thus there is no uncertainty at the interpolation points, which is meaningful because we know the true function values there.

For the P-algorithm,  $\sigma$  is interpreted as the correlation function of a Gaussian stochastic process. The use of a normal distribution, for example, gives  $\sigma(x, y) := e^{-\|x-y\|^2/2}$ , but other choices of  $\sigma$  are also considered in the literature. All of them are positive definite, so we set  $m = -1$ . The conditional mean and the conditional variance can be expressed as ([21])

$$\text{Mean}(y) = (f(x_1), \dots, f(x_n)) \Sigma^{-1} \begin{pmatrix} \sigma(y, x_1) \\ \vdots \\ \sigma(y, x_n) \end{pmatrix}, \quad (5.7)$$



$$\text{Var}(y) = \sigma(y, y) - \left( \sigma(y, x_1), \dots, \sigma(y, x_n) \right) \Sigma^{-1} \begin{pmatrix} \sigma(y, x_1) \\ \vdots \\ \sigma(y, x_n) \end{pmatrix}, \quad (5.8)$$

which agree with (5.4) and (5.6).

For our method, given  $\phi$  and  $m$ , it is suitable to define  $\sigma(x, y) := \phi(\|x - y\|)$ . Thus  $\Sigma = \Phi$ , and the coefficients of the interpolant  $s$  solve (2.13). This gives the form

$$s(y) = v_m(y)^T A^{-1} \begin{pmatrix} F \\ 0 \end{pmatrix}, \quad (5.9)$$

which is the same as the function (5.4).

We have seen already that  $|1/\mu_n|$  can be regarded as a variance in the case of linear splines in one dimension. An expression for it containing the matrix  $A$  and the vector  $v_m(x)$  can be derived in the following way. For any  $y \in \mathcal{D} \setminus \{x_1, \dots, x_n\}$ , consider the cardinal function (3.3). The second cardinality condition from (3.2) implies

$$\frac{1}{\mu_n(y)} = \phi(0) + \sum_{i=1}^n \frac{\alpha_i(y)}{\mu_n(y)} \phi(\|y - x_i\|) + \sum_{j=1}^{\hat{m}} \frac{b_j(y)}{\mu_n(y)} p_j(y). \quad (5.10)$$

The coefficients  $\alpha(y)$ ,  $\mu_n(y)$  and  $b(y)$  solve the system (3.5). Moreover, the vector (5.5) contains the first  $n$  and the last  $\hat{m}$  elements of the  $(n + 1)$ -th column of  $A(y)$ . Therefore  $\alpha(y)$  and  $b(y)$  also solve

$$A \begin{pmatrix} \alpha(y) \\ b(y) \end{pmatrix} = -\mu_n(y)v_m(y),$$

which implies

$$\frac{1}{\mu_n(y)} \begin{pmatrix} \alpha(y) \\ b(y) \end{pmatrix} = -A^{-1}v_m(y).$$

Thus, replacing the terms  $\alpha_i(y)/\mu_n(y)$  and  $b_j(y)/\mu_n(y)$  in (5.10), we find

$$\frac{1}{\mu_n(y)} = \phi(0) - v_m(y)^T A^{-1}v_m(y). \quad (5.11)$$

It follows from  $\sigma(x, x) = \phi(0)$ ,  $x \in \mathcal{D}$ , that expression (5.6) is equivalent to  $|1/\mu_n(y)|$ .

Finally, we consider the selection rule for the next point. For the P-algorithm, it has already been noted that the maximization of (5.1) is equivalent to the maximization of (5.2). For a given target value  $f^*$  define the function  $U : \mathbb{R}^2 \rightarrow \mathbb{R}$  as

$$U(M, V) := \frac{V^2}{(M - f^*)^2}. \quad (5.12)$$

It is increasing in  $V$  and decreasing in  $M$ , if  $f^* \leq M$ . Also, it satisfies the axioms of rational search stated in Žilinskas [22]. Then, employing (5.4) and (5.6), both methods choose the point that maximizes

$$U\left(s(y), \sqrt{\text{Var}(y)}\right), \quad y \in \mathcal{D}.$$

## 6. Search Strategies and Practical Questions

Practical features of our method have received little attention in this paper. Several questions arise concerning the choice of parameters in Algorithm 3.

1. What radial basis function  $\phi$  should be chosen, and what polynomial degree  $m$ ?
2. What is a good strategy for the choice of the target values  $f_n^*$ ?
3. Given a target value, how should the minimization of  $g_n$  in (3.7) (or the maximization of  $h_n$  in (3.8)) be carried out? Should we approximate the global optimum of  $g_n$  (or  $h_n$ ) or compute a (possibly non-global) local minimum?

The first problem has not been investigated thoroughly. Experiments using cubic and thin plate splines on a few test functions suggest that one cannot say in general that one of them is better than the other. Experiments have not been tried yet for the other types.

The choice of target values is crucial for the performance of the method. The interpretation of the two extremal cases has been noted in Section 3. Specifically, the choice  $f_n^* = \min_{y \in \mathcal{D}} s_n(y)$  means that we trust our model and assume that the minimizer of  $s_n$  is close to the global minimizer of  $f$ . In the other case, namely  $f_n^* = -\infty$ , we try to find a point in a region that has not been explored at all. It may be best to employ a mix between values of  $f_n^*$  that are suitable for convergence to a local minimizer and values that provide points in previously unexplored regions of the domain.

Research is going on for the third question. We prefer to maximize  $h_n$ , as this function is defined everywhere on  $\mathcal{D}$ . It might seem strange that we consider computing the global optimum of  $h_n$ , i.e. that a global optimization problem is replaced by another one. However, unlike  $f$ ,  $h_n$  can be evaluated quickly, and derivatives are available as well. Also we know roughly where the local maxima of  $h_n$  lie. Thus the maximization of  $h_n$  is much easier than the minimization of  $f$ . In addition, as the problem (GOP) is very difficult under our assumptions, it would take too long to compute a global minimizer accurately. Therefore, from a practical point of view, we are interested in an approximate solution of (GOP). So it should suffice to determine an approximation to the maximizer of (3.8). As far as the second option in question 3 is concerned, we have to find a way to choose starting points or search regions in order to ensure fast convergence, which is still an open problem.

Experiments show that large differences between function values can cause the interpolant to oscillate very strongly. Thus its minimal value can be much below

Table 1. Dixon-Szegö test functions and their dimension, the domain and the number of local and global minima

Function	Dimension	No. of local minima	No. of global minima	Domain
Branin	2	3	3	$[-5, 10] \times [0, 15]$
Goldstein-Price	2	4	1	$[-2, 2]^2$
Hartman 3	3	4	1	$[0, 1]^3$
Shekel 5	4	5	1	$[0, 10]^4$
Shekel 7	4	7	1	$[0, 10]^4$
Shekel 10	4	10	1	$[0, 10]^4$
Hartman 6	6	4	1	$[0, 1]^6$

the least calculated function value. We have found in numerical computations that these inefficiencies are reduced if large function values are replaced by the median of all available function values.

Some experiments were performed using the test functions proposed by Dixon and Szegö [2]. Table 1 gives the name of each function, the dimension, the domain and the number of local and global minima in that domain.

In all the cases, we use thin plate splines as interpolants, where the additional polynomials are linear. The initial points are chosen to be the corners of the domain. It seems that a different choice gives similar results, but this has not been studied yet. Also, for the computation of the interpolant  $s_n$  and the target value, function values that are larger than the median of all available values are replaced by the median.

The maximization of (3.8) is carried out using a version of the tunneling method (Levy and Montalvo [13]). When this method is at a local maximum, searches are started in each coordinate direction, using an auxiliary function. If all these searches fail to provide a point with a larger function value, the algorithm is stopped. The complexity of function evaluations as well as the number of local maxima of (3.8) increase with  $n$ , the number of iterations. This is a problem, when  $n$  is large, but in our experiments the time needed to solve the auxiliary problem normally was only a few seconds.

Finally, the target values  $f_n^*$  are determined as follows. The idea is to perform cycles of  $N + 1$  iterations for some  $N \in \mathbb{N}$ , where each cycle employs a range of target values, starting with a low one (global search), and ending with a value of  $f_n^*$  that is close to  $\min s_n(y)$  (local search). Then we go back to a global search, starting the cycle again. The results that we report have been obtained using the following

Table 2. Number of function evaluations for our method in comparison to DIRECT, DE, EGO and MCS with two different stopping criteria

Test function	Error < 1%				Error < 0.01%		
	RBF	DIRECT	DE	EGO	RBF	DIRECT	MCS
Branin	44	63	1190	28	64	195	41
Goldstein–Price	63	101	1018	32	76	191	81
Hartman 3	43	83	476	35	158	199	79
Shekel 5	76	103	6400	–	100	155	83
Shekel 7	76	97	6194	–	125	145	129
Shekel 10	51	97	6251	–	112	145	103
Hartman 6	112	213	7220	121	160	571	111

strategy. We choose the cycle length  $N = 5$ . Let the number of initial points be  $n_0$ , let the cycle start at  $n = \tilde{n}$ , and let the function values be ordered, i.e.  $f(x_1) \leq \dots \leq f(x_n)$ . If  $f(x_1) = f(x_n)$ , the interpolant is a constant function, because we pick  $m \geq 0$ , so the maximization of (3.8) is equivalent to the maximization of  $|1/\mu_n|$ , if  $f_n^* < f(x_1)$ . In this case, we choose  $f_n^* = -\infty$ . Otherwise, for  $\tilde{n} \leq n \leq \tilde{n} + N - 1$ , we set

$$f_n^* = \min_{y \in \mathcal{D}} s_n(y) - \left( \frac{N - n + \tilde{n}}{N} \right)^2 \left( f(x_{\sigma(n)}) - \min_{y \in \mathcal{D}} s_n(y) \right),$$

where  $\sigma(\tilde{n}) = \tilde{n}$  and  $\sigma(n) = \sigma(n-1) - \left\lfloor \frac{n - n_0}{N} \right\rfloor$ ,  $\tilde{n} + 1 \leq n \leq \tilde{n} + N - 1$ . When  $n = \tilde{n} + N$  we set  $f_n^* = \min_{y \in \mathcal{D}} s_n(y)$ . However, we have to be careful here since this choice is only admissible if  $x_1$  is not one of the global minimizers of  $s_n$ . Thus, we do not accept this choice, if  $(f(x_1) - \min_{y \in \mathcal{D}} s_n(y)) < 10^{-4}|f(x_1)|$ , provided  $f(x_1)$  is nonzero, or if  $f(x_1) - \min_{y \in \mathcal{D}} s_n(y) < 10^{-4}$ , if  $f(x_1) = 0$ . In these cases, we set  $f_n^* = \min_{y \in \mathcal{D}} s_n(y) - 10^{-2}|f(x_1)|$  and  $f_n^* = \min_{y \in \mathcal{D}} s_n(y) - 10^{-2}$ , respectively, to try to obtain a yet lower function value.

The algorithm is stopped when the relative error  $|f_{\text{best}} - f^*|/|f^*|$  becomes smaller than a fixed  $\epsilon$ , where  $f^*$  is the global optimum and  $f_{\text{best}}$  the current best function value. The optimal values of all test functions in Table 1 are nonzero, so this stopping criterion is valid.

Table 2 reports the number of function evaluations needed to achieve a relative error less than 1% and 0.01%. RBF denotes our method using the target value strategy described above. DIRECT (Jones, Perttunen and Stuckman [9]) and MCS (Multilevel Coordinate Search, Huyer and Neumaier [5]) are recent methods that, according to the results presented in those papers, are more efficient than most of their competitors on the Dixon-Szegö testbed. DE (Differential Evolution, Storn

and Price [18]) is an evolutionary method that operates only at the global level, which explains the large number of function evaluations. EGO (Efficient Global Optimization, Jones, Schonlau and Welch [10]) is the latest method known to us. Unfortunately, no tests are reported on the Shekel test functions. All the results from these methods are quoted from the papers mentioned above. For the DIRECT method, numbers of evaluations for both the 1% and the 0.01% stopping criterion are reported. For DE and EGO only results for the 1% criterion are available, whereas MCS only uses the 0.01% criterion. It should be noted that MCS uses a local search method at some stages of the algorithm, and in all the cases of Table 2 the first local minimum found is the global one.

## 7. Conclusions

Our global optimization method converges to the global minimizer of an arbitrary continuous function  $f$ , if we choose the sequence of target values carefully. If  $f$  is sufficiently smooth, there is even a suitable condition on this sequence that can be checked by the algorithm. However, it is unsatisfactory that the multiquadric and Gaussian cases are excluded from the statement of Theorem 7. It is believed that the convergence result is true also in these cases, although they are not covered by the analysis in [3].

Table 2 shows that the method is able to compete with other global optimization methods on the set of the Dixon-Szegö test functions. The test functions in this testbed, however, are of relatively low dimension, and the number of local and global minima is very small. Therefore, it is necessary to test the method on other sets of test functions and of course on real-world applications.

The maximization of the utility function  $h_n$  is still an unresolved problem. The tunneling method that has been used for the experiments requires too many parameters whose choice is not obvious. Jones, Schonlau and Welch [10] describe a branch-and-bound algorithm to solve an auxiliary problem in the EGO method. It is interesting to investigate how this approach can be used for our method. As mentioned above, another option is to carry out local searches from one or several starting points. The challenge here is to find suitable starting points that yield points with large values of  $h_n$ .

The relation to the P-algorithm is very interesting. It is hoped that the connections can be exploited further. In particular, the choice of the target values and the determination of the point of highest utility are common problems. Solutions to these problems may be developed that are useful for both methods.

## Appendix

The main convergence theorem, Theorem 7, is proved in this appendix. In order to establish it, some lemmas are needed on the behaviour of the coefficients  $\mu_n$ .

LEMMA 12. Let  $\phi$  be any of the radial basis functions in (2.2), and let  $m_0$  and  $m \geq m_0$  be chosen as in Section 2. Let  $\mathcal{D} \subset \mathbb{R}^d$  be compact, and let  $(x_n)_{n \in \mathbb{N}}$  be a convergent sequence in  $\mathcal{D}$  with pairwise different elements. Further, let  $(y_n)_{n \in \mathbb{N}}$  be a sequence in  $\mathcal{D}$  such that  $y_n \neq x_n$ ,  $n \in \mathbb{N}$ , and  $\lim_{n \rightarrow \infty} \|x_n - y_n\| = 0$ . Choose  $k$  points  $z_1, \dots, z_k \in \mathcal{D}$  that satisfy condition (2.11). Assume  $(x_n)$  converges to  $x^* \in \mathcal{D} \setminus \{z_1, \dots, z_k\}$ . For any  $y \in \mathcal{D} \setminus \{z_1, \dots, z_k, y_{n+1}\}$ , let  $\tilde{\ell}_y$  be the cardinal spline that interpolates the data  $(z_1, 0), \dots, (z_k, 0), (y_{n+1}, 0)$  and  $(y, 1)$ , and let  $\tilde{\mu}_n(y)$  be the coefficient of  $\tilde{\ell}_y$  that satisfies  $(-1)^{m_0+1} \tilde{\mu}_n(y) = \langle \tilde{\ell}_y, \ell_y \rangle$ . Then, for  $0 \leq \rho < 1$  in the linear case and  $0 \leq \rho < 2$  in the other cases,

$$\lim_{n \rightarrow \infty} (-1)^{m_0+1} \|y_{n+1} - x_{n+1}\|^\rho \tilde{\mu}_n(x_{n+1}) = \infty. \quad (\text{A.1})$$

*Proof.* Let  $A_n$  and  $A_n(x_{n+1})$  be the matrices of the form (2.10) for the interpolation points  $z_1, \dots, z_k, y_{n+1}$  and  $z_1, \dots, z_k, y_{n+1}, x_{n+1}$ , respectively. For sufficiently large  $n$ , neither  $x_{n+1}$  nor  $y_{n+1}$  is in the set  $\{z_1, \dots, z_k\}$ . Thus,  $A_n$  and  $A_n(x_{n+1})$  are nonsingular. Cramer's Rule implies

$$\tilde{\mu}_n(x_{n+1}) = \frac{\det A_n}{\det A_n(x_{n+1})}. \quad (\text{A.2})$$

Also, let  $A^*$  be the matrix of the form (2.10) for the interpolation points  $z_1, \dots, z_k, x^*$ . By the continuity of the determinant,

$$\lim_{n \rightarrow \infty} \det A_n = \det A^* \neq 0. \quad (\text{A.3})$$

In order to investigate the behaviour of  $\|y_{n+1} - x_{n+1}\|^{-\rho} \det A_n(x_{n+1})$ , we let

$$v(y) := \left( \phi(\|y - z_1\|), \dots, \phi(\|y - z_k\|) \right)^T, \quad y \in \mathcal{D},$$

and

$$p(y) := \left( p_1(y), \dots, p_{\hat{m}}(y) \right)^T, \quad y \in \mathcal{D},$$

where  $\hat{m} = \dim \Pi_m$  and  $p_1, \dots, p_{\hat{m}}$  are as in Section 2. Thus  $A_n(x_{n+1})$  is the matrix

$$\begin{pmatrix} \Phi & v(y_{n+1}) & v(x_{n+1}) & P \\ v(y_{n+1})^T & \phi(0) & \phi(\|y_{n+1} - x_{n+1}\|) & p(y_{n+1})^T \\ v(x_{n+1})^T & \phi(\|y_{n+1} - x_{n+1}\|) & \phi(0) & p(x_{n+1})^T \\ P^T & p(y_{n+1}) & p(x_{n+1}) & 0 \end{pmatrix}, \quad (\text{A.4})$$

where  $\Phi$  and  $P$  correspond to (2.3) and (2.9), respectively, if we set  $n = k$  and  $\{x_1, \dots, x_n\} = \{z_1, \dots, z_k\}$ . Now the rows

$$\begin{aligned} & \left( v(y_{n+1})^T \quad \phi(0) \quad \phi(\|y_{n+1} - x_{n+1}\|) \quad p(y_{n+1})^T \right) \\ \text{and} & \left( v(x_{n+1})^T \quad \phi(\|y_{n+1} - x_{n+1}\|) \quad \phi(0) \quad p(x_{n+1})^T \right) \end{aligned}$$

have the same limit as  $n \rightarrow \infty$ , so  $\det A_n(x_{n+1})$  tends to zero. Hence the properties (A.2) and (A.3) prove the assertion (A.1) for  $\rho = 0$ .

For  $\rho > 0$ , note that the value of the determinant of the matrix (A.4) does not change if we replace the second row by the difference between the second and third rows, and subsequently replace the second column by the difference between the second and third columns. Then the new second column of  $\det A_n(x_{n+1})$  becomes

$$\begin{pmatrix} v(y_{n+1}) - v(x_{n+1}) \\ 2[\phi(0) - \phi(\|y_{n+1} - x_{n+1}\|)] \\ \phi(\|y_{n+1} - x_{n+1}\|) - \phi(0) \\ p(y_{n+1}) - p(x_{n+1}) \end{pmatrix}, \tag{A.5}$$

and the new second row is its transpose. We have to divide the determinant by  $\|y_{n+1} - x_{n+1}\|^\rho$ , so we divide the second row and then the second column by  $\|y_{n+1} - x_{n+1}\|^{\rho/2}$ . Thus all components of (A.5) are multiplied by  $\|y_{n+1} - x_{n+1}\|^{\rho/2}$ , except the second one which is multiplied by  $\|y_{n+1} - x_{n+1}\|^\rho$ . Then the following remarks are helpful.

For each choice of  $\phi$  and  $j = 1, \dots, k$ , the function  $\phi(\|z_j - x\|)$ ,  $x \in \mathcal{D}$ , is Lipschitz continuous, so the components of  $v(y_{n+1}) - v(x_{n+1})$  satisfy

$$|\phi(\|z_j - y_{n+1}\|) - \phi(\|z_j - x_{n+1}\|)| \leq \text{const} \|x_{n+1} - y_{n+1}\|, \quad j = 1, \dots, k.$$

Thus for  $\rho < 2$  we have

$$\lim_{n \rightarrow \infty} \frac{1}{\|y_{n+1} - x_{n+1}\|^{\rho/2}} [\phi(\|z_j - y_{n+1}\|) - \phi(\|z_j - x_{n+1}\|)] = 0, \tag{A.6}$$

Similarly, for  $\rho < 2$ , the components of  $p(y_{n+1}) - p(x_{n+1})$  have the property

$$\lim_{n \rightarrow \infty} \frac{1}{\|y_{n+1} - x_{n+1}\|^{\rho/2}} [p_i(y_{n+1}) - p_i(x_{n+1})] = 0 \tag{A.7}$$

Finally, we deduce

$$\lim_{n \rightarrow \infty} \frac{1}{\|y_{n+1} - x_{n+1}\|^\rho} [\phi(\|y_{n+1} - x_{n+1}\|) - \phi(0)] = 0, \tag{A.8}$$

for  $\rho < 1$  in the linear case, for  $\rho < 2$  in the thin plate spline, multiquadric and Gaussian cases, and for  $\rho < 3$  in the cubic case. This is clear in the linear, thin plate spline and cubic cases. In the other two cases it follows from second order Taylor expansion of  $\phi$ , because  $\phi'(0) = 0$  and  $\phi''$  is bounded on  $\mathbb{R}_+$ . Thus (A.5) – (A.8) provide

$$\lim_{n \rightarrow \infty} \|y_{n+1} - x_{n+1}\|^{-\rho} \det A_n(x_{n+1}) = 0,$$

for the given values of  $\rho$ . Hence (A.2) and (A.3) imply (A.1). □

Now we obtain

LEMMA 13. Let  $\phi$ ,  $m_0$  and  $m$  be chosen as in Lemma 12, and let  $(x_n)_{n \in \mathbb{N}}$  be the sequence generated by Algorithm 3. Further, let  $0 \leq \rho < 1$  in the linear case and  $0 \leq \rho < 2$  in the other cases. Then, for every convergent subsequence  $(x_{n_k})_{k \in \mathbb{N}}$  of  $(x_n)$ ,

$$\lim_{k \rightarrow \infty} (-1)^{m_0+1} \Delta_{n_k}^\rho \mu_{n_k-1}(x_{n_k}) = \infty,$$

where  $\mu_n(\cdot)$  is defined in Section 3 and  $\Delta_{n_k}$  is expression (4.1) for  $n = n_k$ .

*Proof.* For  $n \geq 2$ , define  $j_n$  to be the natural number  $j$  that minimizes  $\|x_n - x_j\|$ ,  $j < n$ , so  $\Delta_n = \|x_n - x_{j_n}\|$ . Further, let  $(y_n)_{n \in \mathbb{N}}$  be the sequence

$$y_n := \begin{cases} x_2, & n = 1, \\ x_{j_n}, & n \geq 2. \end{cases}$$

Let  $(x_{n_k})_{k \in \mathbb{N}}$  be a subsequence of  $(x_n)_{n \in \mathbb{N}}$ , that converges to  $x^*$ , say. Convergence and the choice of  $(y_n)_{n \in \mathbb{N}}$  imply  $\lim_{k \rightarrow \infty} \|x_{n_k} - y_{n_k}\| = 0$ .

The initial step of Algorithm 3 provides a finite number of points that satisfy (2.11), so the initial interpolation matrix (2.10) is nonsingular. If one of these points is  $x^*$ , we pick  $x_{n_{k_0}}$  in a neighbourhood of  $x^*$  so that the interpolation matrix to  $x_{n_{k_0}}$  and the other initial points is also nonsingular. Therefore there exist points  $\hat{x}_1, \dots, \hat{x}_l$  in  $(x_n)_{n \in \mathbb{N}}$  such that their interpolation matrix is nonsingular, and  $x^* \notin \{\hat{x}_1, \dots, \hat{x}_l\}$ .

For sufficiently large  $k \in \mathbb{N}$ , such that  $y_{n_k} \notin \{\hat{x}_1, \dots, \hat{x}_l\}$ , and for any  $y \in \mathcal{D} \setminus \{x_1, \dots, x_{n_k-1}\}$ , we let  $\hat{\ell}_k(y, \cdot)$  be the radial function that interpolates  $(\hat{x}_1, 0), \dots, (\hat{x}_l, 0), (y_{n_k}, 0)$  and  $(y, 1)$ , and we let  $\ell_{n_k-1}(y, \cdot)$  be the interpolant to  $(x_1, 0), \dots, (x_{n_k-1}, 0)$  and  $(y, 1)$ . Because  $\ell_{n_k-1}(y, \cdot)$  interpolates  $(y_{n_k}, 0)$  and  $(\hat{x}_i, 0)$ ,  $i = 1, \dots, l$ , for sufficiently large  $k$ , (3.6) and Theorem 1 imply the inequality

$$\begin{aligned} (-1)^{m_0+1} \hat{\mu}_k(y) &= \langle \hat{\ell}_k(y, \cdot), \hat{\ell}_k(y, \cdot) \rangle \\ &\leq \langle \ell_{n_k-1}(y, \cdot), \ell_{n_k-1}(y, \cdot) \rangle = (-1)^{m_0+1} \mu_{n_k-1}(y) \end{aligned} \quad (\text{A.9})$$

for the coefficients  $\hat{\mu}_k$  and  $\mu_{n_k-1}$ .

Now we apply Lemma 12 in the case when  $\{z_1, \dots, z_k\}$  is the set  $\{\hat{x}_1, \dots, \hat{x}_l\}$  and  $n = n_k - 1$ . It follows that

$$\lim_{k \rightarrow \infty} (-1)^{m_0+1} \Delta_{n_k}^\rho \hat{\mu}_k(x_{n_k}) = \lim_{k \rightarrow \infty} (-1)^{m_0+1} \|x_{n_k} - y_{n_k}\|^\rho \hat{\mu}_k(x_{n_k}) = \infty,$$

with the choice of  $\rho$  stated in Lemma 12. Thus, setting  $y = x_{n_k}$  in (A.9), we obtain that  $(-1)^{m_0+1} \Delta_{n_k}^\rho \mu_{n_k-1}(x_{n_k})$  tends to infinity as  $k \rightarrow \infty$ .  $\square$

Finally we show, using Proposition 6, that the coefficients  $\mu_n(y)$  are uniformly bounded, if  $y$  is bounded away from the points that are generated by the algorithm.

LEMMA 14. Let  $\phi(r) = r$ ,  $\phi(r) = r^2 \log r$  or  $\phi(r) = r^3$ . Further, choose the integer  $m$  such that  $0 \leq m \leq d$  in the linear case,  $1 \leq m \leq d + 1$  in the thin plate



spline case and  $1 \leq m \leq d + 2$  in the cubic case. Let  $(x_n)_{n \in \mathbb{N}}$  be the sequence generated by Algorithm 3, and let  $n_0$  be the number of points chosen in the initial step. Assume that there exist  $y_0 \in \mathcal{D}$  and a neighbourhood  $N_\delta := \{x \in \mathbb{R}^d : \|x - y_0\| < \delta\}$ ,  $\delta > 0$ , that does not contain any point of the sequence. Then there exists  $K > 0$ , that depends only on  $y_0$  and  $\delta$ , such that

$$(-1)^{m_0+1} \mu_n(y_0) \leq K \quad \forall n \geq n_0.$$

*Proof.* For any  $n \geq n_0$ , let  $\ell_n$  be the radial function that is defined by  $\ell_n(x_i) = 0$ ,  $i = 1, \dots, n$ , and  $\ell_n(y_0) = 1$ . There exists a compactly supported infinitely differentiable function  $F$  that takes the value 1 at  $y_0$  and 0 on  $\mathbb{R}^d \setminus N_\delta$ . It follows from Proposition 6 that  $F \in \mathcal{N}_{\phi,m}$ .  $\ell_n$  interpolates  $F$  at  $x_1, \dots, x_n$  and  $y_0$ . Therefore, there is a positive number  $K$ , depending on  $y_0$  and  $\delta$ , such that

$$(-1)^{m_0+1} \mu_n(y_0) = \langle \ell_n, \ell_n \rangle \leq K, \quad n \geq n_0. \quad \square$$

Now we are ready to prove Theorem 7.

*Proof of Theorem 7.* Assume there is  $y_0 \in \mathcal{D}$  and an open neighbourhood  $U = \{x \in \mathbb{R}^d : \|x - y_0\| < \delta\}$ ,  $\delta > 0$ , that does not contain an interpolation point. The iteration step of Algorithm 3 gives

$$g_n(x_{n+1}) \leq g_n(y_0), \quad n \geq n_0,$$

where  $n_0$  is the number of points chosen in the initial step of the algorithm.

By assumption (4.2), there is a subsequence  $(n_k)_{k \in \mathbb{N}}$  of the natural numbers such that

$$\min_{y \in \mathcal{D}} s_{n_k-1}(y) - f_{n_k-1}^* > \tau \Delta_{n_k-1}^{\rho/2} \|s_{n_k-1}\|_\infty \geq 0, \quad k \in \mathbb{N}, \quad (\text{A.10})$$

with  $\tau > 0$ ,  $\Delta_{n_k-1}$  being the expression (4.1) for  $n = n_k - 1$ ,  $0 \leq \rho < 1$  in the linear and  $0 \leq \rho < 2$  in the thin plate spline and cubic cases. The sequence  $(x_{n_k})_{k \in \mathbb{N}}$  is a sequence in a compact set, thus it contains a convergent subsequence. Therefore, without loss of generality, we assume that  $(x_{n_k})_{k \in \mathbb{N}}$  itself converges.

For all  $k \in \mathbb{N}$ ,  $x_{n_k}$  is the minimizer of  $g_{n_k-1}(x)$ . Thus, if  $f_{n_k-1}^* > -\infty$ ,

$$\begin{aligned} & (-1)^{m_0+1} \mu_{n_k-1}(x_{n_k}) [s_{n_k-1}(x_{n_k}) - f_{n_k-1}^*]^2 \\ & \leq (-1)^{m_0+1} \mu_{n_k-1}(y_0) [s_{n_k-1}(y_0) - f_{n_k-1}^*]^2. \end{aligned} \quad (\text{A.11})$$

If  $\|s_{n_k-1}\|_\infty > 0$ , this inequality, condition (A.10) and the definition of  $\|\cdot\|_\infty$  provide

$$\begin{aligned} (-1)^{m_0+1} \mu_{n_k-1}(x_{n_k}) &\leq (-1)^{m_0+1} \mu_{n_k-1}(y_0) \left[ \frac{s_{n_k-1}(y_0) - f_{n_k-1}^*}{s_{n_k-1}(x_{n_k}) - f_{n_k-1}^*} \right]^2 \\ &\leq (-1)^{m_0+1} \mu_{n_k-1}(y_0) \left[ 1 + \frac{|s_{n_k-1}(y_0) - s_{n_k-1}(x_{n_k})|}{s_{n_k-1}(x_{n_k}) - f_{n_k-1}^*} \right]^2 \\ &\leq (-1)^{m_0+1} \mu_{n_k-1}(y_0) \left[ 1 + \frac{1}{\tau \Delta_{n_k}^{\rho/2}} \frac{|s_{n_k-1}(y_0) - s_{n_k-1}(x_{n_k})|}{\|s_{n_k-1}\|_\infty} \right]^2 \\ &\leq (-1)^{m_0+1} \mu_{n_k-1}(y_0) \left[ 1 + \frac{2}{\tau \Delta_{n_k}^{\rho/2}} \right]^2. \end{aligned}$$

If  $\|s_{n_k-1}\|_\infty = 0$ ,  $s_{n_k-1}(y) - f_{n_k-1}^*$  is a positive number independent of  $y$ , thus (A.11) gives

$$\begin{aligned} (-1)^{m_0+1} \mu_{n_k-1}(x_{n_k}) &\leq (-1)^{m_0+1} \mu_{n_k-1}(y_0) \\ &\leq (-1)^{m_0+1} \mu_{n_k-1}(y_0) \left[ 1 + \frac{2}{\tau \Delta_{n_k}^{\rho/2}} \right]^2, \end{aligned}$$

for any positive  $\tau$ , as before. Remark 2 shows that this inequality holds also in the case  $f_{n_k-1}^* = -\infty$ . Multiplying both sides by  $\Delta_{n_k}^\rho$  yields

$$\Delta_{n_k}^\rho (-1)^{m_0+1} \mu_{n_k-1}(x_{n_k}) \leq (-1)^{m_0+1} \mu_{n_k-1}(y_0) \left[ \Delta_{n_k}^{\rho/2} + \frac{2}{\tau} \right]^2. \quad (\text{A.12})$$

By Lemma 13, the left-hand side of (A.12) tends to infinity as  $k$  tends to infinity. However, Lemma 14 states that  $(-1)^{m_0+1} \mu_n(y_0)$  is bounded above by a constant that does not depend on  $n$ . Thus the right-hand side of (A.12) is bounded by a constant that is independent of  $k$  which contradicts (A.12). Therefore there is a point in the sequence that is an element of  $U$ . This implies that in each neighbourhood of an arbitrary  $y_0 \in \mathcal{D}$  there are infinitely many elements of  $(x_n)_{n \in \mathbb{N}}$ , so the sequence is dense in  $\mathcal{D}$ .  $\square$

### Acknowledgements

I am very grateful to my supervisor, Prof. M.J.D. Powell, for his constant help and his guidance. Also, I would like to thank two anonymous referees for their useful suggestions.

## References

1. Alotto, P., Caiti, A., Molinari, G. and Repetto, M. (1996), A Multiquadrics-based Algorithm for the Acceleration of Simulated Annealing Optimization Procedures, *IEEE Transactions on Magnetics* 32(3): 1198–1201.
2. Dixon, L.C.W. and Szegö, G. (1978), The Global Optimization Problem: An Introduction, in: Dixon, L. and Szegö, G. (eds), *Towards Global Optimization 2*, North-Holland, Amsterdam, pp. 1–15.
3. Gutmann, H.-M. On the semi-norm of radial basis function interpolants, Report DAMTP2000/NA04, University of Cambridge.
4. Horst, R. and Pardalos, P.M. (1994), *Handbook of Global Optimization*, Kluwer, Dordrecht.
5. Huyer, W. and Neumaier, A. (1999), Global optimization by multilevel coordinate search, *Journal of Global Optimization* 14(4): 331–355.
6. Ishikawa, T. and Matsunami, M. (1997), An Optimization Method Based on Radial Basis Functions, *IEEE Transactions on Magnetics* 33(2): 1868–1871.
7. Ishikawa, T., Tsukui, Y. and Matsunami, M. (1999), A Combined Method for the Global Optimization Using Radial Basis Function and Deterministic Approach, *IEEE Transactions on Magnetics* 35(3): 1730–1733.
8. Jones, D.R. (1996), Global optimization with response surfaces, presented at the Fifth SIAM Conference on Optimization, Victoria, Canada.
9. Jones, D.R., Perttunen, C. and Stuckman, B.E. (1993), Lipschitz Optimization Without the Lipschitz Constant, *Journal of Optimization Theory and Applications* 78(1): 157–181.
10. Jones, D.R., Schonlau, M. and Welch, W.J. (1998), Efficient Global Optimization of Expensive Black-Box Functions, *Journal of Global Optimization* 13(4): 455–492.
11. Kushner, H.J. (1962), A Versatile Model of a Function of Unknown and Time Varying Form, *Journal of Mathematical Analysis and Applications* 5: 150–167.
12. Kushner, H.J. (1964), A New Method of Locating the Maximum Point of an Arbitrary Multipeak Curve in the Presence of Noise, *Journal of Basic Engineering* 86: 97–106.
13. Levy, A.V. and Montalvo, A. (1985), The Tunneling Algorithm for the Global Minimization of Functions, *SIAM Journal on Scientific and Statistical Computing* 6(1): 15–29.
14. Powell, M.J.D. (1981), *Approximation Theory and Methods*, Cambridge University Press.
15. Powell, M.J.D. (1992), The Theory of Radial Basis Function Approximation in 1990, in: Light, W. (ed.), *Advances in Numerical Analysis, Volume 2: Wavelets, Subdivision Algorithms and Radial Basis Functions*, Oxford University Press, pp. 105–210.
16. Powell, M.J.D. (1999), Recent research at Cambridge on radial basis functions, in: M. Müller, M. Buhmann, D. Mache and M. Felten (eds), *New Developments in Approximation Theory, International Series of Numerical Mathematics, Vol. 132*, Birkhauser Verlag, Basel, pp. 215–232.
17. Schaback, R. (1993), Comparison of radial basis function interpolants, in: K. Jetter and F. Utreras (eds), *Multivariate Approximations: From CAGD to Wavelets*, World Scientific, Singapore, pp. 293–305.
18. Storn, R. and Price, K. (1997), Differential Evolution - A Simple and Efficient Heuristic for Global Optimization over Continuous Spaces, *Journal of Global Optimization* 11(4): 341–359.
19. Törn, A. and Žilinskas, A. (1987), *Global Optimization*, Springer, Berlin.
20. Whitney, H. (1934), Analytic extension of differentiable functions defined in closed sets, *Transactions of the American Mathematical Society* 36: 63–89.
21. Žilinskas, A. (1982), Axiomatic Approach to Statistical Models and their Use in Multimodal Optimization Theory, *Mathematical Programming* 22(1): 104–116.
22. Žilinskas, A. (1985), Axiomatic Characterization of a Global Optimization Algorithm and Investigation of its Search Strategy, *Operations Research Letters* 4(1): 35–39.