

On estimation of surrogate models for multivariate computer experiments[‡]

Benedikt Bauer¹, Felix Heimrich², Michael Kohler¹ and Adam Krzyżak^{3,§}

¹ *Fachbereich Mathematik, Technische Universität Darmstadt, Schlossgartenstr. 7, 64289 Darmstadt, Germany, email: bbauer@mathematik.tu-darmstadt.de, kohler@mathematik.tu-darmstadt.de*

² *Fachbereich Maschinenbau, Technische Universität Darmstadt, Otto-Bernd-Str. 2, 64287 Darmstadt, Germany, email: heimrich@dik.tu-darmstadt.de*

³ *Department of Computer Science and Software Engineering, Concordia University, 1455 De Maisonneuve Blvd. West, Montreal, Quebec, Canada H3G 1M8, email: krzyzak@cs.concordia.ca*

Supplementary material for the referees

Proof of Lemma 3.

In the proof we will use Proposition 3.8 in Mhaskar (1993), which we reformulate here (in a slightly different form) as Lemma 6.

Lemma 6. *Let $K \subseteq \mathbb{R}^d$ be a polytope bounded by hyperplanes $v_j \cdot x + w_j \geq 0$ ($j = 1, \dots, L$), where $v_1, \dots, v_L \in \mathbb{R}^d$ and $w_1, \dots, w_L \in \mathbb{R}$. For $\delta > 0$ set*

$$K_\delta^0 := \left\{ x \in \mathbb{R}^d : v_j \cdot x + w_j \geq \delta \text{ for all } j \in \{1, \dots, L\} \right\}$$

and

$$K_\delta^c := \left\{ x \in \mathbb{R}^d : v_j \cdot x + w_j \leq -\delta \text{ for some } j \in \{1, \dots, L\} \right\}.$$

Let $\sigma : \mathbb{R} \rightarrow [0, 1]$ be a squashing function. Let $\varepsilon, \delta \in (0, 1]$ be arbitrary. Then there exists a neural network of the form

$$f(x) = \sigma \left(\sum_{j=1}^L b_j \cdot \sigma \left(\sum_{k=1}^d a_{j,k} \cdot x^{(k)} + a_{j,0} \right) + b_0 \right)$$

satisfying

$$\begin{aligned} |f(x)| &\leq 1 \text{ for } x \in \mathbb{R}^d, \\ |f(x) - 1| &\leq \varepsilon \text{ for } x \in K_\delta^0, \end{aligned}$$

[‡]Running title: *On estimation of surrogate models*

The online version of this article contains supplementary material.

[§]Corresponding author. Tel: +1-514-848-2424, ext. 3007, Fax: +1-514-848-2830.

$$|f(x)| \leq \varepsilon \text{ for } x \in K_\delta^c. \quad (33)$$

In case that the squashing function satisfies

$$|\sigma(y) - 1| \leq \frac{1}{y} \text{ if } y > 0 \quad \text{and} \quad |\sigma(y)| \leq \frac{1}{|y|} \text{ if } y < 0,$$

the weights above can be chosen such that

$$\begin{aligned} |b_j| &\leq \frac{4L}{\varepsilon} \text{ for all } j = 0, \dots, L \\ |a_{j,k}| &\leq \frac{4L}{\delta} \cdot \max\{\|v_1\|_\infty, |w_1|, \dots, \|v_L\|_\infty, |w_L|\} \text{ for all } j = 1, \dots, L, k = 0, \dots, d. \end{aligned}$$

Proof. Follows from the proof of Proposition 3.8 in Mhaskar (1993). \square

Proof of Lemma 3. We partition $[-a - \frac{2a}{M}, a]^d$ into $(M+1)^d$ equivolume cubes of side length $2a/M$. Approximating m by a piecewise constant approximant with respect to this partition yields (since m is (p, C) -smooth with $p \leq 1$) a function S satisfying

$$\|S - m\|_{\infty, [-a, a]^d} \leq \sqrt{d} \cdot C \cdot \left(\frac{2a}{M}\right)^p. \quad (34)$$

If we choose S suitably, it can be expressed in the form

$$S(x) = m(x_{(1, \dots, 1)}) + \sum_{i \in \{1, \dots, M+1\}^d \setminus \{(1, \dots, 1)\}} d_i \cdot \prod_{j=1}^d \left(x^{(j)} - x_i^{(j)}\right)_+^0,$$

where x_i are the corners of the cubes forming the above partition (indexed in ascending order per component), $0^0 := 0$, $x_+ := \max\{x, 0\}$, and d_i for $i = (i_1, \dots, i_d)$ as above are constants satisfying

$$d_i = \sum_{J \subseteq \{1, \dots, d\} \setminus \{k: i_k=1\}} (-1)^{|J|} \cdot m(x_{i-J}), \quad (35)$$

where $i - J$ symbolizes the index i with i_j replaced by $i_j - 1$ for all $j \in J$. Since for a fixed set with $n > 0$ elements the number of subsets with even and uneven cardinality is 2^{n-1} , respectively, and the corners used in the above expression have a distance of at most $\sqrt{d} \cdot \frac{2a}{M}$, we can conclude from the (p, C) -smoothness of m

$$|d_i| \leq 2^{d - |\{k: i_k=1\}| - 1} \cdot C \cdot d^{\frac{p}{2}} \cdot \left(\frac{2a}{M}\right)^p \leq c_{14} \cdot \left(\frac{2a}{M}\right)^p. \quad (36)$$

Let K_i be the polytope defined by $x^{(j)} - x_i^{(j)} \geq 0$ ($j = 1, \dots, d$). Set $\varepsilon = (M+1)^{-d}$, $\delta = a \cdot \eta / (2 \cdot d \cdot M)$ and apply Lemma 6 for each K_j (i.e., with $L = d$, $v_j = \mathbf{e}_j$ and $w_j = -x_i^{(j)}$, where \mathbf{e}_j denotes the j -th unit vector) to obtain $f_i(x)$ satisfying (33) with K_i instead of K . Let

$$P(x) = m(x_{(1, \dots, 1)}) + \sum_{i \in \{1, \dots, M+1\}^d \setminus \{(1, \dots, 1)\}} d_i \cdot f_i(x).$$

Then we can conclude from (33) and (36)

$$\begin{aligned}
|P(x) - S(x)| &\leq \sum_{i \in \{1, \dots, M+1\}^d \setminus \{(1, \dots, 1)\}} |d_i| \cdot \left| f_i(x) - \prod_{j=1}^d (x^{(j)} - x_i^{(j)})_+^0 \right| \\
&\leq \sum_{i \in \{1, \dots, M+1\}^d \setminus \{(1, \dots, 1)\}} |d_i| \cdot (M+1)^{-d} \\
&\leq c_{14} \cdot \left(\frac{2a}{M} \right)^p
\end{aligned} \tag{37}$$

for all $x \in [-a, a]^d$ which are not contained in

$$\bigcup_{j=1, \dots, d} \bigcup_{i \in \{1, \dots, M+1\}^d} \left\{ x \in \mathbb{R}^d \quad : \quad |x^{(j)} - x_i^{(j)}| < a \cdot \eta / (2 \cdot d \cdot M) \right\}. \tag{38}$$

By shifting the positions of the x_i in the j th component slightly to the right (in the sense of increasing values) we can construct

$$\left\lfloor \frac{2a/M}{2\delta} \right\rfloor = \left\lfloor \frac{2a}{M} \cdot \frac{2 \cdot d \cdot M}{2 \cdot a \cdot \eta} \right\rfloor = \left\lfloor \frac{2 \cdot d}{\eta} \right\rfloor \geq d/\eta$$

different versions of P , that still satisfy (34) and (37) for all $x \in [-a, a]^d$, and corresponding disjoint versions of

$$\bigcup_{i \in \{1, \dots, M+1\}^d} \left\{ x \in \mathbb{R}^d \quad : \quad |x^{(j)} - x_i^{(j)}| < a \cdot \eta / (2 \cdot d \cdot M) \right\},$$

and since the sum of the ν -measures of these sets is less than or equal to one, at least one of them must have measure less than or equal to η/d . Consequently we can shift the x_i such that (38) has ν -measure less than or equal to η . This together with (34) and (37) implies the first assertion of the lemma, because $P(x)$ complies with the structure of the postulated neural network $t(x)$.

In case that σ satisfies the conditions specified in the second part of the lemma, Lemma 6 allows to bound the coefficients of the neural network $t(x) := P(x)$ respecting the values of the parameters we used during the application of this lemma above. This leads to

$$|a_{i,j,k}| \leq \frac{4 \cdot d \cdot 2 \cdot d \cdot M}{a \cdot \eta} \cdot \max \left\{ 1, a + \frac{2a}{M} \right\} \leq 8 \cdot d^2 \cdot \frac{M}{\eta} \cdot \max \left\{ \frac{1}{a}, 3 \right\}$$

for all $i \in \{1, \dots, (M+1)^d\}$, $j \in \{1, \dots, d\}$, $k \in \{0, \dots, d\}$ and

$$|b_{i,j}| \leq 4 \cdot d \cdot (M+1)^d$$

for all $i \in \{1, \dots, (M+1)^d\}$, $j \in \{0, \dots, d\}$. Furthermore, the definition of P and (35) imply

$$|d_i| \leq 2^d \cdot \|m\|_\infty$$

for all $i \in \{0, \dots, (M+1)^d\}$, which leads to the second assertion of the lemma. \square

In order to prove Lemma 5, we introduce the following technical result.

Lemma 7. Let $l \in \mathbb{N}_0$ and let $\sigma_r : \mathbb{R} \rightarrow \mathbb{R}$ for $r = 1, \dots, l+1$ be Lipschitz continuous functions with Lipschitz constant $L \geq 1$, which satisfy

$$|\sigma_r(x)| \leq L \cdot \max\{|x|, 1\} \quad (x \in \mathbb{R}). \quad (39)$$

Let $K_0 = d$, $K_r \in \mathbb{N}$ for $r \in \{1, \dots, l\}$ and $K_{l+1} = 1$. For $r \in \{1, \dots, l+1\}$ and $i \in \{1, \dots, K_r\}$ define recursively

$$f_i^{(r)}(x) = \sigma_r \left(\sum_{j=1}^{K_{r-1}} c_{i,j}^{(r-1)} \cdot f_j^{(r-1)}(x) + c_{i,0}^{(r-1)} \right)$$

and

$$\bar{f}_i^{(r)}(x) = \sigma_r \left(\sum_{j=1}^{K_{r-1}} \bar{c}_{i,j}^{(r-1)} \cdot \bar{f}_j^{(r-1)}(x) + \bar{c}_{i,0}^{(r-1)} \right),$$

where $c_{i,0}^{(r-1)}, \bar{c}_{i,0}^{(r-1)}, \dots, c_{i,K_{r-1}}^{(r-1)}, \bar{c}_{i,K_{r-1}}^{(r-1)} \in \mathbb{R}$, and $f_j^{(0)}(x) = \bar{f}_j^{(0)}(x) = x^{(j)}$. Furthermore, set

$$\bar{C} = \max_{\substack{r=0, \dots, l, \\ j=1, \dots, K_r}} \max_{i=1, \dots, K_{r+1}} \left\{ |c_{i,j}^{(r)}|, |\bar{c}_{i,j}^{(r)}|, 1 \right\}.$$

Then

$$\begin{aligned} & |f_1^{(l+1)}(x) - \bar{f}_1^{(l+1)}(x)| \\ & \leq (l+1) \cdot L^{l+1} \cdot \prod_{r=0}^l (K_r + 1) \cdot \bar{C}^l \cdot \max\{\|x\|_\infty, 1\} \cdot \max_{\substack{r=0, \dots, l, \\ j=0, \dots, K_r}} \left| c_{i,j}^{(r)} - \bar{c}_{i,j}^{(r)} \right| \end{aligned}$$

for any $x \in \mathbb{R}^d$.

Proof. At first, we notice that (39) implies

$$\left| \bar{f}_i^{(r)}(x) \right| \leq L \cdot (K_{r-1} + 1) \cdot \bar{C} \cdot \max_{j=1, \dots, K_{r-1}} \left\{ \left| \bar{f}_j^{(r-1)}(x) \right|, 1 \right\}$$

for $r = 1, \dots, l$ and $i = 1, \dots, K_r$, from which we can conclude

$$\left| \bar{f}_i^{(r)}(x) \right| \leq L^r \cdot \prod_{\bar{r}=1}^r (K_{\bar{r}-1} + 1) \cdot \bar{C}^r \cdot \max\{\|x\|_\infty, 1\}. \quad (40)$$

Using the Lipschitz continuity of σ_r and the triangle inequality in combination with (40) we get

$$\begin{aligned} & \left| f_i^{(r)}(x) - \bar{f}_i^{(r)}(x) \right| \\ & \leq L \cdot \left| \sum_{j=1}^{K_{r-1}} c_{i,j}^{(r-1)} \cdot f_j^{(r-1)}(x) + c_{i,0}^{(r-1)} - \sum_{j=1}^{K_{r-1}} \bar{c}_{i,j}^{(r-1)} \cdot \bar{f}_j^{(r-1)}(x) - \bar{c}_{i,0}^{(r-1)} \right| \end{aligned}$$

$$\begin{aligned}
&\leq L \cdot \left(\sum_{j=1}^{K_{r-1}} |c_{i,j}^{(r-1)}| \cdot \left| f_j^{(r-1)}(x) - \bar{f}_j^{(r-1)}(x) \right| \right. \\
&\quad \left. + \sum_{j=1}^{K_{r-1}} \left| c_{i,j}^{(r-1)} - \bar{c}_{i,j}^{(r-1)} \right| \cdot \left| \bar{f}_j^{(r-1)}(x) \right| + \left| c_{i,0}^{(r-1)} - \bar{c}_{i,0}^{(r-1)} \right| \right) \\
&\leq L \cdot K_{r-1} \cdot \bar{C} \cdot \max_{j=1, \dots, K_{r-1}} \left| f_j^{(r-1)}(x) - \bar{f}_j^{(r-1)}(x) \right| \\
&\quad + L \cdot (K_{r-1} + 1) \cdot L^{r-1} \cdot \prod_{\tilde{r}=1}^{r-1} (K_{\tilde{r}-1} + 1) \cdot \bar{C}^{r-1} \cdot \max \{ \|x\|_\infty, 1 \} \\
&\quad \cdot \max_{j=0, \dots, K_{r-1}} |c_{i,j}^{(r-1)} - \bar{c}_{i,j}^{(r-1)}| \\
&= L \cdot K_{r-1} \cdot \bar{C} \cdot \max_{j=1, \dots, K_{r-1}} \left| f_j^{(r-1)}(x) - \bar{f}_j^{(r-1)}(x) \right| \\
&\quad + L^r \cdot \prod_{\tilde{r}=1}^r (K_{\tilde{r}-1} + 1) \cdot \bar{C}^{r-1} \cdot \max \{ \|x\|_\infty, 1 \} \cdot \max_{j=0, \dots, K_{r-1}} |c_{i,j}^{(r-1)} - \bar{c}_{i,j}^{(r-1)}|
\end{aligned}$$

for all $r = 1, \dots, l+1$. Now we start with the above inequality for $r = l+1$ and plug it in repeatedly for decreasing r in the expression $\left| f_j^{(r-1)}(x) - \bar{f}_j^{(r-1)}(x) \right|$ on the right-hand side of the inequality. Finally, the summand containing $\left| f_j^{(0)}(x) - \bar{f}_j^{(0)}(x) \right|$ vanishes in the case of $r = 1$, which implies the assertion. \square

Proof of Lemma 5. At first, we notice the space $\mathcal{F}_n = \mathcal{H}^{(l)}$ (with $l > 0$) can be expressed as

$$\mathcal{H}^{(l)} = \left\{ h : \mathbb{R}^d \rightarrow \mathbb{R} : h(x) = \sum_{k=1}^K \sigma_{id}(g_k(\sigma_{id}(f_{1,k}(x)), \dots, \sigma_{id}(f_{d^*,k}(x)))) \quad (x \in \mathbb{R}^d) \right. \\
\left. \text{for some } g_k \in \mathcal{F}_{M_n, d^*, d^*, \alpha, \beta, \gamma}^{(\text{neural networks})} \text{ and } f_{j,k} \in \mathcal{H}^{(l-1)} \right\},$$

where $\sigma_{id} : \mathbb{R} \rightarrow \mathbb{R}$ is the identity $\sigma_{id}(x) = x$ for all $x \in \mathbb{R}$. Furthermore, all $g \in \mathcal{F}_{M_n, d^*, d^*, \alpha, \beta, \gamma}^{(\text{neural networks})}$ can be written as

$$\begin{aligned}
g(x) &= \sum_{i=1}^{(M_n+1)^{d^*}} d_i \cdot \sigma \left(\sum_{j=1}^{d^*} b_{i,j} \cdot \sigma \left(\sum_{m=1}^{d^*} a_{i,j,m} \cdot x^{(m)} + a_{i,j,0} \right) + b_{i,0} \right) + d_0 \\
&= \sum_{i=1}^{(M_n+1)^{d^*}} d_i \cdot \sigma \left(\sum_{\substack{j=1, \dots, d^* \\ \bar{i}=1, \dots, (M_n+1)^{d^*}}} b_{i, \bar{i}, j} \cdot \sigma \left(\sum_{m=1}^{d^*} a_{\bar{i}, j, m} \cdot x^{(m)} + a_{\bar{i}, j, 0} \right) + b_{i, \bar{i}, 0} \right) + d_0,
\end{aligned}$$

where the new coefficients are defined by

$$b_{i,\bar{i},j} := \begin{cases} b_{i,j} & \text{if } \bar{i} = i \\ 0 & \text{otherwise} \end{cases}$$

for all $i, \bar{i} \in \{1, \dots, (M_n + 1)^{d^*}\}$ and $j \in \{0, \dots, d^*\}$ (which works analogously for $h \in \mathcal{H}^{(0)}$). Respecting the above representations, all the functions $\sigma_{id}(h) = h$ for $h \in \mathcal{H}^{(l)}$ comply with the structure of the functions $f_1^{(l+1)}$ in Lemma 7, if we use the following specifications of the parameters in that lemma: The Lipschitz constant L is chosen as the maximum of the Lipschitz constants of σ_{id} (which is obviously 1) and the squashing function σ from Theorem 3. Thus, the property (39) is satisfied due to $\|\sigma\|_\infty \leq 1$, $L \geq 1$, and $|\sigma_{id}(x)| = |x|$. The parameter l in Lemma 7 is $4l + 2$ (regarding the l in $\mathcal{H}^{(l)}$ above) and the parameters K_r with $r = 0, \dots, l$ take repeatedly the values $\tilde{d}, d^* \cdot (M_n + 1)^{d^*}, (M_n + 1)^{d^*}, K$ one after another, where \tilde{d} is equal to d^* except for K_0 , where it is d . Since all the coefficients $c_{i,j}^{(r)}$ with $r = 0, \dots, l$, $i = 1, \dots, K_{r+1}$, $j = 1, \dots, K_r$ (using $K_{l+1} = 1$ again) are 0, 1, or one of the $a_{i,j,m}, b_{i,j}, d_i$ in the definition of $\mathcal{F}_{M_n, d^*, d, \alpha, \beta, \gamma}^{(\text{neural networks})}$, we can use $\bar{C} = \max\{\alpha, \beta, \gamma\}$ for n sufficiently large.

Let h and \bar{h} be functions in \mathcal{F}_n . Since they comply with the structure of the functions in Lemma 7 according to the above argumentation, we can conclude

$$\begin{aligned} & \|h - \bar{h}\|_{\infty, [-a_n, a_n]^d} \\ & \leq (4l + 3) \cdot L^{4l+3} \cdot \left(d^* \cdot (M_n + 1)^{d^*} + 1\right)^{4l+3} \cdot \max\{\alpha, \beta, \gamma\}^{4l+2} \\ & \quad \cdot \max\{a_n, 1\} \cdot \max_{\substack{r=0, \dots, \tilde{l}, \\ j=0, \dots, K_r}} \max_{\substack{i=1, \dots, K_{r+1}, \\ j=0, \dots, K_r}} \left|c_{i,j}^{(r)} - \bar{c}_{i,j}^{(r)}\right| \\ & \leq a_n \cdot n^{c_{11}} \cdot \max_{\substack{r=0, \dots, \tilde{l}, \\ j=0, \dots, K_r}} \max_{\substack{i=1, \dots, K_{r+1}, \\ j=0, \dots, K_r}} \left|c_{i,j}^{(r)} - \bar{c}_{i,j}^{(r)}\right| \end{aligned}$$

for n sufficiently large and an adequately chosen $c_{11} > 0$. Thus, if we consider an arbitrary $h \in \mathcal{H}^{(l)}$, it suffices to choose the coefficients $\bar{c}_{i,j}^{(r)}$ of a function $\bar{h} \in \mathcal{H}^{(l)}$ such that

$$\left|c_{i,j}^{(r)} - \bar{c}_{i,j}^{(r)}\right| \leq \frac{\varepsilon_n}{a_n \cdot n^{c_{11}}} \quad (41)$$

holds for all possible indices, in order to satisfy $\|h(x) - \bar{h}(x)\|_{\infty, [-a_n, a_n]^d} \leq \varepsilon_n$. For n sufficiently large, which is assumed permanently in the following, the coefficients $c_{i,j}^{(r)}$ have to take values in $[-\max\{\alpha, \beta, \gamma\}, \max\{\alpha, \beta, \gamma\}]$ and the relations $\max\{\alpha, \beta, \gamma\} \leq \log(n) \cdot n^2 \cdot M_n^{d^*} \leq n^4$ and $a_n \leq M_n \leq n$ hold. Then due to $\varepsilon_n = \frac{a_n^p}{M_n^p}$ a number of

$$\left\lceil \frac{2 \cdot \max\{\alpha, \beta, \gamma\} \cdot a_n \cdot n^{c_{11}}}{2 \cdot \varepsilon_n} \right\rceil \leq n^{c_{12}}$$

different $\bar{c}_{i,j}^{(r)}$ suffices to guarantee, that at least one of them satisfies the relation (41) for any $c_{i,j}^{(r)}$ with fixed indices. Furthermore, the coefficients $c_{i,j}^{(r)}$, which can actually differ

regarding different $h \in \mathcal{H}^{(l)}$, are the ones originating from the coefficients $a_{i,j,m}, b_{i,j}, d_i$ in the definition of $\mathcal{F}_{M_n, d^*, d, \alpha, \beta, \gamma}^{(neural\ networks)}$. Using (22), their number can be bounded by $c_{13} \cdot M_n^{d^*}$. So the logarithm of the covering number $\mathcal{N}(\varepsilon_n, \mathcal{F}_n, \|\cdot\|_{\infty, [-a_n, a_n]^d})$ can be bounded by

$$\log \left(\mathcal{N}(\varepsilon_n, \mathcal{F}_n, \|\cdot\|_{\infty, [-a_n, a_n]^d}) \right) \leq \log \left((n^{c_{12}})^{c_{13} \cdot M_n^{d^*}} \right) \leq c_{10} \cdot \log(n) \cdot M_n^{d^*},$$

which proves the assertion. □