

Data mining with sparse grids using simplicial basis functions

Jochen Garcke and Michael Griebel
Institut für Angewandte Mathematik
Abteilung für wissenschaftliches Rechnen und numerische Simulation
Universität Bonn
Wegelerstraße 6
D-53115 Bonn, Germany
(garcke, griebel)@iam.uni-bonn.de

ABSTRACT

Recently we presented a new approach [18] to the classification problem arising in data mining. It is based on the regularization network approach but, in contrast to other methods which employ ansatz functions associated to data points, we use a grid in the usually high-dimensional feature space for the minimization process. To cope with the curse of dimensionality, we employ sparse grids [49]. Thus, only $O(h_n^{-1} n^{d-1})$ instead of $O(h_n^{-d})$ grid points and unknowns are involved. Here d denotes the dimension of the feature space and $h_n = 2^{-n}$ gives the mesh size. We use the sparse grid combination technique [28] where the classification problem is discretized and solved on a sequence of conventional grids with uniform mesh sizes in each dimension. The sparse grid solution is then obtained by linear combination. In contrast to our former work, where d -linear functions were used, we now apply linear basis functions based on a simplicial discretization. This allows to handle more dimensions and the algorithm needs less operations per data point.

We describe the sparse grid combination technique for the classification problem, give implementational details and discuss the complexity of the algorithm. It turns out that the method scales linearly with the number of given data points. Finally we report on the quality of the classifier built by our new method on data sets with up to 10 dimensions. It turns out that our new method achieves correctness rates which are competitive to that of the best existing methods.

Keywords

data mining, classification, approximation, sparse grids, combination technique, simplicial discretization

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 2001 ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

1. INTRODUCTION

Data mining is the process of finding patterns, relations and trends in large data sets. Examples range from scientific applications like the post-processing of data in medicine or the evaluation of satellite pictures to financial and commercial applications, e.g. the assessment of credit risks or the selection of customers for advertising campaign letters. For an overview on data mining and its various tasks and approaches see [5, 12].

In this paper we consider the classification problem arising in data mining. Given is a set of data points in a d -dimensional feature space together with a class label. From this data, a classifier must be constructed which allows to predict the class of any newly given data point for future decision making. Widely used approaches are, besides others, decision tree induction, rule learning, adaptive multivariate regression splines, neural networks, and support vector machines. Interestingly, some of these techniques can be interpreted in the framework of regularization networks [21]. This approach allows a direct description of the most important neural networks and it also allows for an equivalent description of support vector machines and n -term approximation schemes [20]. Here, the classification of data is interpreted as a scattered data approximation problem with certain additional regularization terms in high-dimensional spaces.

In [18] we presented a new approach to the classification problem. It is also based on the regularization network approach but, in contrast to the other methods which employ mostly global ansatz functions associated to data points, we use an independent grid with associated local ansatz functions in the minimization process. This is similar to the numerical treatment of partial differential equations. Here, a uniform grid would result in $O(h_n^{-d})$ grid points, where d denotes the dimension of the feature space and $h_n = 2^{-n}$ gives the mesh size. Therefore the complexity of the problem would grow exponentially with d and we encounter the curse of dimensionality. This is probably the reason why conventional grid-based techniques are not used in data mining up to now.

However, there is the so-called sparse grids approach which allows to cope with the complexity of the problem to some extent. This method has been originally developed for the solution of partial differential equations [2, 8, 28, 49] and

is now used successfully also for integral equations [14, 27], interpolation and approximation [3, 26, 39, 42], eigenvalue problems [16] and integration problems [19]. In the information based complexity community it is also known as 'hyperbolic cross points' and the idea can even be traced back to [41]. For a d -dimensional problem, the sparse grid approach employs only $O(h_n^{-1}(\log(h_n^{-1}))^{d-1})$ grid points in the discretization. The accuracy of the approximation however is nearly as good as for the conventional full grid methods, provided that certain additional smoothness requirements are fulfilled. Thus a sparse grid discretization method can be employed also for higher-dimensional problems. The curse of the dimensionality of conventional 'full' grid methods affects sparse grids much less.

In this paper, we apply the sparse grid combination technique [28] to the classification problem. For that the regularization network problem is discretized and solved on a certain sequence of conventional grids with uniform mesh sizes in each coordinate direction. In contrast to [18], where d -linear functions stemming from a tensor-product approach were used, we now apply linear basis functions based on a simplicial discretization. In comparison, this approach allows the processing of more dimensions and needs less operations per data point. The sparse grid solution is then obtained from the solutions on the different grids by linear combination. Thus the classifier is built on sparse grid points and not on data points. A discussion of the complexity of the method gives that the method scales linearly with the number of instances, i.e. the amount of data to be classified. Therefore, our method is well suited for realistic data mining applications where the dimension of the feature space is moderately high (e.g. after some preprocessing steps) but the amount of data is very large. Furthermore the quality of the classifier build by our new method seems to be very good. Here we consider standard test problems from the UCI repository and problems with huge synthetical data sets in up to 10 dimensions. It turns out that our new method achieves correctness rates which are competitive to those of the best existing methods. Note that the combination method is simple to use and can be parallelized in a natural and straightforward way.

The remainder of this paper is organized as follows: In Section 2 we describe the classification problem in the framework of regularization networks as minimization of a (quadratic) functional. We then discretize the feature space and derive the associated linear problem. Here we focus on grid-based discretization techniques. Then, we introduce the sparse grid combination technique for the classification problem and discuss its properties. Furthermore, we present a new variant based on a discretization by simplices and discuss complexity aspects. Section 3 presents the results of numerical experiments conducted with the sparse grid combination method, demonstrates the quality of the classifier build by our new method and compares the results with the ones from [18] and with the ones obtained with different forms of SVMs [33]. Some final remarks conclude the paper.

2. THE PROBLEM

Classification of data can be interpreted as traditional scattered data approximation problem with certain additional regularization terms. In contrast to conventional scattered data approximation applications, we now encounter quite high-dimensional spaces. To this end, the approach of

regularization networks [21] gives a good framework. This approach allows a direct description of the most important neural networks and it also allows for an equivalent description of support vector machines and n -term approximation schemes [20].

Consider the given set of already classified data (the training set)

$$S = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}\}_{i=1}^M.$$

Assume now that these data have been obtained by sampling of an unknown function f which belongs to some function space V defined over \mathbb{R}^d . The sampling process was disturbed by noise. The aim is now to recover the function f from the given data as good as possible. This is clearly an ill-posed problem since there are infinitely many solutions possible. To get a well-posed, uniquely solvable problem we have to assume further knowledge on f . To this end, regularization theory [43, 47] imposes an additional smoothness constraint on the solution of the approximation problem and the regularization network approach considers the variational problem

$$\min_{f \in V} R(f)$$

with

$$R(f) = \frac{1}{M} \sum_{i=1}^m C(f(\mathbf{x}_i), y_i) + \lambda \Phi(f). \quad (1)$$

Here, $C(\cdot, \cdot)$ denotes an error cost function which measures the interpolation error and $\Phi(f)$ is a smoothness functional which must be well defined for $f \in V$. The first term enforces closeness of f to the data, the second term enforces smoothness of f and the regularization parameter λ balances these two terms. Typical examples are

$$C(x, y) = |x - y| \text{ or } C(x, y) = (x - y)^2,$$

and

$$\Phi(f) = \|Pf\|_2^2 \quad \text{with} \quad Pf = \nabla f \text{ or } Pf = \Delta f,$$

with ∇ denoting the gradient and Δ the Laplace operator. The value of λ can be chosen according to cross-validation techniques [13, 22, 37, 44] or to some other principle, such as structural risk minimization [45]. Note that we find exactly this type of formulation in the case $d = 2, 3$ in scattered data approximation methods, see [1, 31], where the regularization term is usually physically motivated.

2.1 Discretization

We now restrict the problem to a finite dimensional subspace $V_N \in V$. The function f is then replaced by

$$f_N = \sum_{j=1}^N \alpha_j \varphi_j(\mathbf{x}). \quad (2)$$

Here the ansatz functions $\{\varphi_j\}_{j=1}^N$ should span V_N and preferably should form a basis for V_N . The coefficients $\{\alpha_j\}_{j=1}^N$ denote the degrees of freedom. Note that the restriction to a suitably chosen finite-dimensional subspace involves some additional regularization (regularization by discretization) which depends on the choice of V_N .

In the remainder of this paper, we restrict ourselves to the choice

$$C(f_N(\mathbf{x}_i), y_i) = (f_N(\mathbf{x}_i) - y_i)^2$$

and

$$\Phi(f_N) = \|Pf_N\|_{L_2}^2 \quad (3)$$

for some given linear operator P . This way we obtain from the minimization problem a feasible linear system. We thus have to minimize

$$R(f_N) = \frac{1}{M} \sum_{i=1}^M (f_N(\mathbf{x}_i) - y_i)^2 + \lambda \|Pf_N\|_{L_2}^2, \quad (4)$$

with f_N in the finite dimensional space V_N . We plug (2) into (4) and obtain after differentiation with respect to α_k , $k = 1, \dots, N$

$$0 = \frac{\partial R(f_N)}{\partial \alpha_k} = \frac{2}{M} \sum_{i=1}^M \left(\sum_{j=1}^N \alpha_j \varphi_j(\mathbf{x}_i) - y_i \right) \cdot \varphi_k(\mathbf{x}_i) + 2\lambda \sum_{j=1}^N \alpha_j (P\varphi_j, P\varphi_k)_{L_2} \quad (5)$$

This is equivalent to ($k = 1, \dots, N$)

$$\sum_{j=1}^N \alpha_j \left[M\lambda (P\varphi_j, P\varphi_k)_{L_2} + \sum_{i=1}^M \varphi_j(\mathbf{x}_i) \cdot \varphi_k(\mathbf{x}_i) \right] = \sum_{i=1}^M y_i \varphi_k(\mathbf{x}_i). \quad (6)$$

In matrix notation we end up with the linear system

$$(\lambda C + B \cdot B^T) \alpha = B y. \quad (7)$$

Here C is a square $N \times N$ matrix with entries $C_{j,k} = M \cdot (P\varphi_j, P\varphi_k)_{L_2}$, $j, k = 1, \dots, N$, and B is a rectangular $N \times M$ matrix with entries $B_{j,i} = \varphi_j(\mathbf{x}_i)$, $i = 1, \dots, M$, $j = 1, \dots, N$. The vector y contains the data labels y_i and has length M . The unknown vector α contains the degrees of freedom α_j and has length N .

Depending on the regularization operator we obtain different minimization problems in V_N . For example if we use the gradient $\Phi(f_N) = \|\nabla f_N\|_{L_2}^2$ in the regularization expression in (1) we obtain a Poisson problem with an additional term which resembles the interpolation problem. The natural boundary conditions for such a partial differential equation are Neumann conditions. The discretization (2) gives us then the linear system (7) where C corresponds to a discrete Laplacian. To obtain the classifier f_N we now have to solve this system.

2.2 Grid based discrete approximation

Up to now we have not yet been specific what finite-dimensional subspace V_N and what type of basis functions $\{\varphi_j\}_{j=1}^N$ we want to use. In contrast to conventional data mining approaches which work with ansatz functions associated to data points we now use a certain grid in the attribute space to determine the classifier with the help of these grid points. This is similar to the numerical treatment of partial differential equations.

For reasons of simplicity, here and in the remainder of this paper, we restrict ourselves to the case $\mathbf{x}_i \in \Omega = [0, 1]^d$. This situation can always be reached by a proper rescaling of the data space. A conventional finite element discretization would now employ an equidistant grid Ω_n with mesh size $h_n = 2^{-n}$ for each coordinate direction, where n is the

refinement level. In the following we always use the gradient $P = \nabla$ in the regularization expression (3). Let \mathbf{j} denote the multi-index $(j_1, \dots, j_d) \in \mathbb{N}^d$. A finite element method with piecewise d -linear, i.e. linear in each dimension, test- and trial-functions $\phi_{n,\mathbf{j}}(\mathbf{x})$ on grid Ω_n now would give

$$(f_N(\mathbf{x}) =) f_n(\mathbf{x}) = \sum_{j_1=0}^{2^n} \dots \sum_{j_d=0}^{2^n} \alpha_{n,\mathbf{j}} \phi_{n,\mathbf{j}}(\mathbf{x})$$

and the variational procedure (4) - (6) would result in the discrete linear system

$$(\lambda C_n + B_n \cdot B_n^T) \alpha_n = B_n y \quad (8)$$

of size $(2^n + 1)^d$ and matrix entries corresponding to (7). Note that f_n lives in the space

$$V_n := \text{span}\{\phi_{n,\mathbf{j}}, j_t = 0, \dots, 2^n, t = 1, \dots, d\}.$$

The discrete problem (8) might in principle be treated by an appropriate solver like the conjugate gradient method, a multigrid method or some other suitable efficient iterative method. However, this direct application of a finite element discretization and the solution of the resulting linear system by an appropriate solver is clearly not possible for a d -dimensional problem if d is larger than four. The number of grid points is of the order $O(h_n^{-d}) = O(2^{nd})$ and, in the best case, the number of operations is of the same order. Here we encounter the so-called curse of dimensionality: The complexity of the problem grows exponentially with d . At least for $d > 4$ and a reasonable value of n , the arising system can not be stored and solved on even the largest parallel computers today.

2.3 The sparse grid combination technique

Therefore we proceed as follows: We discretize and solve the problem on a certain sequence of grids $\Omega_1 = \Omega_{l_1, \dots, l_d}$ with uniform mesh sizes $h_t = 2^{-l_t}$ in the t -th coordinate direction. These grids may possess different mesh sizes for different coordinate directions. To this end, we consider all grids Ω_1 with

$$l_1 + \dots + l_d = n + (d - 1) - q, \quad q = 0, \dots, d - 1, \quad l_t > 0. \quad (9)$$

For the two-dimensional case, the grids needed in the combination formula of level 4 are shown in Figure 1. The finite element approach with piecewise d -linear test- and trial-functions

$$\phi_{1,\mathbf{j}}(\mathbf{x}) := \prod_{t=1}^d \phi_{l_t, j_t}(x_t) \quad (10)$$

on grid Ω_1 now would give

$$f_1(\mathbf{x}) = \sum_{j_1=0}^{2^{l_1}} \dots \sum_{j_d=0}^{2^{l_d}} \alpha_{1,\mathbf{j}} \phi_{1,\mathbf{j}}(\mathbf{x})$$

and the variational procedure (4) - (6) would result in the discrete system

$$(\lambda C_1 + B_1 \cdot B_1^T) \alpha_1 = B_1 y \quad (11)$$

with the matrices

$$(C_1)_{\mathbf{j},\mathbf{k}} = M \cdot (\nabla \phi_{1,\mathbf{j}}, \nabla \phi_{1,\mathbf{k}}) \quad \text{and} \quad (B_1)_{\mathbf{j},i} = \phi_{1,\mathbf{j}}(\mathbf{x}_i),$$

$j_t, k_t = 0, \dots, 2^{l_t}, t = 1, \dots, d, i = 1, \dots, M$, and the unknown vector $(\alpha_1)_{\mathbf{j}}, j_t = 0, \dots, 2^{l_t}, t = 1, \dots, d$. We then solve these

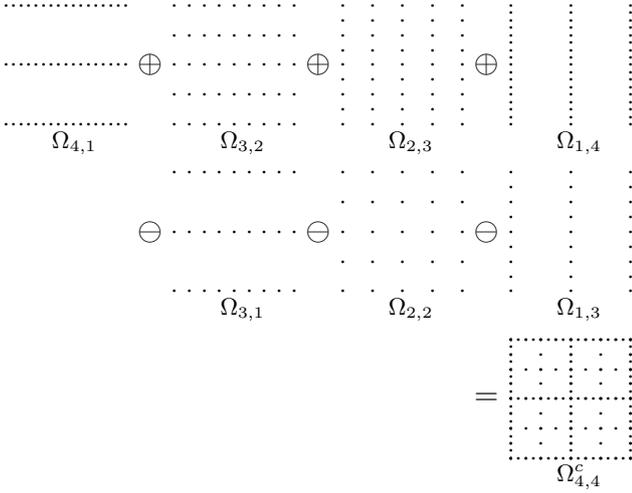


Figure 1: Combination technique with level $n = 4$ in two dimensions

problems by a feasible method. To this end we use here a diagonally preconditioned conjugate gradient algorithm. But also an appropriate multigrid method with partial semi-coarsening can be applied. The discrete solutions f_1 are contained in the spaces

$$V_1 := \text{span}\{\phi_{1,j}, j_t = 0, \dots, 2^{l_t}, t = 1, \dots, d\}, \quad (12)$$

of piecewise d -linear functions on grid Ω_1 .

Note that all these problems are substantially reduced in size in comparison to (8). Instead of one problem with size $\dim(V_n) = O(h_n^{-d}) = O(2^{nd})$, we now have to deal with $O(dn^{d-1})$ problems of size $\dim(V_1) = O(h_n^{-1}) = O(2^n)$. Moreover, all these problems can be solved independently, which allows for a straightforward parallelization on a coarse grain level, see [23]. There is also a simple but effective static load balancing strategy available [25].

Finally we linearly combine the results $f_1(\mathbf{x}) \in V_1$, $f_1 = \sum_j \alpha_{1,j} \phi_{1,j}(\mathbf{x})$, from the different grids Ω_1 as follows:

$$f_n^{(c)}(\mathbf{x}) := \sum_{q=0}^{d-1} (-1)^q \binom{d-1}{q} \sum_{|\mathbf{l}|_1 = n + (d-1) - q} f_1(\mathbf{x}). \quad (13)$$

The resulting function $f_n^{(c)}$ lives in the sparse grid space

$$V_n^{(s)} := \bigcup_{\substack{l_1 + \dots + l_d = n + (d-1) - q \\ q = 0, \dots, d-1 \quad l_t > 0}} V_1.$$

This space has $\dim(V_n^{(s)}) = O(h_n^{-1}(\log(h_n^{-1}))^{d-1})$. It is spanned by a piecewise d -linear hierarchical tensor product basis, see [8].

Note that the summation of the discrete functions from different spaces V_1 in (13) involves d -linear interpolation which resembles just the transformation to a representation in the hierarchical basis. For details see [24, 28, 29]. However we never explicitly assemble the function $f_n^{(c)}$ but keep instead the solutions f_1 on the different grids Ω_1 which arise in the combination formula. Now, any linear operation F on $f_n^{(c)}$ can easily be expressed by means of the combination

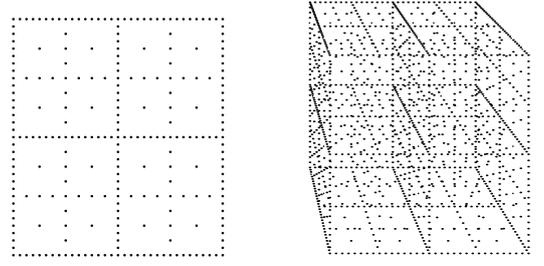


Figure 2: Two-dimensional sparse grid (left) and three-dimensional sparse grid (right), $n = 5$

formula (13) acting directly on the functions f_1 , i.e.

$$F(f_n^{(c)}) = \sum_{q=0}^{d-1} (-1)^q \binom{d-1}{q} \sum_{l_1 + \dots + l_d = n + (d-1) - q} F(f_1). \quad (14)$$

Therefore, if we now want to evaluate a newly given set of data points $\{\tilde{\mathbf{x}}_i\}_{i=1}^{\tilde{M}}$ (the test or evaluation set) by

$$\tilde{y}_i := f_n^{(c)}(\tilde{\mathbf{x}}_i), \quad i = 1, \dots, \tilde{M}$$

we just form the combination of the associated values for f_1 according to (13). The evaluation of the different f_1 in the test points can be done completely in parallel, their summation needs basically an all-reduce/gather operation.

For second order elliptic PDE model problems, it was proven that the combination solution $f_n^{(c)}$ is almost as accurate as the full grid solution f_n , i.e. the discretization error satisfies

$$\|e_n^{(c)}\|_{L_p} := \|f - f_n^{(c)}\|_{L_p} = O(h_n^2 \log(h_n^{-1})^{d-1})$$

provided that a slightly stronger smoothness requirement on f than for the full grid approach holds. We need the seminorm

$$|f|_\infty := \left\| \frac{\partial^{2d} f}{\prod_{j=1}^d \partial x_j^2} \right\|_\infty \quad (15)$$

to be bounded. Furthermore, a series expansion of the error is necessary for the combination technique. Its existence was shown for PDE model problems in [10].

The combination technique is only one of the various methods to solve problems on sparse grids. Note that there exist also finite difference [24, 38] and Galerkin finite element approaches [2, 8, 9] which work directly in the hierarchical product basis on the sparse grid. But the combination technique is conceptually much simpler and easier to implement. Moreover it allows to reuse standard solvers for its different subproblems and is straightforwardly parallelizable.

2.4 Simplicial basis functions

So far we only mentioned d -linear basis functions based on a tensor-product approach, this case was presented in detail in [18]. But on the grids of the combination technique linear basis functions based on a simplicial discretization are also possible. For that we use the so-called Kuhn's triangulation [15, 32] for each rectangular block, see Figure 3. Now, the summation of the discrete functions for the different spaces V_1 in (13) only involves linear interpolation.

Table 1: Complexities of the storage, the assembly and the matrix-vector multiplication for the different matrices arising in the combination method on one grid Ω_1 for both discretization approaches. C_1 and G_1 can be stored together in one matrix structure.

	d -linear basis functions			linear basis functions		
	C_1	$G_1 := B_1 \cdot B_1^T$	B_1	C_1	$G_1 := B_1 \cdot B_1^T$	B_1
storage	$O(3^d \cdot N)$	$O(3^d \cdot N)$	$O(2^d \cdot M)$	$O((2 \cdot d + 1) \cdot N)$	$O(2^d \cdot N)$	$O((d + 1) \cdot M)$
assembly	$O(3^d \cdot N)$	$O(d \cdot 2^{2d} \cdot M)$	$O(d \cdot 2^d \cdot M)$	$O((2 \cdot d + 1) \cdot N)$	$O((d + 1)^2 \cdot M)$	$O((d + 1) \cdot M)$
mv-multiplication	$O(3^d \cdot N)$	$O(3^d \cdot N)$	$O(2^d \cdot M)$	$O((2 \cdot d + 1) \cdot N)$	$O(2^d \cdot N)$	$O((d + 1) \cdot M)$



Figure 3: Kuhn's triangulation of a three-dimensional unit cube

The theoretical properties of this variant of the sparse grid technique still has to be investigated in more detail. However the results which are presented in section 3 warrant its use. We see, if at all, just slightly worse results with linear basis functions than with d -linear basis functions and we believe that our new approach results in the same approximation order.

Since in our new variant of the combination technique the overlap of supports, i.e. the regions where two basis functions are both non-zero, is greatly reduced due to the use of a simplicial discretization, the complexities scale significantly better. This concerns both the costs of the assembly and the storage of the non-zero entries of the sparsely populated matrices from (8), see Table 1. Note that for general operators P the complexities for C_1 scale with $O(2^d \cdot N)$. But for our choice of $P = \nabla$ structural zero-entries arise, which need not to be considered, and which further reduce the complexities, see Table 1 (right), column C_1 . The actual iterative solution process (by a diagonally preconditioned conjugate gradient method) scales independent of the number of data points for both approaches.

Note however that both the storage and the run time complexities still depend exponentially on the dimension d . Presently, due to the limitations of the memory of modern workstations (512 MByte - 2 GByte), we therefore can only deal with the case $d \leq 8$ for d -linear basis functions and $d \leq 11$ for linear basis functions. A decomposition of the matrix entries over several computers in a parallel environment would permit more dimensions.

3. NUMERICAL RESULTS

We now apply our approach to different test data sets. Here we use both synthetic data and real data from practical data mining applications. All the data sets are rescaled to $[0, 1]^d$. To evaluate our method we give the correctness rates on testing data sets, if available, or the ten-fold cross-validation results otherwise. For further details and a criti-

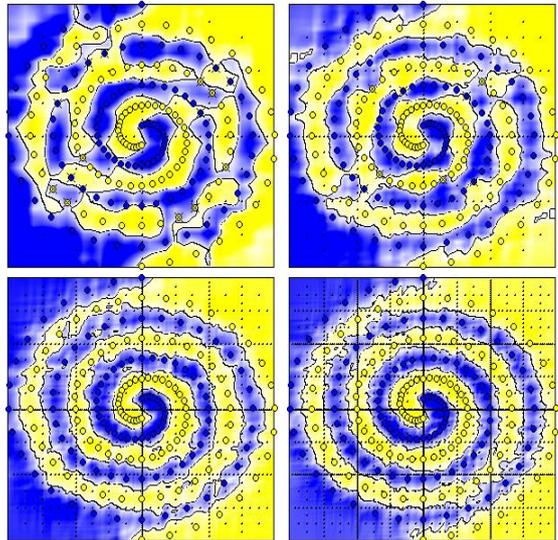


Figure 4: Spiral data set, sparse grid with level 5 (top left) to 8 (bottom right)

cal discussion on the evaluation of the quality of classification algorithms see [13, 37].

3.1 Two-dimensional problems

We first consider synthetic two-dimensional problems with small sets of data which correspond to certain structures.

3.1.1 Spiral

The first example is the spiral data set, proposed by Alexis Wieland of MITRE Corp [48]. Here, 194 data points describe two intertwined spirals, see Figure 4. This is surely an artificial problem which does not appear in practical applications. However it serves as a hard test case for new data mining algorithms. It is known that neural networks can have severe problems with this data set and some neural networks can not separate the two spirals at all [40].

In Table 2 we give the correctness rates achieved with the leave-one-out cross-validation method, i.e. a 194-fold cross-validation. The best testing correctness was achieved on level 8 with 89.18% in comparison to 77.20% in [40].

In Figure 4 we show the corresponding results obtained with our sparse grid combination method for the levels 5 to 8. With level 7 the two spirals are clearly detected and resolved. Note that here 1281 grid points are contained in the sparse grid. For level 8 (2817 sparse grid points) the shape of the two reconstructed spirals gets smoother and

Table 3: Results for the Ripley data set

level	linear basis			d -linear basis	best possible %	
	ten-fold test %	λ	test data %	test data %	linear	d -linear
1	85.2	0.0020	89.9	89.8	90.6	90.3
2	85.2	0.0065	90.3	90.4	90.4	90.9
3	88.4	0.0020	90.9	90.6	91.0	91.2
4	87.2	0.0035	91.4	90.6	91.4	91.2
5	88.0	0.0055	91.3	90.9	91.5	91.1
6	86.8	0.0045	90.7	90.8	90.7	90.8
7	86.8	0.0008	89.0	88.8	91.1	91.0
8	87.2	0.0037	91.0	89.7	91.2	91.0
9	87.7	0.0015	90.1	90.9	91.1	91.0
10	89.2	0.0020	91.0	90.6	91.2	91.1

level	λ	training correctness	testing correctness
5	0.0003	94.87 %	82.99 %
6	0.0006	97.42 %	84.02 %
7	0.00075	100.00 %	88.66 %
8	0.0006	100.00 %	89.18 %
9	0.0006	100.00 %	88.14 %

Table 2: Leave-one-out cross-validation results for the spiral data set

the reconstruction gets more precise.

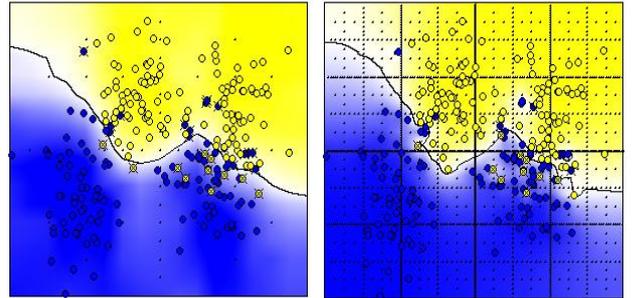
3.1.2 Ripley

This data set, taken from [36], consists of 250 training data and 1000 test points. The data set was generated synthetically and is known to exhibit 8 % error. Thus no better testing correctness than 92 % can be expected.

Since we now have training and testing data, we proceed as follows: First we use the training set to determine the best regularization parameter λ per ten-fold cross-validation. The best test correctness rate and the corresponding λ are given for different levels n in the first two columns of Table 3. With this λ we then compute the sparse grid classifier from the 250 training data. The column three of Table 3 gives the result of this classifier on the (previously unknown) test data set. We see that our method works well. Already level 4 is sufficient to obtain results of 91.4 %. The reason is surely the relative simplicity of the data, see Figure 5. Just a few hyperplanes should be enough to separate the classes quite properly. We also see that there is not much need to use any higher levels, on the contrary there is even an overfitting effect visible in Figure 5.

In column 4 we show the results from [18], there we achieve almost the same results with d -linear functions.

To see what kind of results could be possible with a more sophisticated strategy for determining λ we give in the last two columns of Table 3 the testing correctness which is achieved for the best possible λ . To this end we compute for *all* (discrete) values of λ the sparse grid classifiers from the 250 data points and evaluate them on the test set. We then pick the best result. We clearly see that there is not much of a difference. This indicates that the approach to determine the value of λ from the training set by cross-validation works well. Again we have almost the same results with linear and d -linear basis functions. Note that a testing correctness of

**Figure 5: Ripley data set, combination technique with linear basis functions. Left: level 4, $\lambda = 0.0035$. Right: level 8, $\lambda = 0.0037$**

90.6 % and 91.1 % was achieved in [36] and [35], respectively, for this data set.

3.2 6-dimensional problems

3.2.1 BUPA Liver

The BUPA Liver Disorders data set from Irvine Machine Learning Database Repository [6] consists of 345 data points with 6 features and a selector field used to split the data into 2 sets with 145 instances and 200 instances respectively. Here we have no test data and therefore can only report our ten-fold cross-validation results.

We compare with our d -linear results from [18] and with the two best results from [33], the therein introduced smoothed support vector machine (SSVM) and the classical support vector machine ($SVM_{\|\cdot\|_2}$) [11, 46]. The results are given in Table 4.

As expected, our sparse grid combination approach with linear basis functions performs slightly worse than the d -linear approach. The best test result was 69.60% on level 4. The new variant of the sparse grid combination technique performs only slightly worse than the SSVM, whereas the d -linear variant performs slightly better than the support vector machines. Note that the results for other SVM approaches like the support vector machine using the 1-norm approach ($SVM_{\|\cdot\|_1}$) were reported to be somewhat worse in [33].

Table 4: Results for the BUPA liver disorders data set

		linear		d -linear		For comparison with other methods	
		λ	%	λ	%		
level 1	10-fold train. correctness	0.012	76.00	0.020	76.00	SVM [33]	
	10-fold test. correctness		69.00		67.87	SSVM SVM $_{\ \cdot\ _2}$	
level 2	10-fold train. correctness	0.040	76.13	0.10	77.49	70.37	70.57
	10-fold test. correctness		66.01		67.84	70.33	69.86
level 3	10-fold train. correctness	0.165	78.71	0.007	84.28		
	10-fold test. correctness		66.41		70.34		
level 4	10-fold train. correctness	0.075	92.01	0.0004	90.27		
	10-fold test. correctness		69.60		70.92		

3.2.2 Synthetic massive data set in 6D

To measure the performance on a massive data set we produced with DatGen [34] a 6-dimensional test case with 5 million training points and 20 000 points for testing. We used the call `datgen -r1 -X0/100,R,O:0/100,R,O:0/100,R,O:0/100,R,O:0/200,R,O:0/200,R,O -R2 -C2/4 -D2/5 -T10/60 -O5020000 -p -e0.15`.

The results are given in Table 5. Note that already on level 1 a testing correctness of over 90 % was achieved with just $\lambda = 0.01$. The main observation on this test case concerns the execution time, measured on a Pentium III 700 MHz machine. Besides the total run time, we also give the CPU time which is needed for the computation of the matrices $G_1 = B_1 \cdot B_1^T$.

We see that with linear basis functions really huge data sets of 5 million points can be processed in reasonable time. Note that more than 50 % of the computation time is spent for the data matrix assembly only and, more importantly, that the execution time scales linearly with the number of data points. The latter is also the case for the d -linear functions, but, as mentioned, this approach needs more operations per data point and results in a much longer execution time, compare also Table 5. Especially the assembly of the data matrix needs more than 96 % of the total run time for this variant. For our present example the linear basis approach is about 40 times faster than the d -linear approach on the same refinement level, e.g. for level 2 we need 17 minutes in the linear case and 11 hours in the d -linear case. For higher dimensions the factor will be even larger.

3.3 10-dimensional problems

3.3.1 Forest cover type

The forest cover type dataset comes from the UCI KDD Archive [4], it was also used in [30], where an approach similar to ours was followed. It consists of cartographic variables for 30 x 30 meter cells and a forest cover type is to be predicted. The 12 originally measured attributes resulted in 54 attributes in the data set, besides 10 quantitative variables there are 4 binary wilderness areas and 40 binary soil type variables. We only use the 10 quantitative variables. The class label has 7 values, Spruce/Fir, Lodgepole Pine, Ponderosa Pine, Cottonwood/Willow, Aspen, Douglas-fir and Krummholz. Like [30] we only report results for the classification of Ponderosa Pine, which has 35754 instances out of the total 581012.

Since far less than 10 % of the instances belong to Ponderosa Pine we weigh this class with a factor of 5, i.e. Ponderosa Pine has a class value of 5, all others of -1 and the

threshold value for separating the classes is 0. The data set was randomly separated into a training set, a test set, and an evaluation set, all similar in size.

In [30] only results up to 6 dimensions could be reported. In Table 6 we present our results for the 6 dimensions chosen there, i.e. the dimensions 1,4,5,6,7, and 10, and for all 10 dimensions as well. To give an overview of the behavior over several λ 's we present for each level n the overall correctness results, the correctness results for Ponderosa Pine and the correctness result for the other class for three values of λ . We then give results on the evaluation set for a chosen λ .

We see in Table 6 that already with level 1 we have a testing correctness of 93.95 % for the Ponderosa Pine in the 6 dimensional version. Higher refinement levels do not give better results. The result of 93.52% on the evaluation set is almost the same as the corresponding testing correctness. Note that in [30] a correctness rate of 86.97 % was achieved on the evaluation set.

The usage of all 10 dimensions improves the results slightly, we get 93.81 % as our evaluation result on level 1. As before higher refinement levels do not improve the results for this data set.

Note that the forest cover example is sound enough as an example of classification, but it might strike forest scientists as being amusingly superficial. It has been known for 30 years that the dynamics of forest growth can have a dominant effect on which species is present at a given location [7], yet there are no dynamic variables in the classifier. This one can see as a warning that it should never be assumed that the available data contains all the relevant information.

3.3.2 Synthetic massive data set in 10D

To measure the performance on a still higher dimensional massive data set we produced with DatGen [34] a 10-dimensional test case with 5 million training points and 50 000 points for testing. We used the call `datgen -r1 -X0/200,R,O:0/200,R,O:0/200,R,O:0/200,R,O:0/200,R,O:0/200,R,O:0/200,R,O:0/200,R,O:0/200,R,O -R2 -C2/6 -D2/7 -T10/60 -O5050000 -p -e0.15`.

Like in the synthetical 6-dimensional example the main observations concern the run time, measured on a Pentium III 700 MHz machine. Besides the total run time, we also give the CPU time which is needed for the computation of the matrices $G_1 = B_1 \cdot B_1^T$. Note that the highest amount of memory needed (for level 2 in the case of 5 million data points) was 500 MBytes, about 250 MBytes for the matrix and about 250 MBytes for keeping the data points in memory.

More than 50 % of the run time is spent for the assembly

Table 5: Results for a 6D synthetic massive data set, $\lambda = 0.01$

	# of points	training correctness	testing correctness	total time (sec)	data matrix time (sec)	# of iterations
linear basis functions						
level 1	50 000	90.4	90.5	3	1	23
	500 000	90.5	90.5	25	8	25
	5 million	90.5	90.6	242	77	28
level 2	50 000	91.4	91.0	12	5	184
	500 000	91.2	91.1	110	55	204
	5 million	91.1	91.2	1086	546	223
level 3	50 000	92.2	91.4	48	23	869
	500 000	91.7	91.7	417	226	966
	5 million	91.6	91.7	4087	2239	1057
<i>d</i> -linear basis functions						
level 1	500 000	90.7	90.8	597	572	91
	5 million	90.7	90.7	5897	5658	102
level 2	500 000	91.5	91.6	4285	4168	656
	5 million	91.4	91.5	42690	41596	742

of the data matrix and the time needed for the data matrix scales linearly with the number of data points, see Table 7. The total run time seems to scale even better than linear.

4. CONCLUSIONS

We presented the sparse grid combination technique with linear basis functions based on simplices for the classification of data in moderate-dimensional spaces. Our new method gave good results for a wide range of problems. It is capable to handle huge data sets with 5 million points and more. The run time scales only linearly with the number of data. This is an important property for many practical applications where often the dimension of the problem can substantially be reduced by certain preprocessing steps but the number of data can be extremely huge. We believe that our sparse grid combination method possesses a great potential in such practical application problems.

We demonstrated for the Ripley data set how the best value of the regularization parameter λ can be determined. This is also of practical relevance.

A parallel version of the sparse grid combination technique reduces the run time significantly, see [17]. Note that our method is easily parallelizable already on a coarse grain level. A second level of parallelization is possible on each grid of the combination technique with the standard techniques known from the numerical treatment of partial differential equations.

Since not necessarily all dimensions need the maximum refinement level, a modification of the combination technique with regard to different refinement levels in each dimension along the lines of [19] seems to be promising.

Note furthermore that our approach delivers a continuous classifier function which approximates the data. It therefore can be used without modification for regression problems as well. This is in contrast to many other methods like e.g. decision trees. Also more than two classes can be handled by using isolines with just different values.

Finally, for reasons of simplicity, we used the operator $P = \nabla$. But other differential (e.g. $P = \Delta$) or pseudo-differential operators can be employed here with their associated regular finite element ansatz functions.

5. ACKNOWLEDGEMENTS

Part of the work was supported by the German Bundesministerium für Bildung und Forschung (BMB+F) within the project 03GRM6BN. This work was carried out in cooperation with Prudential Systems Software GmbH, Chemnitz. The authors thank one of the referees for his remarks on the forest cover data set.

6. REFERENCES

- [1] E. Arge, M. Dæhlen, and A. Tveito. Approximation of scattered data using smooth grid functions. *J. Comput. Appl. Math.*, 59:191–205, 1995.
- [2] R. Balder. *Adaptive Verfahren für elliptische und parabolische Differentialgleichungen auf dünnen Gittern*. Dissertation, Technische Universität München, 1994.
- [3] G. Baszenski. *N*-th order polynomial spline blending. In W. Schempp and K. Zeller, editors, *Multivariate Approximation Theory III*, ISNM 75, pages 35–46. Birkhäuser, Basel, 1985.
- [4] S. D. Bay. The UCI KDD archive. <http://kdd.ics.uci.edu>, 1999.
- [5] M. J. A. Berry and G. S. Linoff. *Mastering Data Mining*. Wiley, 2000.
- [6] C. L. Blake and C. J. Merz. UCI repository of machine learning databases, 1998. <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
- [7] D. Botkin, J. Janak, and J. Wallis. Some ecological consequences of a computer model of forest growth. *J. Ecology*, 60:849–872, 1972.
- [8] H.-J. Bungartz. *Dünne Gitter und deren Anwendung bei der adaptiven Lösung der dreidimensionalen Poisson-Gleichung*. Dissertation, Institut für Informatik, Technische Universität München, 1992.
- [9] H.-J. Bungartz, T. Dornseifer, and C. Zenger. Tensor product approximation spaces for the efficient numerical solution of partial differential equations. In *Proc. Int. Workshop on Scientific Computations, Konya, 1996*. Nova Science Publishers, 1997.
- [10] H.-J. Bungartz, M. Griebel, D. Rösche, and

Table 6: Results for forest cover type data set using 6 and 10 attributes

	λ	testing correctness		
		overall	Ponderosa Pine	other class
6 dimensions				
level 1	0.0005	92.68	93.87	92.59
	0.0050	92.52	93.95	92.42
	0.0500	92.45	93.43	92.39
on evaluation set	0.0050	92.50	93.52	92.43
level 2	0.0001	93.34	92.08	93.42
	0.0010	93.20	92.30	93.25
	0.0100	92.31	88.95	92.52
on evaluation set	0.0010	93.19	91.73	93.28
level 3	0.0010	92.78	90.90	92.90
	0.0100	93.10	91.74	93.18
	0.1000	93.50	87.97	93.86
on evaluation set	0.0100	93.02	91.42	93.13
10 dimensions				
level 1	0.0025	93.64	94.03	93.62
	0.0250	93.56	94.19	93.52
	0.2500	93.64	92.30	93.72
on evaluation set	0.0250	93.53	93.81	93.51
level 2	0.0050	92.95	92.36	92.98
	0.0500	93.67	92.96	93.71
	0.5000	93.10	91.81	93.18
on evaluation set	0.0500	93.72	92.89	93.77

Table 7: Results for a 10D synthetic massive data set, $\lambda = 0.01$

	# of points	training correctness	testing correctness	total time (sec)	data matrix time (sec)	# of iterations
level 1	50 000	98.8	97.2	19	4	47
	500 000	97.6	97.4	104	49	50
	5 million	97.4	97.4	811	452	56
level 2	50 000	99.8	96.3	265	45	592
	500 000	98.6	97.8	1126	541	635
	5 million	97.9	97.9	7764	5330	688

C. Zenger. Pointwise convergence of the combination technique for the Laplace equation. *East-West J. Numer. Math.*, 2:21–45, 1994. also as SFB-Bericht 342/16/93A, Institut für Informatik, TU München, 1993.

- [11] V. Cherkassky and F. Mulier. *Learning from Data - Concepts, Theory and Methods*. John Wiley & Sons, 1998.
- [12] K. Cios, W. Pedrycz, and R. Swiniarski. *Data Mining Methods for Knowledge Discovery*. Kluwer, 1998.
- [13] T. G. Dietterich. Approximate statistical tests for comparing supervised classification learning algorithms. *Neural Computation*, 10(7):1895–1924, 1998.
- [14] K. Frank, S. Heinrich, and S. Pereverzev. Information Complexity of Multivariate Fredholm Integral Equations in Sobolev Classes. *J. of Complexity*, 12:17–34, 1996.
- [15] H. Freudenthal. Simplicialzerlegungen von beschränkter Flachheit. *Annals of Mathematics*, 43:580–582, 1942.
- [16] J. Garcke and M. Griebel. On the computation of the

eigenproblems of hydrogen and helium in strong magnetic and electric fields with the sparse grid combination technique. *Journal of Computational Physics*, 165(2):694–716, 2000. also as SFB 256 Preprint 670, Institut für Angewandte Mathematik, Universität Bonn, 2000.

- [17] J. Garcke and M. Griebel. On the parallelization of the sparse grid approach for data mining. SFB 256 Preprint 721, Universität Bonn, 2001. <http://wissrech.iam.uni-bonn.de/research/pub/garcke/psm.pdf>.
- [18] J. Garcke, M. Griebel, and M. Thess. Data mining with sparse grids. 2000. Submitted, also as SFB 256 Preprint 675, Institut für Angewandte Mathematik, Universität Bonn, 2000.
- [19] T. Gerstner and M. Griebel. Numerical Integration using Sparse Grids. *Numer. Algorithms*, 18:209–232, 1998. (also as SFB 256 preprint 553, Univ. Bonn, 1998).
- [20] F. Girosi. An equivalence between sparse approximation and support vector machines. *Neural Computation*, 10(6):1455–1480, 1998.

- [21] F. Girosi, M. Jones, and T. Poggio. Regularization theory and neural networks architectures. *Neural Computation*, 7:219–265, 1995.
- [22] G. Golub, M. Heath, and G. Wahba. Generalized cross validation as a method for choosing a good ridge parameter. *Technometrics*, 21:215–224, 1979.
- [23] M. Griebel. The combination technique for the sparse grid solution of PDEs on multiprocessor machines. *Parallel Processing Letters*, 2(1):61–70, 1992. also as SFB Bericht 342/14/91 A, Institut für Informatik, TU München, 1991.
- [24] M. Griebel. Adaptive sparse grid multilevel methods for elliptic PDEs based on finite differences. *Computing*, 61(2):151–179, 1998. also as Proceedings Large-Scale Scientific Computations of Engineering and Environmental Problems, 7. June - 11. June, 1997, Varna, Bulgaria, Notes on Numerical Fluid Mechanics 62, Vieweg-Verlag, Braunschweig, M. Griebel, O. Iliev, S. Margenov and P. Vassilevski (editors).
- [25] M. Griebel, W. Huber, T. Störtkuhl, and C. Zenger. On the parallel solution of 3D PDEs on a network of workstations and on vector computers. In A. Bode and M. D. Cin, editors, *Parallel Computer Architectures: Theory, Hardware, Software, Applications*, volume 732 of *Lecture Notes in Computer Science*, pages 276–291. Springer Verlag, 1993.
- [26] M. Griebel and S. Knappek. Optimized tensor-product approximation spaces. *Constructive Approximation*, 16(4):525–540, 2000.
- [27] M. Griebel, P. Oswald, and T. Schiekofer. Sparse grids for boundary integral equations. *Numer. Mathematik*, 83(2):279–312, 1999. also as SFB 256 report 554, Universität Bonn.
- [28] M. Griebel, M. Schneider, and C. Zenger. A combination technique for the solution of sparse grid problems. In P. de Groen and R. Beauwens, editors, *Iterative Methods in Linear Algebra*, pages 263–281. IMACS, Elsevier, North Holland, 1992. also as SFB Bericht, 342/19/90 A, Institut für Informatik, TU München, 1990.
- [29] M. Griebel and V. Thurner. The efficient solution of fluid dynamics problems by the combination technique. *Int. J. Num. Meth. for Heat and Fluid Flow*, 5(3):251–269, 1995. also as SFB Bericht 342/1/93 A, Institut für Informatik. TU München, 1993.
- [30] M. Hegland, O. M. Nielsen, and Z. Shen. High dimensional smoothing based on multilevel analysis. Technical report, Data Mining Group, The Australian National University, Canberra, November 2000. Submitted to SIAM J. Scientific Computing.
- [31] J. Hoschek and D. Lasser. *Grundlagen der geometrischen Datenverarbeitung*, chapter 9. Teubner, 1992.
- [32] H. W. Kuhn. Some combinatorial lemmas in topology. *IBM J. Res. Develop.*, 4:518–524, 1960.
- [33] Y. J. Lee and O. L. Mangasarian. SSVM: A smooth support vector machine for classification. *Computational Optimization and Applications*, 20(1), 2001. to appear.
- [34] G. Melli. Datgen: A program that creates structured data. Website. <http://www.datasetgenerator.com>.
- [35] W. D. Penny and S. J. Roberts. Bayesian neural networks for classification: how useful is the evidence framework ? *Neural Networks*, 12:877–892, 1999.
- [36] B. D. Ripley. Neural networks and related methods for classification. *Journal of the Royal Statistical Society B*, 56(3):409–456, 1994.
- [37] S. L. Salzberg. On comparing classifiers: Pitfalls to avoid and a recommended approach. *Data Mining and Knowledge Discovery*, 1:317–327, 1997.
- [38] T. Schiekofer. *Die Methode der Finiten Differenzen auf dünnen Gittern zur Lösung elliptischer und parabolischer partieller Differentialgleichungen*. Doktorarbeit, Institut für Angewandte Mathematik, Universität Bonn, 1999.
- [39] W. Sickel and F. Sprengel. Interpolation on sparse grids and Nikol'skij-Besov spaces of dominating mixed smoothness. *J. Comput. Anal. Appl.*, 1:263–288, 1999.
- [40] S. Singh. 2d spiral pattern recognition with possibilistic measures. *Pattern Recognition Letters*, 19(2):141–147, 1998.
- [41] S. A. Smolyak. Quadrature and interpolation formulas for tensor products of certain classes of functions. *Dokl. Akad. Nauk SSSR*, 148:1042–1043, 1963. Russian, Engl. Transl.: Soviet Math. Dokl. 4:240–243, 1963.
- [42] V. N. Temlyakov. Approximation of functions with bounded mixed derivative. *Proc. Steklov Inst. Math.*, 1, 1989.
- [43] A. N. Tikhonov and V. A. Arsenin. *Solutios of ill-posed problems*. W.H. Winston, Washington D.C., 1977.
- [44] F. Utreras. Cross-validation techniques for smoothing spline functions in one or two dimensions. In T. Gasser and M. Rosenblatt, editors, *Smoothing techniques for curve estimation*, pages 196–231. Springer-Verlag, Heidelberg, 1979.
- [45] V. N. Vapnik. *Estimation of dependences based on empirical data*. Springer-Verlag, Berlin, 1982.
- [46] V. N. Vapnik. *The Nature of Statistical Learning Theory*. Springer, 1995.
- [47] G. Wahba. *Spline models for observational data*, volume 59 of *Series in Applied Mathematics*. SIAM, Philadelphia, 1990.
- [48] A. Wieland. Spiral data set. <http://www.cs.cmu.edu/afs/cs.cmu.edu/project/ai-repository/ai/areas/neural/bench/cmu/0.html>.
- [49] C. Zenger. Sparse grids. In W. Hackbusch, editor, *Parallel Algorithms for Partial Differential Equations, Proceedings of the Sixth GAMM-Seminar, Kiel, 1990*, volume 31 of *Notes on Num. Fluid Mech.* Vieweg-Verlag, 1991.