

Data Mining with Sparse Grids

Seminar Computational Aspects of Machine Learning

Sebastian Kreisel

Department for Informatics

Technische Universität München

Email: sebastian.kreisel@tum.de

Zusammenfassung—516.0pt**TODO:** Abstract **TODO:** Ab-
abstract **TODO:** Abstract **TODO:** Abstract **TODO:** Abstract
TODO: Abstract **TODO:** Abstract **TODO:** Abstract **TODO:**
Abstract **TODO:** Abstract **TODO:** Abstract **TODO:** Abstract
TODO: Abstract **TODO:** Abstract **TODO:** Abstract **TODO:**
Abstract **TODO:** Abstract **TODO:** Abstract **TODO:** Abstract
TODO: Abstract **TODO:** Abstract **TODO:** Abstract **TODO:**
Abstract **TODO:** Abstract **TODO:** Abstract **TODO:** Abstract
TODO: Abstract **TODO:** Abstract **TODO:** Abstract **TODO:**

Index Terms—Sparse grids; Data mining; Hierarchical discretization; Curse of dimensionality

I. INTRODUCTION

Large datasets and high dimensional data remain challenging aspects of data mining. Even with growing computational power, many problems require specialized algorithms to archive accurate results within the given time and cost restraints.

Sparse grids belong to a more general class of *grid-based* discretization methods. These methods are primarily applied to tackle scenarios with large amount of data points and high-dimensional feature spaces, posing these problems:

Often, algorithms scale quadratic or worse in the number of data points and thus quickly leading to time and cost related issues. High dimensionality introduces a problem widely known as the *Curse of Dimensionality*, denoting an exponential dependency between computational effort and the number of dimensions of the data.

By focusing on grid points instead of the data points themselves grid-based methods allow for better handling of large amounts of data. Sparse grids specifically combat the curse of dimensionality and mitigate the exponential dependency. Note also, that grid-based approaches are not applicable to data mining exclusively, but are also suited for a number of different areas including PDA, model order reduction or numerical quadrature.

In this report the sparse grid technique is applied to data mining, investigating the mitigation-capabilities to the previously mentioned issues. First this report introduces grid discretization and sparse grids in general as well as related topics like spatial adaptivity.

Then, sparse grids will be applied machine learning through *least square estimation*. To confirm the capabilities of sparse grids, the results of employing difficult, test-datasets for regression and classification (i.e. checker-board dataset) will

be shown.

Lastly efficient implementations on modern systems and parallelization for sparse grids will be examined.

II. GRID DISCRETIZATION

In machine learning, algorithms usually focus on a given training dataset X , for instance

$$X = \{x^{(i)} \mid x^{(i)} \in [0, 1]^d\}_{i=1}^M, \quad Y = \{y^{(i)} \mid y^{(i)} \in \mathbb{R}\}_{i=1}^M$$

with, in case of supervised learning, an associated solution set Y .

Grid-based approaches introduce an additional set G of N *grid points* with

$$G = \{1, 2, \dots, N\} \text{ .}$$

For each dimension of the feature space a separate G (with possibly different N) is constructed dividing the space into a grid. This, by the grid *discretized*, space will then be used instead of working with the datapoints in the original feature space directly.

A. Full grid discretization

In the following functions will be restricted to the unit hypercube

$$f : [0, 1]^d \rightarrow \mathbb{R} .$$

To construct a *full* grid we chose the grid points G equidistant, without grid points lying on the borders.

We first consider the case of a one-dimensional f being discretized. Around each gridpoint i we center a one-dimensional *basis function*

$$\phi_i(x) = \max\{0, 1 - |(N+1)x - i|\} \text{ .}$$

$\phi_i(x)$ is a standard hat function centered around i and dilated to have local support between the grid points $i - 1$ and $i + 1$. Fig. 1 shows $G = \{1, 2, \dots, 7\}$ and the related basis-functions.

To discretize a function $f(x)$ we introduce a coefficient (surplus) α_i for each grid point i . This coefficient is defined to be f evaluated at the grid point i

$$\alpha_i = f\left(\frac{i}{N+1}\right) .$$

Taking the sum

$$f(x) \approx \hat{f}(x) = \sum_{i \in G} \alpha_i \phi_i(x)$$

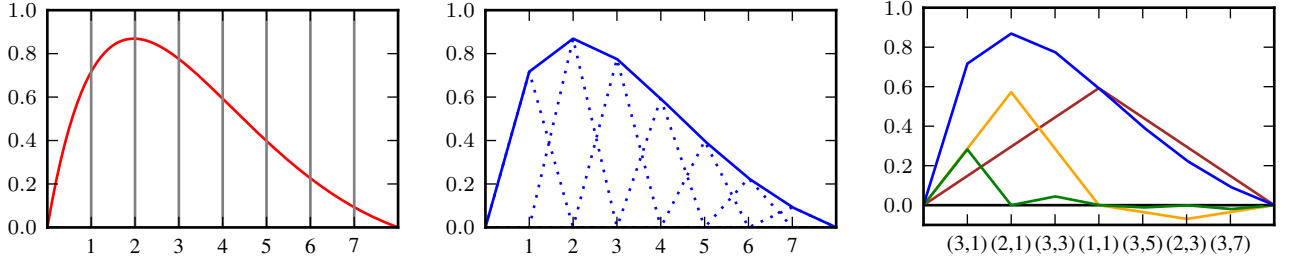


Abbildung 1: A function f , red, to be discretized. Seven grid points discretizing the space (left). Full grid discretization by nodal basis (middle) and hierarchical basis, scaled by α_i and $\alpha_{l,i}$ respectively (right).

over all weighted basis-functions ϕ_i discretizes (approximates) f . Fig. 1 illustrates this.

For $f(\vec{x})$ with $d > 1$, grid point representation is extended to a d -tuple of indices, i.e. $(1, 3, 1)$ denoting the grid point with position $x = 1$, $y = 3$, $z = 3$ in the dimensions x, y, z . The related basis function

$$\phi_i(\vec{x}) = \prod_{j=1}^d \phi_{i,j}(x_j)$$

gets extended to d dimensions using the tensor product over the previously defined one-dimensional hat functions $\phi_{i,j}(x_j)$ with x_j being the j -th element of \vec{x} and $\phi_{i,j}$ denoting the basis-function of grid point i in the dimension j . To improve readability the dimension-related index j of $\phi_{i,j}$ will be omitted in the following.

B. Hierarchical basis

Besides constructing the grid in the simple way described in Sec. II-A, more sophisticated methods are available. In order to make a grid sparse and still keep a sufficient accuracy the following *hierarchical basis* is introduced.

We first examine the case $d = 1$. Let $l \in \{1, 2, \dots\}$ be the *level* with $|G| = 2^{(l-1)}$ associated grid points on each level. This level hierarchy groups grid points into sets

$$G_l = \{i \in \mathbb{N} \mid 1 \leq i \leq 2^l, i \text{ odd}\},$$

omitting every second grid point. Together the adjusted hat function

$$\phi_{l,i}(x) = \max\{0, 1 - |2^l x - i|\}$$

this forms the hierarchical basis in one dimension up to a level n . By disregarding every even grid point the local supports of basis functions on the *same* level are mutually exclusive and for each value of x exactly one basis function is not zero.

A grid point is now referred to as grid point in the level l with index i . Note, that the indices alone do not uniquely define a grid point (i.e. grid points of index $i = 1$ are element in every G_l). The same applies to the corresponding ϕ and α .

Taking the weighted sum over all levels and all grid points in one dimension

$$f(x) \approx \hat{f}(x) = \sum_{l \leq n, i \in G_l} \alpha_{l,i} \phi_{l,i}(x)$$

discretizes f on a full grid.

In contrast to the conventional approach from Sec. II-A the surpluses now are calculated differently. Let $x_{l,i}$ be the x -value of the grid point given by l and i . Then

$$\alpha_{l,i} = f(x_{l,i}) - \frac{f(x_{l,i} - 2^{-l}) + f(x_{l,i} + 2^{-l})}{2}$$

is the *hierarchical surplus* for the grid point (l, i) . The function value at the grid point is taken and the function values at neighbouring grid points are subtracted (“neighbouring” is disregarding l). For instance, for $\alpha_{2,1}$ we get $x_{2,1} = \frac{1}{4}$ and

$$\alpha_{2,1} = f\left(\frac{1}{4}\right) - \frac{f\left(\frac{1}{4} - 2^{-2}\right) + f\left(\frac{1}{4} + 2^{-2}\right)}{2}.$$

For $d > 1$ we combine the one dimensional basis functions to d -dimensional basis functions using the tensor product, analogous to Sec. II-A. This is done for all possible combinations of l and i in all dimensions. This process of building d -dimensional basis functions through combining over the level in different dimensions leads to a subspaces defined by the level-vector $\vec{l} = l_x, l_y, \dots$ as shown for $d = 2$ in Fig. 2. Taking the sum $\sum_{l,i} \alpha_{l,i} \phi_{l,i}(x)$ over all

$$\phi_{\vec{l},i}(\vec{x}) = \prod_j \phi_{l_j,i_j}(x_j)$$

discretizes $f(\vec{x})$ on a hierarchical structured grid.

However, this does not lead to a sparse grid immediately. So far the gridpoints only got regrouped and for a the maximum level n this results in $|G| = 2^n - 1$ basis functions for each dimension. This further leads to an exponential dependency of the number of grid points and d , thus having no effect on mitigating the curse of dimensionality.

C. Sparse grid discretization

In order to make the hierarchical grid *sparse*, we now disregard certain subspaces with their associated grid points. The goal is to reduced the total number of grid points by finding and disregarding those that contribute the least to the discretization of f . Which gridpoints that are is a *a-priori* solvable optimization problem. Thus, independent of f all $\phi_{l,i}$

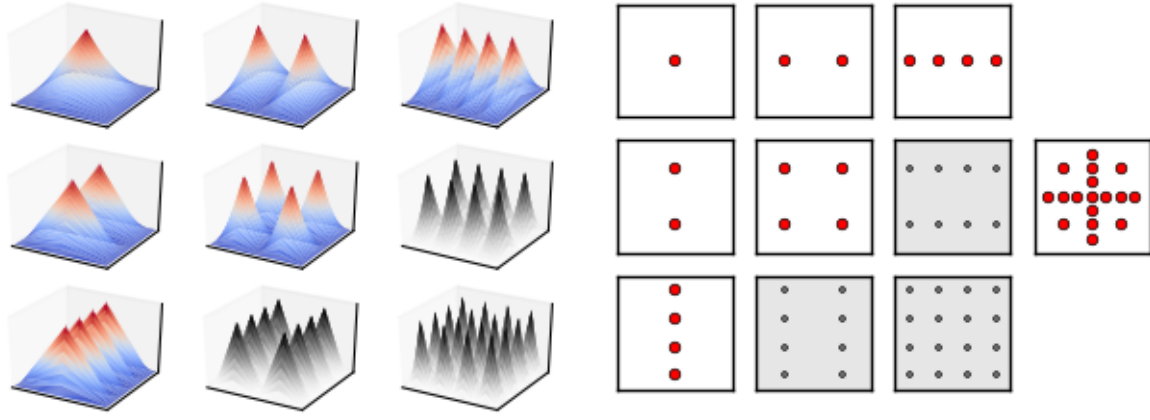


Abbildung 2: Hierarchical subspaces (left) and corresponding grid points (right) ($d = 2$) with the last column being the complete sparse grid. For the sparse grid omitted subspaces/grid points in grey.

related to the subspaces in the lower right of the diagonal in Fig. 2 will be left out of the sum

$$\hat{f}(x) = \sum_{l \leq n, i \in G_l} \alpha_i \phi_{l,i}(x)$$

from Sec. II-A.

By doing so, we reduced the total number of grid points drastically from $\mathcal{O}(2^{nd})$ to $\mathcal{O}(2^n \cdot n^{d-1})$ where n is the maximal level in the hierarchical structure. The asymptotic error on the other hand only slightly increases from $\mathcal{O}(2^{-2n})$ to $\mathcal{O}(2^{-2n} \cdot n^{d-1})$. Tab. I illustrates the quickly growing gap between the number of grid points in a full and sparse grid.

d	1	2	3	5	10	20
Full	15	225	3375	$> 10^5$	$> 10^{11}$	$> 10^{23}$
Sparse	15	49	111	351	2001	13201

Tabelle I: Number of grid points in a full and sparse grid (without points on the boundaries) with maximal level $n = 4$ and growing dimension d .

It is important to note that the numbers in Tab. I are taken for a sparse grid without points on the boundaries. In case the function to be discretized is not zero on the boundaries, such points become necessary. Treatment of the boundaries might require considerably more grid points. For further information please refer to REF.

D. Adaptive sparse grids

Even though sparse grid offer a optimal choice for a general function f the discretization error is in praxis often unacceptable due to the properties of f itself. If f exhibits steep, complex or discontinuous areas, which the *a-priori* distribution of grid points can not capture, additional grid points have to be added locally. Two different approaches are possible:

III. SPARSE GRIDS IN MACHINE LEARNING

- Quick note on classification/regression
- Least squares
- Least square with sparse grids
- Matrix formulation
- Notes on matrix solving etc.

IV. SOMETHING SOMETHING IMPLEMENATION

-
-

V. CONCLUSION

LITERATUR

- [1] D. Pflüger, "Spatially adaptive sparse grids for high-dimensional problems," Ph.D. dissertation, Technische Universität München, 2010.
- [2] A. Heinecke, "Boosting scientific computing applications through leveraging data parallel architectures," Ph.D. dissertation, Technische Universität München, 2013.
- [3] B. Peherstorfer, "Model order reduction of parametrized systems with sparse grid learning techniques," Ph.D. dissertation, Technische Universität München, 2013.
- [4] H. Bungartz and M. Griebel, "Sparse grids," *Acta Numerica*, pp. 1–43, 2004.
- [5] H. Bungartz, D. Pflüger, and S. Zimmer, "Adaptive sparse grid techniques for data mining," *Modelling, Simulation and Optimization of Complex Processes*, 2006.