

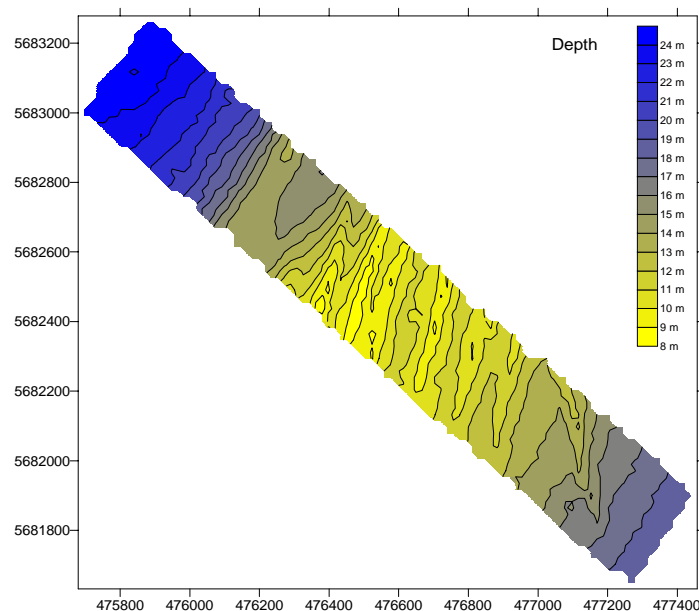
# MAPPING WITH MULTIBEAM DATA : ARE THERE IDEAL MODEL SETTINGS ?

Vande Wiele Tom  
University of Ghent  
Geography Department  
Krijgslaan 281 – S8  
9000 Gent, Belgium  
[tom.vandewiele@rug.ac.be](mailto:tom.vandewiele@rug.ac.be)

## 1. Introduction.

In shallow waters, like the Belgian continental shelf, a high density of data is available from a multibeam survey. Processing this huge amount of data in order to obtain maps is thus a time consuming task. But do we need all this data to make acceptable maps ?

For answering this question a case study has been made. The selected area of study is the central part of the Kwintebank, part of the Flemisch Banks, in front of the Belgian coast. This sandbank is characterized by his steepest slope faced north and the presence of sand waves. Depth values are varying between –25 m and –5 m.



**Figure 1 - Depth map.**

[Data available from the Belgian Geological Survey.]

## 2. Methodology.

A hydrographic survey with a multibeam echosounder consists out of irregularly spaced data. As most contouring algorithms are based on a regular grid, a local interpolator (model) is used to construct this grid. Our research is about the settings of this model and the behaviour of the model in relation to the different data density patterns. This means that we need a method to compare the results of models with different settings, and models based on different data sets.

The method we will use to calculate the model performance is given in the following diagram (fig. 2). Out of the available data we select two different types of data sets. The first one, called the construction set, is used by the model to make the estimates, while the second data set, the test set or validation set, is used to validate the model (fig. 3). With different data densities, several construction sets are available.

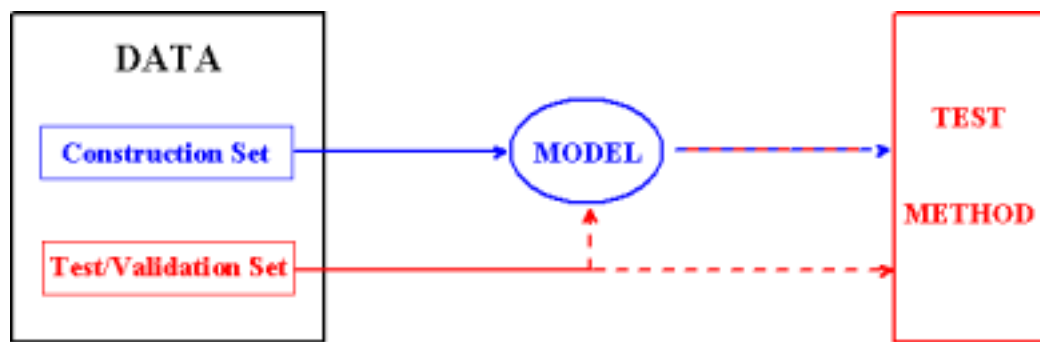


Figure 2 - Method – construction set/validation set.

The test or validation set forms the basis for the test method. This test method will compare the values predicted by the model ( $Z_p$ ) with the observed values in the validation set ( $Z_o$ ). The predicted values are acquired by importing the  $x$ ,  $y$  - coordinates into the model.

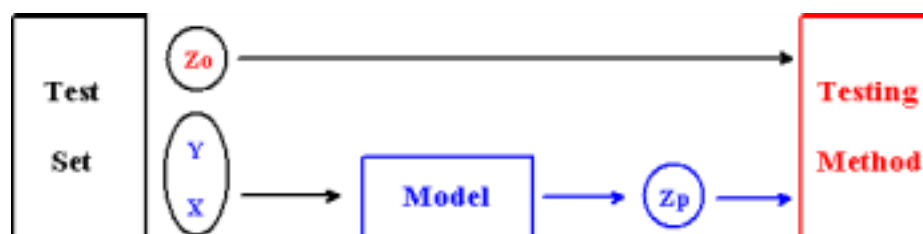


Figure 3 - Method – use of the test set.

The test method itself exists out of three main topics : general descriptive coefficients, correlation coefficients and graphical techniques [VANDE WIELE T, VERNEMMEN. C -

1999] [WILLMOTT, C. 1984]. In this paper we will restrict ourselves to use only a few of these parameters or methods.

As the main purpose of our research is to find out how much data we need to construct acceptable maps, we need at least a definition of an acceptable map. To do this we will base ourselves on an average error expressed as the Root Mean Square Error (RMSE).

### 3. The model.

The model we will use for the gridding process is a simple local interpolator. For estimating the value of a certain point we can make use of global as well as local models. As global models make use of all data from the area of study, local models only make use of the data in a small neighbourhood. The model we consider in our research is called the inverse distance (ID) model and is a local linear interpolator.

Meaning that the estimate exists out of a linear combination of local surrounding values ( $z_j$ ), each with a certain weighing factor ( $w_j$ ).

$$Z_{IP} = \sum_{j=1}^n w_j z_j$$

The determination of these weighing factors leads to different models. For the inverse distance model the weighing factor is determined by the distance from the data point to the estimated point ( $d_{IPj}$ ). To make nearby data points more important (meaning giving a greater weight) an inverse proportional function is used, based on the distance to a certain power. (DAVIS, J.C. 1973)

To meet the unbiasedness condition the sum of all weights has to be one [ISAAKS, E.H. and SRIVASTAVA, R.M. 1989].

$$\sum_{i=1}^n w_i = 1$$

Meaning that the weighing function will look like :

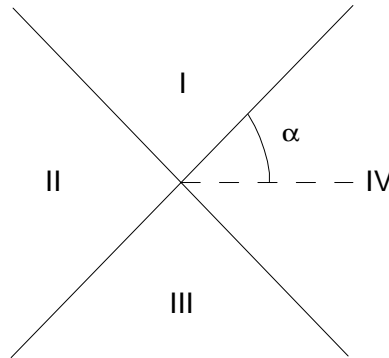
$$w_j = \frac{d_{IPj}^{-p}}{\sum_{j=1}^n d_{IPj}^{-p}}$$

and the model is described by the following formula :

$$Z_{IP} = \frac{\sum_{j=1}^n z_j d_{IPj}^{-p}}{\sum_{j=1}^n d_{IPj}^{-p}}$$

This brings us to the model settings. As can be derived from the basic formula, a number of possibilities are available like the number of points used for the estimate ( $n$ ) and the power of the weighing function ( $p$ ).

Concerning the data points for the interpolation not only the number of data points is important, but also the spreading of these data points in space. It is generally accepted that using data points surrounding the estimating point leads to a better estimate than clustered data points. Therefor the following search structures are provided : a simple, a quadrant and an octant search. In the first case the  $n$  nearest data points are selected, in the second case these  $n$  points are divided over four quadrants (meaning  $n/4$  data points in each quadrant) and in the last case a further refinement is made by a division into eight octants. In this way a better spreading of the data points is achieved. Besides a subdivision into sectors, the orientation of these sectors ( $\alpha$ ) can be important.



**Figure 4 - Quadrant search structure.**

The last setting concerns the anisotropy ratio of the area. The orientation as well as the anisotropy degree are possible settings. [TOMCZAK, M. 1998]

#### 4. Analyses and results.

The first three graphs (fig. 5, 6 and 7) are showing the results (RMSE) for the following model settings : the three main clusters along the x-axis (each existing out of three blue lines and three red lines) are differing in weighing functions. The first cluster (closest to the y-axis) is characterized by a power of one, the second cluster has a power of two while the third cluster has a power of three. With each cluster existing out of six lines there are three blue lines for the isotropic settings and three red lines for the anisotropic settings (with an anisotropy ratio of two). As can be seen in the graph legend, each line represents a different search structure, varying from simple over quadrant to octant. The different points on each line represent the varying number of data points used for the interpolation, ranging from 4 over 8 and 16 to 32. The first graph (fig. 5) gives the result for a point density of only  $0.4/100 \text{ m}^2$ , the second graph (fig. 6) represents a density of  $4.2/100 \text{ m}^2$  and the last graph  $8.5/100 \text{ m}^2$  (fig. 7).

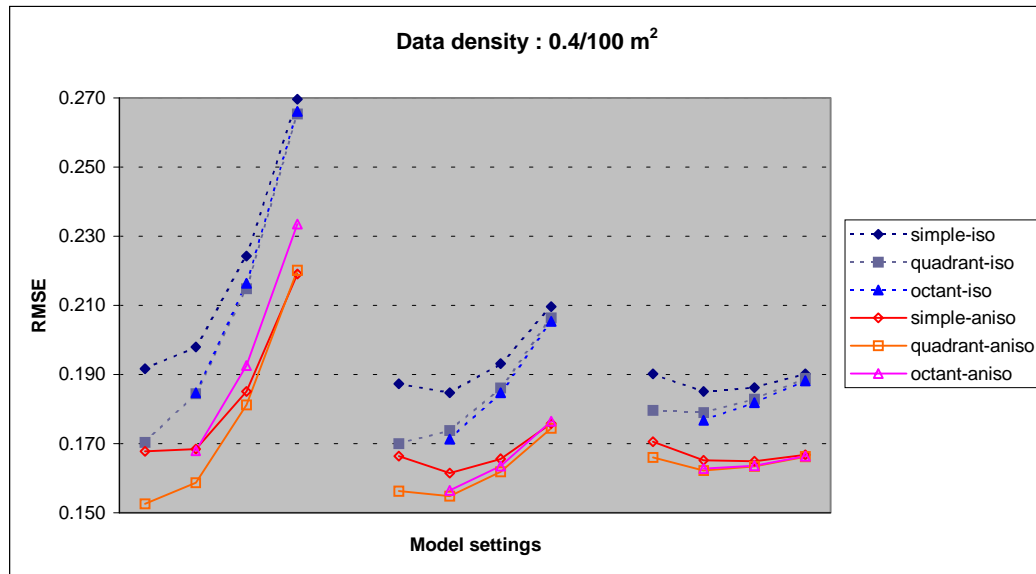


Figure 5 - Model results - density  $0.4/100 \text{ m}^2$ .

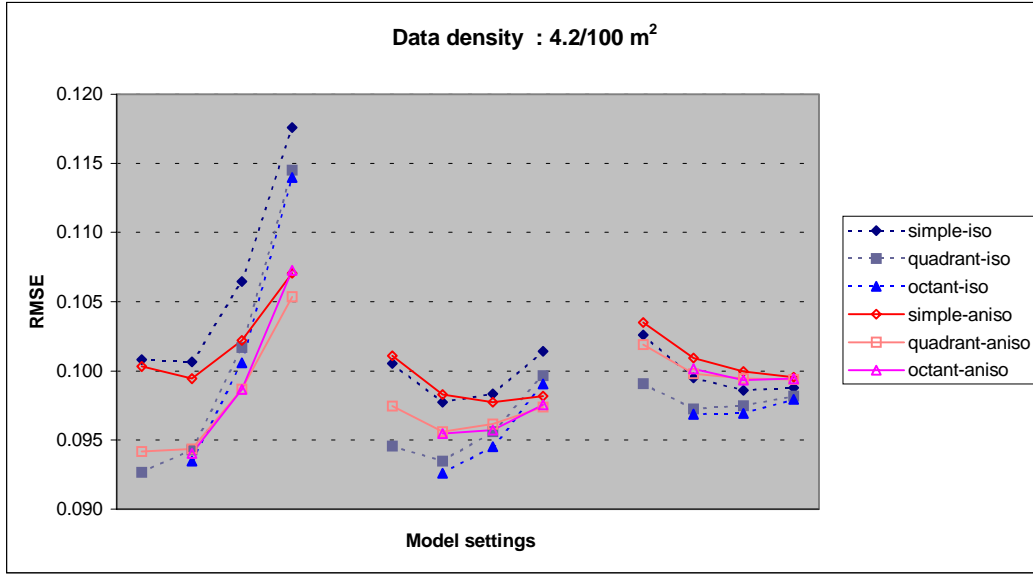


Figure 6 - Model results - density 4.2/100 m<sup>2</sup>.

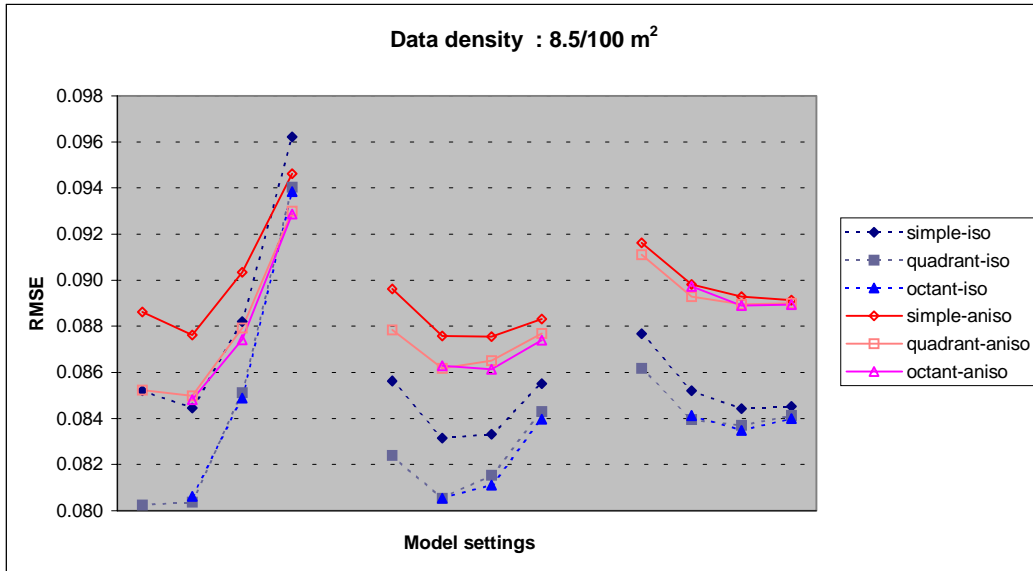


Figure 7 - Model results - density 8.5/100 m<sup>2</sup>.

Based on these diagrams the following observations can be made :

The anisotropy setting has a positive influence (lower RMSE) at low data density but changes to a negative influence at high data density. An other general fact concerns the search structure. Using the quadrant or octant search structure leads to better estimates compared to a simple search, with a small mutual difference between the quadrant and the octant search. When using a weighing function with a power of one, we notice that the more data is used for the interpolation the more the model performance decreases. Remarkable is the big difference between the best and worst performer in this cluster.

With a more pronounced weighing function (power equals two) the influence of the number of data points used for the interpolation changes. The best performer is now characterized by the use of slightly more data point (eight to be precise), while the tendency of a performance decrease with the use of many data points stays.

Using a weighing function with a power of three provides us with still an other influence of the number of data points used for the estimate. Here we see that using more data points tends to give a better result. This cluster of models seems to perform worse than a less pronounced weighing function, especially with a higher data density.

Defining the 'best or ideal' model settings as the model with the lowest RMSE, we must conclude that the quadrant search method with one point in each quadrant and a weighing function with a power of one does the job. For low data densities an additional setting of the anisotropy ratio further improves the model performance.

The next step in our research is to find out why the model behaves like this. Lets start with some basic thoughts. Making a good estimate is all about picking out the 'right' data points and involve them in the estimating process in an adequate manner. Some common sense will help us to achieve this goal. First of all, it is important to select data points around the estimating point, possible by imposing a quadrant or octant search structure. Secondly, the closer a data point is the more important it gets. Expressing the importance of this closest point is achieved by means of the weighing function and the anisotropy ratio. In certain cases the two theses, translated into model settings, could be contradictory.

Leaving us still with a model setting like the number of data points used for the interpolation. For example : using more data for the estimate gives us a (probably) better spreading, but it also implies points farther away. Other important factors are of course the relief itself and the data density.

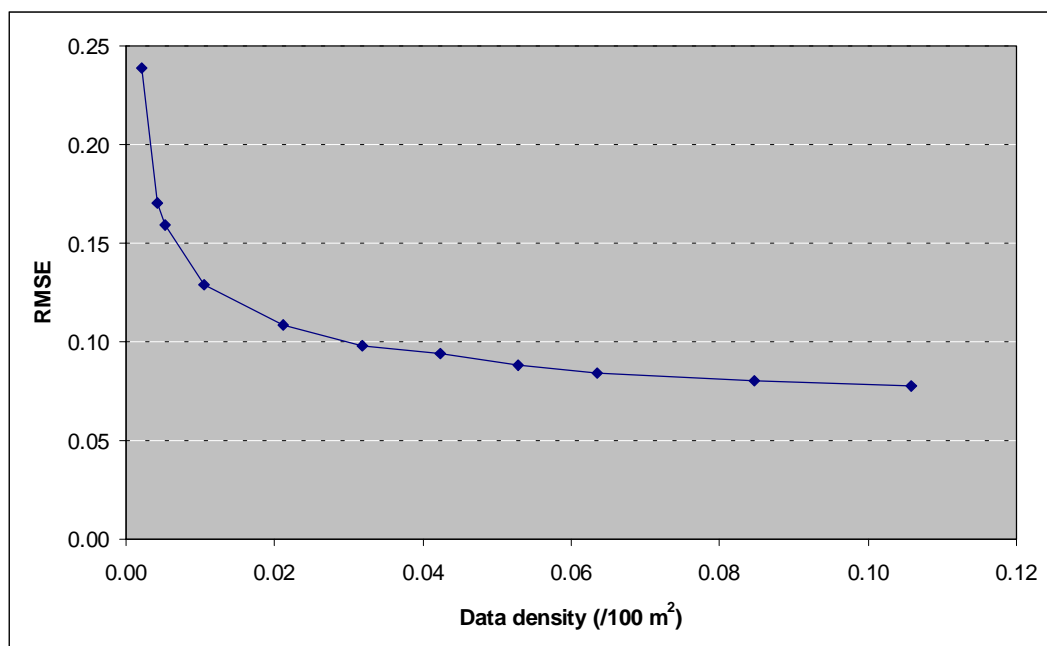
From the preceding graphs we can see that an interpolation based on surrounding data points (imposed by a quadrant or octant search structure) results in better estimates. As already mentioned the influences of the other settings are not only interrelated, but also related to the data density. A more pronounced weighing function implies that the closest data point gains importance. Making one point more important (and thus also one sector) causes a spreading impoverishment. Compensating this effect is possible by the use of more data, making the nearest point again less important. This effect is clearly shown for the weighing function based on a power of three and to a lesser degree for a power of two. Using more data points also implies longer distances, meaning data containing a lower amount of relevant information, explaining the RMSE rise for less pronounced weighing functions.

With a high density of data points (e.g.  $\approx 400$  points in an area of 100 m x 100 m) and sand waves characterized by a cross-section of about 100 m and a slope of roughly 10 %, the relief

can be approximated by a linear function. A high density makes it possible to detect small local variations, explaining the poor performance of the anisotropy models for these high densities. The setting of the anisotropy ratio is based on the general orientation of the sand waves. With low data densities the anisotropy setting is necessary to be sure that the used data describes the same structure (in this case study the sand waves). For high data densities, the closest data points always meet this requirement. But it is also possible to describe even smaller variations, not characterized by an anisotropy ratio, explaining the lower model performance of the anisotropy models (as a ‘wrong’ set anisotropy ratio causes a spreading impoverishment).

Fig. 8 shows the influence of the data density for the described ‘ideal’ model. Like we expected, the model performance increases with a higher data density, in a logarithmic manner. As more data is available, smaller (or more local) variations can be mapped, reducing the RMSE.

Using the RMSE for expressing map accuracy we only need a limit value to derive the minimum necessary data density. For example : the limit value for a 95 % probability interval of  $\pm 20$  cm is a RMSE of 0.1; for an interval of  $\pm 15$  cm the limit value amounts to 0.075.



**Figure 8 - Influence data density.**



## **5. Conclusion.**

Selecting the ideal model settings is not easy. It implies a good knowledge of the model itself and how it reacts on different settings, but it also implies a certain knowledge about the relief. The scale and the nature of the different structures play an important role in the mapping process. In this case study three different structures or variations are available, namely the general bank structure, the sand waves and the local variations on the sand waves, all differing in scale, with the first two structures clearly anisotropic (each in a different direction) and the last variation as a local isotropic component. Mapping these differing structures depends on the availability of data. If sufficient data is available for mapping a certain structure, it also implies that the information of the lower scale (or bigger) structure is included in the data. Making the model settings depending on the data density and the scale of the structures.

## **6. References.**

- DAVIS, J.C. 1973. Statistics and data analysis in geology. John Wiley & Sons, inc. New York. (p 549)
- ISAACS, E.H. and SRIVASTAVA, R.M. 1989. An Introduction to Applied Geostatistics. Oxford University Press, New York.
- TOMCZAK, M. 1998. Spatial Interpolation and its Uncertainty Using Automated Anisotropic Inverse Distance Weighting (IDW) - Cross-Validation/Jackknife Approach. Journal of Geographic Information and Decision Analysis, vol. 2, no. 2, pp. 18-33.
- VANDE WIELE, T. & VERNEMMEN, C. 1999. The influence of Different Point Patterns on Relief Models Applied to Sandbanks : General Construction and First Results. The Hydrographic Journal. No 91. January 1999. (p 21 - 27)
- WILLMOTT, C. 1984. On the evaluation of model performance in physical geography. Spatial statistics and models. D. Reidel publishing company, p 443-460.