# Data-driven probabilistic predictions of sand wave bathymetry

J.C. Wüst

North Sea Directorate, P.O. Box 5807, 2280 HV, Rijswijk, the Netherlands

Email: j.c.wust@dnz.rws.minvenw.nl

## Abstract

A state space model assembled of spatially distributed local linear growth models, serves as a reduced dimension representation of the sea floor and its dynamics. The model state is mapped to continuous output space through an interpolation kernel. Using a time series of gridded multi-beam data, probabilistic model state estimates are recursively made through linear Kalman filtering. The model results demonstrate the utility of this approach. The flexible state space framework offers ample room for improvement. The probabilistic bathymetry predictions of the model will be used to support the channel maintenance decision process of the North Sea Directorate.

## 1. Introduction

One of the tasks of the North Sea Directorate is to keep the Euro channel to the port of Rotterdam and its approaches safely accessible for deep drafted ships. Figure 1 shows the "maintained depth" areas.
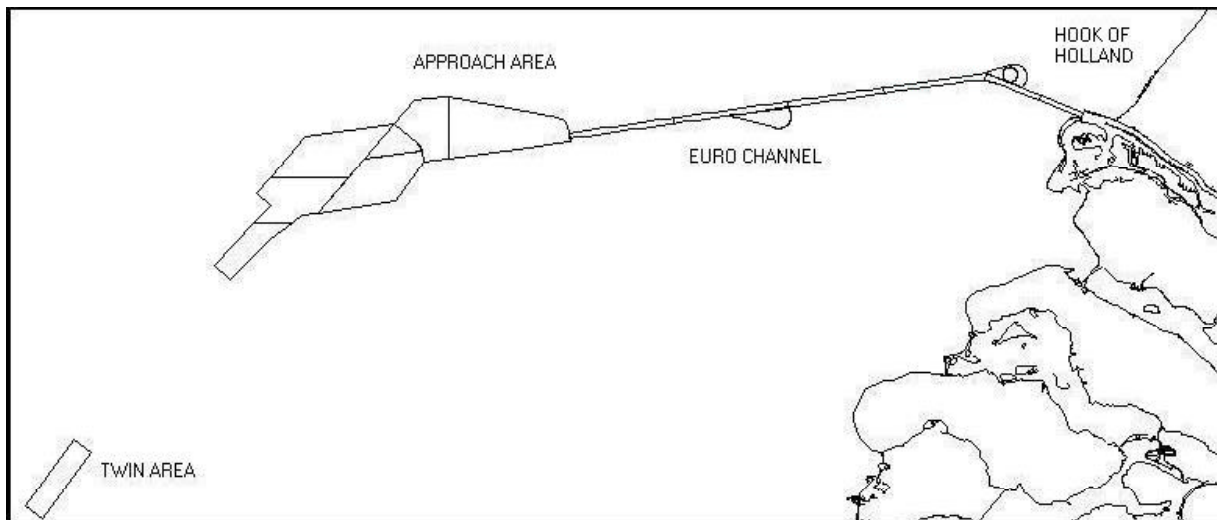


**Figure 1: The Euro channel and its approaches**

For the quantitative justification of the long term scheduling of dredging and surveying activities, predictions of sand wave evolution are needed including reliable uncertainty estimates. This paper describes an attempt to achieve this goal using linear growth models for inferring probabilistic predictive sea floor information from a time series of multi-beam surveys.

## 2. The model

The choice of modelling principle is inspired by the flexible class of convolution method type space-time models [Stroud et al., Lee et al., 2001]. The basic type of model convolves a white noise process with a convolution kernel. The model class allows convolution with many types of random process, like e.g. regression surfaces or temporally dynamic processes.

## 2.1 The model principle

The model consists of a set of spatially distributed support points. To each support point a linear growth model is attributed which is of linear state space form. This is a spatial extension of the linear growth model of [West and Harrison, 1997]. The local models are combined in a joint state space model allowing state predictions through temporal state evolution. The state consists of the local level and trend values of the support points and a single element representing the average depth of the analysed sea floor area. A measurement equation based on kernel interpolation, maps the local level elements of the model state to continuous measurement space. The model states are recursively estimated on the basis of measurement data using linear Kalman filtering [West and Harrison, 1997], the state estimates being described by a Gaussian (or normal) distribution.

## 2.2 Model specification

Application of the state space Kalman filtering framework requires specifying the following equations.

1) The system equation $\mathbf{q}_{t|t-1} = \mathbf{G}\mathbf{q}_{t-1|t-1} + \mathbf{w}_t$,
   with system matrix $\mathbf{G}$ describing the time evolution of the state vector $\mathbf{q}$. $\mathbf{q}_{t|t-1}$ is the mean vector of the time t (Gaussian) state estimate, based on information up to and including time t-1, with state variance matrix $\mathbf{P}_{t|t-1}$. The $\mathbf{w}_t$ vector represents the state evolution error when the state evolves from time t-1 to time t. $\mathbf{w}_t$ is assumed to have zero mean and variance matrix $\mathbf{W}_t$, implying that the variance of the state estimate evolves as $\mathbf{P}_{t|t-1} = \mathbf{G}\mathbf{P}_{t-1|t-1}\mathbf{G}^T + \mathbf{W}_t$.

2) The measurement equation $\mathbf{y}_t = \mathbf{F}_t\mathbf{q}_{t|t-1} + \mathbf{v}_t$,
   with $\mathbf{y}_t$ the measurement vector, $\mathbf{F}_t$ the deterministic kernel interpolation matrix and $\mathbf{v}_t$ the time t zero mean measurement noise vector, with measurement noise matrix $\mathbf{V}_t$ being independent of $\mathbf{W}_t$.

During the time t Kalman filter measurement update step, $\mathbf{q}_{t|t-1}$ and $\mathbf{P}_{t|t-1}$ are updated to $\mathbf{q}_{t|t}$ and $\mathbf{P}_{t|t}$. For starting this recursive updating process an initial state estimate must be specified having mean $\mathbf{q}_{0|0}$ and variance matrix $\mathbf{P}_{0|0}$, $\mathbf{P}_{0|0}$ being independent of $\mathbf{W}_t$ and $\mathbf{V}_t$. Since $\mathbf{W}_t$, $\mathbf{V}_t$ and $\mathbf{P}_{0|0}$ are assumed Gaussian, due to the linearity of the model all future states and model predictions will be Gaussian as well.

The time-independent system matrix $\mathbf{G}$ represents the state dynamics. $\mathbf{G}$ consists of an identity matrix with a few non-zero off-diagonal elements representing the linear growth dynamics of the local level state elements as governed by the values of the local trend state elements. This implies that each $\mathbf{G}$-row representing a local level state element has a one added at the corresponding trend element position. The local trend and average area depth state elements are assigned Gaussian random walk dynamics. This implies that they deterministically evolve unchanged through the identity elements of the $\mathbf{G}$-matrix. They are only adapted stochastically through the Kalman filter measurement updates.

Since there is little a priori knowledge of the structure of the system evolution variance matrix $\mathbf{W}_t$, the discounting method is adopted [West and Harrison, 1997]. According to this proved method the system evolution variance is assessed as $\mathbf{W}_t = \mathbf{G}\mathbf{P}_{t-1|t-1}\mathbf{G}^T(1-\delta)/\delta$, with discount factor $\delta$ ($0 < \delta \leq 1$). $\mathbf{W}_t$ is also used for state evolution time steps beyond time t.

Because of their possibly different dynamical nature the level and trend state elements are assigned distinct discount factors. The area average depth state element is assigned an explicit evolution variance.

The measurement noise $\mathbf{v}_t$ comprises many sources of inaccuracy and may show significant spatial and temporal variation. Some recognised sources are e.g. positioning errors, affected by local sea floor slope effects, measurement equipment calibration and sea level estimation. For this first modelling attempt the simplifying assumption is made to model the measurement noise as being spatially and temporally white and stationary. Since being unknown, its estimation is incorporated in the calibration process. Future analysis of model performance will learn what measurement noise elements deserve special attention.

### 2.3 Increasing computational efficiency

The computational burden of Kalman filtering easily becomes prohibitive with increasing measurement and state dimension. Among many possibilities for reducing this burden the following two methods are applied as a first step.

1) Due to the limited range of the kernels the measurement matrix F contains a significant number of negligible contributions of distant support points which may be set to zero, allowing the computationally efficient application of sparse matrix techniques. A cut-off distance of 3.5 times the kernel standard deviation is applied.

2) Processing the measurements in blocks instead of in one go, reduces computation times and memory demands during the Kalman filter measurement update step [Anderson and Moore, 1979].

## 3. The data

The available data consist of a time series of multi-beam surveys as collected yearly by the North Sea Directorate. For most areas multi-beam data are available from 1991 onwards. The data are gridded to a 5x5 m grid. The model processes the grid cell average depths and their corresponding standard deviations. These spatially varying standard deviations range from 1 to 5 dm and this range remains quite constant over the years.

The accuracy of the multi-beam data has increased over the years mainly due to improvements in positioning and measuring equipment. The average survey density increased from 5-8 measurements per grid cell for the first survey years up to about 50 values from 2001 onwards. Incidentally some grid cell values are missing. For many areas no multi-beam data are available for the years 1998 and 1999.

## 4. Calibration

The applied calibration method is described in the appendix. In summary this method approximates the posterior distribution of the parameter vector conditionally on the observations, by comparing the posterior probabilities of a set of models, run over a range of parameter settings. Since this is actually a process of model selection, it allows the joint calibration of model architecture and model parameters.

### 4.1 Calibration strategy

The calibration is carried out using data of three sand waves. These sand waves are chosen to lie far apart in order to gain some insight in spatial calibration sensitivity. They are located in the central part of the approach area, the western part of the Euro channel and the southern part of the Twin area, see figure 1. For reasons of computational feasibility, of each sand wave a limited area of about 25,000 $m^2$ is selected. Each area is aligned along the sand wave crest. The lengths of the selected crests vary from 300 to 400 m.

The calibration method as described in the appendix focuses on comparing one year lead time predictions. Testing the calibration results clearly indicated the sub-optimality of predictions for longer lead times. Extending the calibration method so that it simultaneously compares the predictions for all lead times up to five years, clearly improved results.

Another adjustment made to the calibration method concerns focussing exclusively on the (nautically critical) positive part of the prediction error distribution instead of on the full distribution. The reason for this is that it was conjectured that the asymmetry of the negative tail induces multi-modality of the posterior distribution causing the calibration results to be erratic. This adjustment improved the robustness of the calibration method and the quality of the calibration results (from the nautical point of view).

### 4.2 Initial state

The initial state distribution is partly calibrated using the above-mentioned calibration method and partly derived directly from the measurement data of the first time step. Table 1 summarises the choices made.

This initial state specification strategy is held fixed in the continuation of the calibration process.

Table 1: Initial state distributions

| state element | value | unit |
|---|---|---|
| initial local levels | average of first survey | dm |
| variance of initial levels | variance of first survey | $dm^2$ |
| initial local trend values | 0 | dm/yr |
| variance of initial trends | 10 | $(dm/yr)^2$ |
| initial average area depth | average of first survey | dm |
| variance of initial area depth | 50 | $dm^2$ |

### 4.3 Model architecture

Since the model architecture is unknown a priori, its specification is considered part of the calibration process. In order to limit the large number of degrees of freedom of the model principle, all kernels are constrained to be spherical Gaussians with identical variances and the model support points are located on a north-south oriented square grid (the model domain may however be bounded by any quadrangle).

Joint calibration of the remaining architectural variables and the main model parameters, resulted in a value of 17.5 m for model grid size as well as for kernel width (standard deviation).

These architectural choices are used for further calibration of the model parameters.

### 4.4 Model parameters

The measurement noise and the local level and trend discount factors (see section 2.2) are the decisive calibration parameters. Calibration results indicate that plausible measurement noise values range from 25 to 50 $dm^2$ and that both discount factors range from 0.95 to 0.99.

Table 2 shows the final parameter values which are assessed after visual inspection of normal probability plots of the normalised prediction errors for all calibration sand waves and all prediction lead times (1 to 5 years, see figure 2). This visual inspection revealed a lead time dependency of the error distributions with a tendency to underestimate the prediction uncertainty for the longer lead times. To compensate for this, a slightly conservative final parameter choice was made at the cost of introducing a slight overestimation of the prediction uncertainty for the shorter lead times.
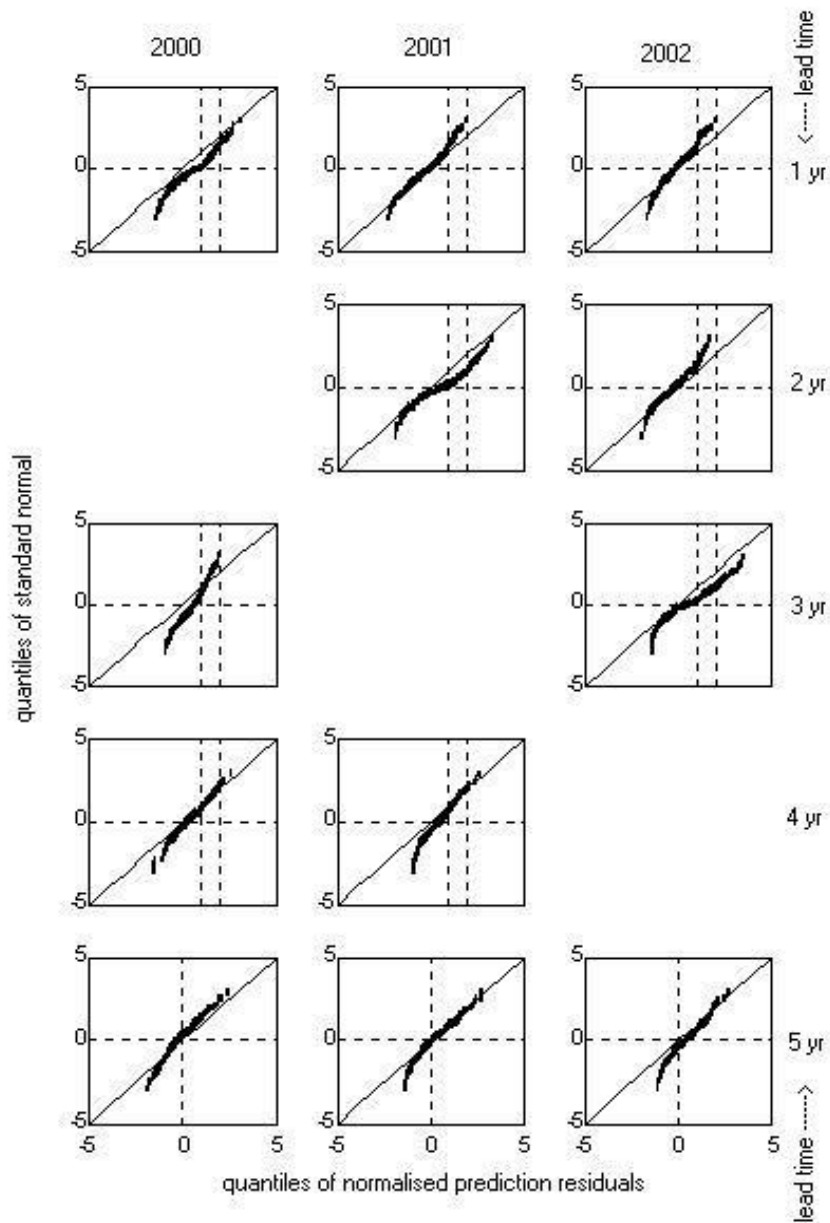
Table 2: Model parameters

| parameter | value | unit |
|---|---|---|
| local level discount factor | 0.93 | dimensionless |
| local trend discount factor | 0.93 | dimensionless |
| system noise area depth | 0.5 | $dm^2$ |
| measurement noise | 23 | $dm^2$ |

## 5 Model results

The model performance is tested for two sand waves other than the ones used for calibration, located in the southern part of the Twin area and in the far western part of the Euro channel.

Figure 2 depicts typical normal probability plots of the prediction residuals of the Euro channel sand wave for three prediction years and all available lead times (of 1998 no data are available). The residuals are normalized with their corresponding prediction standard deviations. The nautically interesting errors are located in the upper right-hand corner of the plots, deviations below the diagonal being nautically unfavourable. The lower left-hand corners of the plots show the asymmetry of the error distribution as mentioned in section 4.1. The plots also reveal the tendency to overestimate the prediction uncertainty for the shorter lead times due to the conservative parameter choice as discussed in section 4.4.

**Figure 2: Quantile-quantile plots of predictions of a sand wave crest area in the Euro Channel**

The sensitivity of the model to measurement outliers was assessed by assigning one randomly located outlier on every north-south grid line of the measurement grid. Outliers sampled from a Gaussian with a standard deviation of up to 20 dm exhibit an almost negligible effect on the model output.

## 6. Discussion

The results of the adopted model principle demonstrate its utility from a probabilistic-predictive point of view. Despite being somewhat conservative for the shorter lead times, the prediction uncertainties appear to be of an economically interesting magnitude. In 2004 the utility of the model will be further tested by applying it to a few dozen sand waves.

From a morphological point of view the model principle is very limited. Incorporating physical principles is expected to lead to improvements. The flexibility of the state space modelling framework offers many opportunities for exploring potential model extensions. Potential directions of improvement are the incorporation of transport model principles obeying the continuity equation and imputing model output information from morphological models like e.g. a model for predicting the sand wave recovery rate after dredging [Knaapen and Hulscher, 2002]. The latter information source will significantly improve the adaptation capabilities of the statistical model to post-dredging conditions and will also allow the inclusion of information on the sand wave growth rate and its saturation behaviour.

The applied predictive calibration method has proved its effectiveness. The successful adjustments forced by the asymmetry of the error distribution and the need to include multiple lead times, demonstrate its flexibility. Future improvements may include the use of automatic calibration methods allowing for easier adaptation to local dynamics and local topography.

The fact that the model output is rather insensitive to measurement outliers emphasises the robustness of the modelling principle.

A significant part of the measurement noise consists of morphological phenomena that are considered as unpredictable for the lead times of interest (e.g. migrating mega-ripples). The use of filtering techniques for extracting this noise component from the data before running the model, may lead to improved model performance. Because it actually concerns morphological noise instead of measurement noise, the variance of this extracted noise component must finally be added to the prediction variance explicitly. Knowledge of the 'real' measurement noise may be economically beneficial in reducing the operational prediction uncertainty by a justified subtraction of this noise component from the prediction variance.

The Gaussian assumption of the prediction error density is clearly violated. Gaining insight into the causes of the perceived asymmetry of the prediction error distribution may help improve the calibration process. Potential causes are the presence of mega-ripples and intrinsic aspects of the hydrographical data collection process.

## Acknowledgments

## References

B.D.O. Anderson and J.B. Moore, Optimal filtering, Prentice Hall, New Jersey 1979.
J.A. Hoeting, D. Madigan, A.E. Raftery and C.T. Volinsky, Bayesian model averaging: A tutorial, *Statistical Science*, Vol. 14 (4), pp. 382-417, 1999.
M.A.F. Knaapen and S.J.M.H. Hulscher, Regeneration of sand waves after dredging, *Coastal Engineering*, 46, pp. 277-289, 2002.
H.K.H. Lee, C.H. Holloman, C.A. Calder and D.M. Higdon, Flexible Gaussian processes via convolution, *Technical report* 02-09, Duke University, 2002.
J. Stroud, P. Müller and B. Sansó, Dynamic models for spatio-temporal data, *Journal of the Royal Statistical Society*, Series B (63), pp 673-689, 2001.
M. West and P. Harrison, Bayesian forecasting and dynamic models, Springer, New York, 2nd ed., 1997.

# Appendix

## Calibration method

For calibrating the space-time model at hand, a full Bayesian approach is pursued [West and Harrison, 1997, Hoeting et al., 1999]. This approach consists of estimating the probability distribution of the parameter vector q, conditionally on the observations y, using Bayes' rule: $P(q|y)=p(y|q)p(q)/p(y)$.

The posterior distribution $P(q|y)$ is proportional to the product of the parameter prior $p(q)$ and the likelihood $p(y|q)$. $p(y)$ is the normalisation factor of the resulting probability distribution.

Since a closed form solution is infeasible the problem is discretised. Within the parameter space a discretely distributed set $\alpha$ is chosen. The probability masses of each parameter combination $\alpha_i$ (i=1:N, with N the number of combinations) of this multi-dimensional set are recursively estimated using Bayes' rule: $p(\alpha_i|D_t) = p(y_t|\alpha_i,D_{t-1}) \, p(\alpha_i |D_{t-1})/p(y_t|D_{t-1})$.

$D_t$ represents all data up to and including time t. $p(y_t|D_{t-1})$ is the normalising constant obtained through summation over all $\alpha_i$ terms in the numerator. The elements of the initial prior $p(\alpha_i |D_0)$ are assigned equal (uniformly distributed) probability masses.

For selecting the optimal model the maximum a posteriori approach is followed. In this discretised approach this implies choosing the parameter combination with highest posterior probability mass.