

# **On identification of nonlinear systems**

CENTRALE LANDBOUWCATALOGUS



0000 0576 1610

Promotoren: dr. ir. J. Grasman  
hoogleraar in de wiskunde, inclusief de numerieke wiskunde  
dr. D. Rasch  
hoogleraar in de wiskundige statistiek  
Co-promotor: dr. M. de Gee  
universitair hoofddocent bij de vakgroep wiskunde

Sixtus Leonardus Jacobus Mous

## On identification of nonlinear systems

### Proefschrift

ter verkrijging van de graad van doctor  
in de landbouw- en milieuwetenschappen  
op gezag van de rector magnificus,  
dr. C.M. Karssen,  
in het openbaar te verdedigen  
op maandag 11 april 1994  
des namiddags te vier uur in de Aula  
van de Lanbouwniversiteit te Wageningen.

**BIBLIOTHEEK  
LANDBOUWUNIVERSITEIT  
WAGENINGEN**

CIP-DATA Koninklijke Bibliotheek, Den Haag

Mous, Sixtus Leonardus Jacobus

On identification of nonlinear systems /  
Sixtus Leonardus Jacobus Mous. - [S.l.:s.n.]

Thesis Wageningen. - With ref. -

With summary in Dutch.

ISBN 90 - 5485 - 225 - 9

Subject headings: system identification / chaotic behavior in systems.

## Stellingen

- 1 Bij het ontwerpen van een experiment ten behoeve van inverse modellering dient over de identificeerbaarheid van de model-structuur grondig vóórgedacht te worden.
- 2 Met behulp van Kalman filtering kan de bijdrage van de verschillende foutenbronnen aan de voorspellingsfout beter worden geschat. Hierdoor kan voorkomen worden dat bij de verbetering van weervoorspellingen onnodig veel geld wordt geïnvesteerd om een kleine bijdrage nog verder te verkleinen.
- 3 De bepaling van de massa-balans, ter controle van de nauwkeurigheid van numerieke oplossingschema's, zoals voorgesteld door Celia and Bouloutas, is slechts van beperkte waarde omdat de berekening van deze maat ook aan onnauwkeurigheden onderhevig is.  
  
Celia, M.A. and E.T. Bouloutas, A general mass-conservative numerical solution for the unsaturated flow equation, Water Resour. Res., 1990.
- 4 Omdat numerieke wiskunde een steeds belangrijkere rol speelt bij het modelleren van veel bodemfysische processen, verdient het aanbeveling om mathematische en statistische software libraries een inherent onderdeel te laten uitmaken van het standaard gereedschap van de modelbouwer.
- 5 Het optimaliseren van een experiment vraagt een goede communicatie tussen experimentele onderzoekers, modelbouwers en statistici, omdat het theoretisch optimale ontwerp van het experiment praktisch niet uitvoerbaar kan zijn.
- 6 De positieve effecten van zowel resonantie-therapie en transcendente meditatie zijn gebaseerd op "wishfull thinking" en zijn wetenschappelijk niet te onderbouwen.
- 7 Om tegen wateroverlast verzekerd te zijn moet men een drijvend huis bouwen.
- 8 Het invoeren van een ISO-9000 certificaat ter beoordeling van de kwaliteit van een opleiding kan de huidige prioriteitsstelling bij veel onderwijsinstellingen, te weten kwantiteit boven kwaliteit, weer omdraaien.
- 9 De huidige spelregel bij zaalvoetbal, waarbij een terugspeelbal op de keeper leidt tot een directe vrije trap vanaf de zes-meterlijn, leidt tot prijsschieten en draagt daarom niet bij aan de primaire doelstelling van sport, namelijk verbroedering.
- 10 Het zou de samenleving ten goede komen als notabelen hun plaats weer zouden innemen aan de stamtafel van de lokale kroeg.

voor mirriam en  
mijn ouders

voor mirriam en  
mijn ouders

# CONTENTS

## Chapter 1

1.	Introduction	1
1.1	Background	1
1.2	System identification	3
1.3	Objectives of this study	6
1.4	Outline of the thesis	8
	References	10

## Chapter 2

2.	Numerical solutions for the one-step experiment	13
2.1	Introduction	14
2.2	Problem formulation	16
2.3	Numerical approximations	21
2.4	Conclusions	34
	References	35

## Chapter 3

3.	Identification of the movement of water in unsaturated soils; the problem of identifiability of the model	37
3.1	Introduction	38
3.2	One-dimensional transient water flow	40
3.3	Identifiability	42
3.4	Experiment design	47
3.5	Results	49
3.5	Summary and discussion	56
	References	57

Chapter 4	
4. Two methods for assessing the size of external perturbations in chaotic processes	61
4.1 Introduction	62
4.2 Sentinels for detecting perturbations	64
4.2.1 The Rössler attractor	69
4.3 Sentinels for estimating uncertain parameters	72
4.3.1 A low-order spectral model of the atmospheric circulation with a perturbed equator-pole temperature gradient	73
4.4 Extended Kalman filtering for estimating uncertain parameters	78
4.5 Conclusions	84
References	86
Appendix A	87
Chapter 5	
5. Estimating uncertain model and noise parameters in chaotic models with applications in meteorology	89
5.1 Introduction	90
5.2 Extended Kalman filtering	94
5.3 Numerical experiments with the Lorenz attractor	100
5.4 Hierarchic optimization	105
5.5 Numerical experiments with the barotropic vorticity equation	108
5.6 Conclusions	114
References	116
Summary	119
Samenvatting	123
Dankwoord	127
Curriculum Vitae	129



# Chapter 1

## INTRODUCTION

### 1.1 BACKGROUND

Models of physical, chemical or biological processes are probably more important than non-specialists may realize. One can use models for various purposes, e.g., simulation, prediction, control. With simulation models one can analyze the effect of different inputs, moreover, different scenarios of control policies for a process can be evaluated. Prediction models on the other hand may provide very useful information for making decisions. For example, the weather forecast can be very important for the decisions a farmer has to make. In control problems a model of a process is essential for the design of a control system: to maintain a constant temperature in a stirred tank, it is necessary to know the transfer function between feed flow rate and stirred tank temperature in order for the control system to compensate for disturbances in the feed flow rate or feed temperature.

System identification deals with the problem of building accurate models of processes. Since many processes are nonlinear, complex numerical models may be necessary to describe them. An introduction into nonlinear regression techniques for estimating parameters in nonlinear dynamical models can be found in e.g., Bates and Watts (1988). Nonlinear regression methods are suitable for the identification of simulation models, since these models aim to predict the output of the process based on input data. For the identification of prediction models, nonlinear filtering techniques may be more appropriate. Nonlinear filtering techniques are used to estimate the future state of the process based on input and output data of the past. A criterion based on the difference between the predicted state and the observation of this state is then a natural choice. An introduction into nonlinear filtering theory can be found in e.g., Jazwinski (1970) and Anderson and Moore (1979).

In this thesis we will focus on models that are used to describe nonlinear processes in hydrology and meteorology. The first process that will be analyzed deals with the movement of water in unsaturated soils. The knowledge of the displacement of water in unsaturated soils is important for the development of water management systems: water from precipitation, from irrigation or from an influent river, infiltrates through the ground surface and percolates downward through the unsaturated zone into a phreatic aquifer. Information about the displacement of water in the unsaturated zone is needed in order to determine the replenishment of a phreatic aquifer as part of the groundwater system. A description of the flow of water in the unsaturated zone is also needed to predict the spread and accumulation of dissolved pollutants in the unsaturated zone and the rate and concentration at which these pollutants reach the water table (Bear and Verruijt, 1987).

In meteorology models are used in the prediction of the weather for a number of days ahead. Models of the atmospheric circulation are analyzed in the second part of this thesis. Due to sensitive dependence on the initial state

it is not possible to produce accurate forecasts beyond a range of about five days. Although the accuracy of weather models has been improved substantially in the last decades, it is very difficult to compare the quality of different weather models, because the predictability range may depend on the initial state. In the evaluation of the model it may be meaningful to distinguish here between the mathematical model and the numerical integration scheme. Numerical errors add up to errors due to neglecting certain physical processes in the mathematical model. It is known how numerical errors should be estimated. However, there does not exist yet a systematic approach to evaluate the error that is caused by not incorporating certain physical processes or by incorrect parameterization of such processes. In the second part of this thesis we will address to this problem.

## 1.2 SYSTEM IDENTIFICATION

In some literature system identification is sometimes denoted by 'inverse modeling' or 'inverse problem' (e.g., Kool and Parker, 1988). It refers to the determination of (differential) equations that describe the physical process, given the input and output signals. Zadeh (1962) gives a more precise definition of system identification. It is the determination given input and output data of a system within a specified class of systems to which the system under test is equivalent. In this formulation 'the system under test' is the process and the elements of 'the class of systems' are the models. Due to the systems complexity, as well as the incomplete availability of observations and the limited a priori knowledge, it is generally impossible to try to obtain an exact mathematical description of the physical process. Therefore, mathematical models only can describe the system approximately. In this light it is more natural to consider system identification as approximate modeling on the basis of observed data and a priori knowledge

(Janssen, 1988).

In general, one can say that the identification process contains three essential ingredients:

- selection of the class of systems
- definition of a criterion of best fit
- experimental design

The problem of selecting the class of systems is highly influenced by the a priori knowledge of the process. If one has sufficient insight into the process, one may give a detailed "physical" model structure. However, often one has to use some type of empirical relations or even a "black box" model structure. A model structure is a description of a process, expressed in terms of mathematical equations. Sometimes one has to decide between two or more model structures. The choice of a specific model structure may be based on statistical selection criteria. In Rasch et al. (1992) some of these statistical selection criteria are discussed. Besides these statistical selection criteria other criteria can play a role also. For example, in predicting the future behavior of a process, it may be more suitable to use a simple model than a detailed description of the process, since the calculation of the prediction with the detailed model may cost so much time that the prediction becomes worthless.

After the choice of the model structure, the identification problem is reduced to a parameter estimation problem. Parameter estimation may be defined as the experimental determination of values of parameters that govern the dynamic behavior assuming that the structure of the model is known. Eykhoff (1974) pointed out that the distinction between knowledge of the model structure and parameters is not as straightforward as it may appear on first sight. The change from a non-zero to zero value of a parameter may represent a simplification in the structure, as that 'branch' of the model may be deleted. Thus, the result of the parameter estimation problem can be that

one has to reconsider the structure of the model.

To determine the best approximating model, we have to define a criterion of best fit. For the identification of simulation models a reasonable choice would be to use the sum of squares of the differences between observations and model values as object function. If the purpose is not simulation but prediction it may be more useful to define an object function based on prediction errors. In the literature of system analysis the first method is called an output-error method whereas the second method is called a (one-step ahead) prediction error method (e.g., Ljung, 1987).

Some statistical methods, such as the Maximum Likelihood method, require a priori knowledge of the probability distribution of the error terms. Often it is assumed that the observations are contaminated by white Gaussian noise and that the differential equation is perturbed by a white Gaussian noise process. For linear systems the Maximum Likelihood estimates can then easily be obtained using the Kalman-Bucy filter (cf., Harvey, 1981). For nonlinear systems we may linearize near some solution and make an approximation in this way.

An important aspect of identification is experimental design. Experimental design means the specification of the input signals, the sampling rate and the number of observations to be taken (Rasch, 1990). In this thesis only the specification of the input signals is considered. The advantage in contriving a well-designed experiment is that we may obtain richer and more informative output signals. Of course we cannot manipulate the input signals freely, because the experimental conditions may not deviate too much from the conditions in the final application. Sometimes there are no controllable input signals, for example with the weather system. For the identification of such systems the determination can only be based on (passive) observations of the process.

### 1.3 OBJECTIVES OF THIS STUDY

Not every meaningful looking combination of model structure, criterion of best fit and experimental design will lead to a unique estimate of the unknown parameters. The problem of not finding a unique solution of the identification problem was already noticed by Bellman and Åström (1970). To analyze this non-uniqueness problem, they have introduced the concept of structural identifiability. A model is called (globally) identifiable if the model structure with different parameters values and identical input signals yields different output signals. However, an identifiable model structure is not sufficient to ensure a unique parameter estimate because uniqueness may also depend on the chosen criterion, non-linearities, the dynamic behavior of the system, and the design of the experiment (e.g., Walter, 1982, Ljung, 1987). One can distinguish several cases in which problems occur in the unique estimation of the parameters. In this study, we first have considered the case where the object function is insensitive to certain parameters or linear combinations of parameters. It is clear that then the optimization problem will not have a unique solution, because many parameter combinations will give an (almost) equal value of the object function. As a second case, we have considered a combination of criterion and model structure which makes the problem ill-posed. This type of optimization problems cannot be solved by ordinary optimization algorithms. An example of such a problem is the identification of a chaotic system (Baake et al., 1992).

Possible causes of an insensitive object function are overparameterization of the model structure, non-informative input signals and non-linearity (Ljung, 1987). Noise in the system may mask the insensitivity to parameters because the optimization problem may have many local minima. One can easily be misled when an optimization algorithm converges to such a local minimum. In many practical situations it can be very difficult to prove that the model structure is not identifiable. In practice we have to restrict our-

selves to a specific combination of model structure, criterion and experimental design and analyze the existence of a unique solutions for this combination.

An example of such a combination is studied in the first part of this thesis. There we analyze the identifiability of the model used in the ONE-STEP method, describing the movement of water in unsaturated soils (Kool et al., 1985, 1988). The water movement in the unsaturated zone is described by the Richards equation in conjunction with the Mualem-van Genuchten relations (Mualem, 1976, van Genuchten, 1980). The study of this problem was motivated from the statement made by several authors that not all parameters could be obtained uniquely (e.g., Hornung, 1983, Kool et al., 1985 and van Dam et al., 1990). The object was to find the cause of this non-uniqueness problem and to suggest designs for new experiments. In a recent paper, Toorman et al. (1992) describe a new experimental setup, in which pressure head measurements have been included. The improved sensitivity of the object function with respect to the parameters, which was reported in their paper, could also be predicted as a result of this study.

In the second part of this thesis the identification of chaotic systems is studied. As mentioned above the identification of these systems may lead to an ill-posed optimization problem (Breedon and Hübner, 1990, Farmer and Sidorowich, 1991, Baake et al., 1992). The optimization method with an output-error criterion for this class of problems is ill-posed, because the model's solution depends sensitively on its initial states. The observed values and the model values will then diverge due to the limited accuracy of the initial state. Consequently, the optimization problem with this criterion will also have many local optima. Therefore, one may be tempted to say that the model is also not-identifiable. In the previous case, an other experimental setup that gives more informative output signals may solve the problem of non-identifiability. Here, this approach is not necessarily to solve the non-identifiability problem since in essence the output signals contains enough

information (Baake et al., 1992). An approach to solve the problem of non-identifiability may be to use a different criterion of best fit. We have analyzed the performance of criteria that are less sensitive to disturbances in the initial state.

In meteorology this is a highly actual problem because the atmospheric circulation behaves chaotic. Parameters in models for the atmospheric circulation are often known with a limited accuracy. The effect of a small perturbation of a parameter may be significant (de Swart, 1988, Grasman and Houtekamer, 1992). An example of a small perturbation is a possible deviation in the equator-pole temperature gradient. The equator-pole temperature gradient can be seen as the driving force of our atmosphere. Therefore, detecting a systematic change in this driving force may give more insight into the greenhouse effect.

## 1.4 OUTLINE OF THE THESIS

This thesis contains four self-contained chapters. In chapter 2 and 3 the main topic is the identification of displacement of water in the unsaturated zone. The accuracy of some numerical schemes is analyzed in chapter 2. Special attention is given to the size of truncation errors due to spatial discretization. These errors may be large because of the presence of a steep "drying front". This makes it difficult to approximate the fluxes accurately. One may obtain accurate approximations of the fluxes by using a variable step-size scheme. This scheme is made more efficient by choosing optimal locations of the nodes.

In chapter 3 the connection between the numerical accuracy of the model output and the identifiability of the model is studied. Identifiability analysis shows that not all parameters of the Mualem-van Genuchten relations can be estimated using the ONE-STEP method. Furthermore, the question



has been raised whether other experimental designs may result in richer and more informative output signals. This chapter has been presented as a paper at the XVI General Assembly of the European Geophysical Society and has been published in the *Journal of Hydrology*.

In chapter 4 and 5 applications in meteorology are treated. In chapter 4 two methods for assessing the size of an external perturbation in a chaotic model are investigated. The first method is based on Lions' sentinel functions (Lions, 1988, 1990). This method is originally developed for analyzing processes that are described with partial differential equations. In this chapter a modified sentinel method is presented that can be used to detect perturbations in processes that are described by systems of ordinary differential equations. The second method is an adaptive extended Kalman filter. In this method the unknown parameters are regarded as random variables and the state vector is augmented with these variables. Since the extended Kalman filter yields estimates of the state vector, it provides an estimate of the unknown parameters as well. This chapter has been published in *Mathematical Models and Methods in Applied Sciences*, with J. Grasman as co-author.

Since the use of the adaptive extended Kalman filter was rather successful for these meteorological problems, this method has been developed further in chapter 5. The advantage of the adaptive extended Kalman filter over the sentinel method is that it provides more accurate estimates in a shorter time. An additional advantage is that it is an on-line method. This gives the possibility to estimate also a parameter that slowly changes with time (Mous, 1993). There are also some disadvantages: reduction in performance because the filter uses initially the wrong parameter values, filter divergence due to the chaotic behavior of the process and finally the computational costs in case of high dimensional systems.

The first problem may be overcome by repeating the process using the estimated parameter values from the first run as initial estimates for the parameters in the second run. Although this approach may have a better

performance, the advantage of an on-line method disappears. The second problem, the divergence of the extended Kalman filter, can be solved by adding an artificial noise term to the state equations. With this noise term we can control the memory of the extended Kalman filter (Jazwinski, 1970). However, for an optimally working extended Kalman filter the noise parameters that describe the artificial noise term have to be estimated as well. In chapter 5 an approximated maximum likelihood method is used to estimate the unknown model parameters and the noise parameters. The last problem of high computational costs has not been solved yet. A solution may be found in using fast sub-optimal filters, such as the simplified Kalman filter of Dee (1991) or in using other criteria, which can make the function and gradient evaluation much faster.

## REFERENCES

- Anderson, B.D.O. and J.B. Moore, *Optimal Filtering*, Prentice-Hall, Englewood Cliffs, 1979.
- Baake, E., M. Baake, H.G. Bock and K.M. Briggs, Fitting ordinary differential equations to chaotic data, *Phys. Rev. A*, vol. 45, 8, 5524-5529, 1992.
- Bates, D.M. and D.G. Watts, *Non-linear regression analysis and its applications*, Wiley, New York, 1988.
- Bear, J. and A. Verruijt, *Modeling Groundwater Flow and Pollution*, D. Reidel Publishing Company, Dordrecht, 1987
- Bellman, R. and K.J. Åström, On structural identifiability, *Math. Biosci.*, 7, 329-339, 1970.
- Breeden, J.L. and A. Hübler, Reconstructing equations of motion from experimental data with unobserved variables, *Phys. Rev. A*, 42, 10, 5817-5826, 1990.

- Dee, D.P., Simplification of the Kalman filter for meteorological data assimilation, *Quart. J. Roy. Meteor. Soc.*, 9-16, 1991.
- De Swart, H.E., Low-order spectral models of the atmospheric circulation: a survey. *Acta Appl. Math.*, 11, 49-96, 1988.
- Eykhoff, P., *System identification, parameters and state estimation*, Wiley, London, 1974.
- Farmer, J.D. and J.J. Sidorowich, Optimal shadowing and noise reduction, *Physica D*, 47, 373-392, 1992.
- Grasman, J., and P. Houtekamer, Methods for improving the prediction of dynamical processes with special reference to the atmospheric circulation, *Proc. of IUGG Symp. "Nonlinear Dynamics and Predictability of critical Geophysical Phenomena*, W.I. Newman and A.M. Gabrielov (eds.), Am. Geophysical Union, 1992.
- Harvey, A.C., *Time Series Models*, Wiley, New York, 1981.
- Hornung, U., Identification of nonlinear soil physical parameters from input-output experiments, in: P. Deuflhard and E. Hairer (ed.), *Workshop on numerical treatments of inverse problems in differential and integral equations*, Birkhäuser, Bosten, 227-237, 1983.
- Janssen, P., *On model parametrization and model structure selection for identification of MIMO-systems*, Ph.D.Thesis, Technische Universiteit Eindhoven, 1988.
- Jazwinski, A.H., *Stochastic Processes and Filtering Theory*, Academic Press, Paris, 1970.
- Kool, J.B., J.C. Parker and M.T. van Genuchten, Determining soil hydraulic properties from one-step outflow measurements by parameter estimation, I, theory and numerical studies, *Soil Sci. Soc. Am. J.*, 49, 1348-1354, 1985.
- Kool, J.B. and J.C. Parker, Analysis of the inverse problem for transient unsaturated flow, *Water Resour. Res.*, 24, 513-522, 1988.
- Lions, J.L., *Sur les sentinelles de systèmes distribuées*, C.R.A.S., Paris, 1988.

- Lions, J.L., Sentinels and stealthy perturbations, *Proc. Int. Symp. on Assimilation of Observations in Meteorology and Oceanography*, Clermont-Ferrand, World Meteorological Organization, 13-18, 1990.
- Ljung, L., *System identification: theory for the user*, Prentice-Hall, Englewood Cliffs, 1987.
- Mous, S.L.J., Detection of a perturbed equator-pole temperature gradient in a spectral model of the atmospheric circulation, to appear in the *Proc. of the 75th anniv. conf. of the Wageningen Agricultural University "Predictability and nonlinear modeling in natural science and economy"*, 1993.
- Mualem, Y., A new model for predicting the hydraulic conductivity of unsaturated porous media, *Water Resour. Res.*, 12, 513-522, 1976.
- Rasch, D., Optimum experiment design in nonlinear regression, *Comm. Statist. Theory Meth.*, 19(12), 4789-4806, 1990.
- Rasch, D., V. Guiard and G. Nürnberg, *Statistische Versuchsplanung - Einführung in die Methoden und Anwendung des Dialogsystems CADEMO*, Gustav Fischer Verlag, Stuttgart-Jena-New York, 1992.
- Toorman, A.F., P.J. Wierenga and R.G. Hills, Parameter estimation of hydraulic properties from one-step outflow data, *Water Resour. Res.*, 28, 3021-3028, 1992.
- Van Dam, J.C., J.N.M. Stricker and P. Droogers, From One-Step to Multi-Step, determination of soil hydraulic functions by outflow experiments, Rep. 7, Dep. of Hydrology, Soil Physics and Hydraulics, Agric. Univ. Wageningen, 1990.
- Van Genuchten, M.T., A closed-form equation for predicting the hydraulic conductivity of unsaturated soils. *Soil Sci. Soc. Am. J.*, 44, 892-898, 1980.
- Walter E., *Identifiability of state space models*, Springer-Verlag, New York, 1982.
- Zadeh, L.A., From circuit theory to system theory, *Proc. IRE*, 50, 856-865, 1962.

# Chapter 2

## NUMERICAL SOLUTIONS FOR THE ONE-STEP EXPERIMENT

### Abstract

The ONE-STEP experiment (Kool et al., 1985) describes an experimental procedure for determining the hydraulic properties of a soil sample. It involves measurement of cumulative outflow with time from a soil core placed in a pressure cell. The water flow in this soil core is one-dimensional and can be described by the Richards equation.

In this paper some numerical solutions of this equation are compared. The effort that is required to obtain accurate solutions, for both the finite element and the finite difference approach, strongly depends on the experimental conditions. In case the pneumatic head applied to the pressure cell is low, e.g., 250 cm, both approaches are sufficiently accurate. However, if the pneumatic head is high, say 1000 cm, a small spatial discretization step size

is required to obtain accurate solutions. Since such a small step size is only needed in a small part of the solution space, a variable step size scheme will improve the efficiency. The optimal positions of the nodes are determined using estimates of the truncation errors. It turns out that this variable step-size scheme yields the same order of accuracy using only 1/4 of the number of nodes.

## 2.1 INTRODUCTION

Recently, a considerable effort has been put in the determination of soil hydraulic properties from transient flow data. Knowledge of these hydraulic properties is important for the calculation of groundwater movement in unsaturated soils. In most models it is assumed that the movement of water in unsaturated soils satisfies the classical Richards equation. The main subject of this paper is the analysis of the discretization errors that are made by the numerical approximation of the solution of this equation. This topic plays an important role in identification methods such as the ONE-STEP method (e.g., Kool et al., 1985, 1988). It is therefore that in this paper we study the problem of approximating numerically the solution of the Richards equation as a separate problem.

Several authors have given numerical approximation schemes for the Richards equation. Hanks et al. (1969) presented a finite difference approach to solve the Richards equation, while van Genuchten (1982) presented a finite element approach. The integration method originally used in the ONE-STEP method was adopted from van Genuchten. More recently, several papers have appeared on so-called mass-conservative schemes (Milly, 1985, Ceila et al., 1990). These schemes have smaller errors due to time discretization compared with schemes where the implicit Euler method is used for time stepping. Another approach to make the time discretization errors small is to use a

high-order integration scheme, for example a linear multistep method or a Runge-Kutta method. For these high-order schemes a larger time step can be taken than the one for the Euler scheme without losing accuracy in the approximations.

As mentioned above, the discretization errors are not only caused by time discretization, spatial discretization may also play a role. The main question is: what is the most important cause of the errors, the time discretization or the spatial discretization? We know from Milly (1985) and Ceila et al. (1990) that the errors due to time discretization can be significant and that to this point special attention has to be given. Ceila et al. (1990) have shown that their mass-conservative scheme, based on the mixed-formulation of the Richards equation, yields smaller errors due to time discretization. However, they also mention that a numerical scheme that conserves mass is not sufficient to guarantee accurate solutions of the mass-balance differential equation. In this paper the attention is focused on numerical errors due to spatial discretization. The importance of the spatial discretization can be evaluated by estimating the magnitude of the spatial discretization errors; this can be done by making the time discretization errors comparatively small, so that the errors only depend on the spatial discretization.

The paper begins with an outline of the outflow experiment. Special attention is given to the movement of water in the porous plate, since the movement of water in this plate is described by a degenerated Richards equation. The movement of water in the soil has to be solved numerically whereas the movement of water in the porous plate can be simply solved analytically. Numerical solutions using either a finite difference, a finite element or a finite element with variable step-size method are compared for two experimental setups of the outflow experiment. These experimental setups differ in the pneumatic head that is applied to a soil core. With a high pneumatic head the solution is characterized by a "steep drying front", which makes the solution sensitive to numerical errors. Finally in the last section we

make some concluding remarks.

## 2.2 PROBLEM FORMULATION

The ONE-STEP method (Kool et al., 1985) is developed for identification of one-dimensional transient water-flow in porous media in a simple and fast way. In this section an improved mathematical description of this experiment is presented.

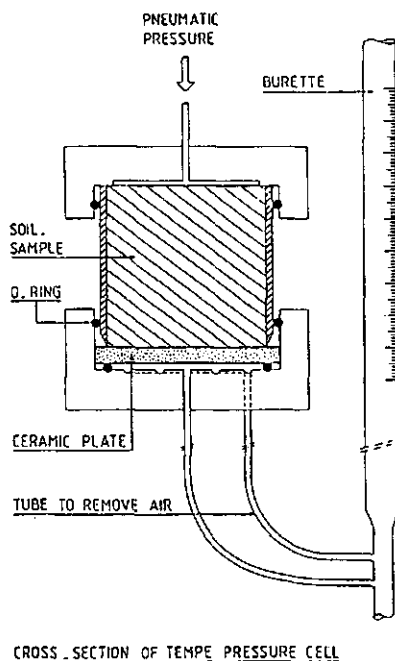


Figure 2.1. The measurement system.

Figure 2.1. shows schematically the measurement system. It consists of a soil sample and a porous plate placed in a pressure cell. Initially the soil



sample is almost saturated. As a result of the pneumatic head that is applied to the cell water seeps out. Comparing the measured cumulative outflow  $Q(t)$  with the outflow that is calculated by a model, one may obtain the unknown parameters in the model using an inverse modeling procedure. The mathematical model that is used to calculate the outflow is based on the Richards equation:

$$C(h) \frac{\partial h}{\partial t} = \frac{\partial}{\partial x} \left( K(h) \left( \frac{\partial h}{\partial x} - 1 \right) \right), \quad 0 \leq x \leq l, \quad t \geq t_0, \quad (2.1)$$

$$h(x, t_0) = h_0(x), \quad (2.2)$$

$$\frac{\partial h(0, t)}{\partial x} = 1, \quad (2.3)$$

$$h(l, t) = -h_c = -\frac{\Delta p}{\rho g}, \quad (2.4)$$

where

- $x$  = the vertical coordinate, with  $x = 0$  at the top of the core and  $x = l$  at the bottom of the porous plate,
- $h(x, t)$  = pressure head at time  $t$  at point  $x$  in the medium,
- $C(h)$  = the differential water capacity of the soil at a pressure head  $h$ ,
- $K(h)$  = hydraulic conductivity of the soil at a pressure head  $h$ ,
- $\Delta p$  = gauge gas pressure applied to the cell,
- $\rho$  = density of water,
- $g$  = gravitation acceleration.

In the above formulation it is assumed that the initial conditions (2.2) are

known. Furthermore it is assumed that the pressure at the bottom of the plate is atmospheric and that no water infiltrates into the soil sample.

The mathematical functions for the hydraulic properties of the soil,  $C(h)$  and  $K(h)$  are modeled by Mualem (1976), van Genuchten (1980) and Wösten and van Genuchten (1988):

$$C(h) = \frac{d\theta}{dh} = \alpha m (\theta_s - \theta_r) S_e^{1/m} (1 - S_e^{1/m})^m (1 - m)^{-1}, \quad (2.5)$$

$$K(h) = K_s S_e^\gamma (1 - (1 - S_e^{1/m})^m)^2, \quad (2.6)$$

$$\theta = \theta_r + S_e (\theta_s - \theta_r) \quad (2.7)$$

with

$$S_e = (1 + |\alpha h|^n)^{-m},$$

$$m = 1 - 1/n$$

and where  $\theta$  is the volumetric water content. The unknown model parameters are  $\alpha$ ,  $n$ ,  $\theta_r$ ,  $\theta_s$ ,  $K_s$ ,  $\gamma$ . In the numerical examples presented later we use numerical values for these parameters according to table 2.1.

The volumetric water content  $\theta$  is used to calculate the cumulative outflow  $Q(t)$  according to

$$Q(t) = A \int_0^l \{ \theta(h(x, t_0)) - \theta(h(x, t)) \} dx \quad (2.8)$$

where  $A$  is the core area in a horizontal cross-section. Since the porous plate is saturated in the beginning of the experiment and remains saturated during the experiment, one may assume that the differential water capacity  $C(h)$  and

Table 2.1. Parameter values used in the Mualem-van Genuchten model.

Parameter	Value(unit)
$\alpha$	0.01 (cm <sup>-1</sup> )
$n$	2.0 (-)
$\theta_r$	0.17 (-)
$\theta_s$	0.47 (-)
$K_s$	3.0 (cm/h)
$\gamma$	2.0 (-)

the hydraulic conductivity  $K(h)$  in the porous plate satisfy

$$C(h) = 0,$$

$$K(h) = K_p.$$

Consequently, the Richards equation degenerates to a simple ordinary differential equation in the fully saturated porous plate:

$$\frac{\partial^2 h}{\partial x^2} = 0, \quad b \leq x \leq l, \quad t \geq t_0, \quad (2.9)$$

where  $x = b$  is the top and  $x = l$  is the bottom of the porous plate. So, the pressure head in the porous plate is given by:

$$h(x,t) = h(b,t) + \left( \frac{x-b}{d} \right) (h_l - h(b,t)) \quad (2.10)$$

with  $h_l = h(l,t)$  the pressure head at the bottom of the porous plate and  $d = l-b$  the thickness of the porous plate. In the numerical examples we will use

as thickness of the porous plate,  $d = 0.57$  cm and as length of the soil column,  $b = 4.00$  cm.

At the interface between soil and porous plate the pressure head and the flux are continuous functions of  $x$ . So at  $x = b$ :

$$h(b^-, t) = h(b^+, t) = h(b, t), \quad (2.11)$$

$$K(h(b, t)) \left( \frac{\partial h(b^-, t)}{\partial x} - 1 \right) = K_p \left( \frac{h_i - h(b, t)}{d} - 1 \right). \quad (2.12)$$

Equation (2.12) can be used as a boundary condition for the unsaturated flow equation describing the movement of water in the soil sample. This leads to the new problem formulation:

$$C(h) \frac{\partial h}{\partial t} = \frac{\partial}{\partial x} \left( K(h) \left( \frac{\partial h}{\partial x} - 1 \right) \right), \quad 0 \leq x \leq b, \quad (2.13)$$

$$h(x, t_0) = h_0(x), \quad (2.14)$$

$$\frac{\partial h(0, t)}{\partial x} = 1, \quad (2.15)$$

$$\frac{\partial h(b, t)}{\partial x} = \frac{\left( K_p \frac{(h_i - h(b, t))}{d} - K_p + K(h(b, t)) \right)}{K(h(b, t))}. \quad (2.16)$$

Because the differential equation and one of its boundary conditions are nonlinear, the solution of this equation is approximated numerically. In the next sections some numerical solutions of this equation are given and the results are compared.

### 2.3 NUMERICAL APPROXIMATIONS

The numerical solutions of eqs. (2.13-2.16) are based on either a finite element or a finite difference approach for the space discretization. This semi-discretization of the Richards equation leads to a system of ordinary differential equations (ODE's) which is then solved by an ODE-solver. In previous work the Euler method is often used as ODE-solver (e.g., van Genuchten, 1982). A way to improve the efficiency and reliability of the numerical integration is to use a higher order ODE-solver, for example a linear multistep method or a Runge-Kutta method. Many software libraries supply routines in which these methods are implemented. These routines also include algorithms to optimize the order of the method and the step size (IMSL, 1987). Since our problem is likely to be stiff, we have chosen the implicit Gear method (backward differences up to order five).

In the finite element (FE) approach the solution  $h(x,t)$  of eqs. (2.13-2.16) is approximated by

$$h(x,t) = \sum_{j=0}^M h_j(t) v_j(x) \quad (2.17)$$

where  $v_j(x)$  are the selected basic functions. As basic functions the simple chapeau functions are used. The nonlinear coefficients  $C(h)$  and  $K(h)$  are also expanded in terms of the chapeau functions. Furthermore the mass matrix is lumped to guarantee non-oscillatory solutions (Ceila et al., 1990). Evaluation of the integrals that occur in this formulation leads to the following system of nonlinear ODE's:

$$u' = A(u)u + D(u) \quad (2.18)$$

with



$$\frac{\partial}{\partial x} \left( K(h_j) \left( \frac{\partial h_j}{\partial x} - 1 \right) \right) = \frac{1}{\Delta x} \left( K_{j+\frac{1}{2}} \left( \frac{h_{j+1} - h_j}{\Delta x} - 1 \right) - K_{j-\frac{1}{2}} \left( \frac{h_j - h_{j-1}}{\Delta x} - 1 \right) \right) + O(\Delta x^2), \quad j=0, \dots, M \quad (2.21)$$

$$\frac{\partial h_0}{\partial x} = \frac{h_1 - h_{-1}}{2\Delta x} + O(\Delta x^2), \quad (2.22)$$

$$\frac{\partial h_M}{\partial x} = \frac{h_{M+1} - h_{M-1}}{2\Delta x} + O(\Delta x^2). \quad (2.23)$$

Neglecting the higher order terms and eliminating the pressure heads  $h_{-1}$  and  $h_{M+1}$  with the boundary condition, gives the following system of ODE's:

$$u' = A(u)u + D(u) \quad (2.24)$$

with

$$A(u) = \begin{pmatrix} -\frac{K_{-1/2} + K_{1/2}}{C_0 \Delta x^2} & \frac{K_{-1/2} + K_{1/2}}{C_0 \Delta x^2} & & & & \\ \frac{K_{1-1/2}}{C_1 \Delta x^2} & -\frac{K_{1-1/2} + K_{1+1/2}}{C_1 \Delta x^2} & \frac{K_{1+1/2}}{C_1 \Delta x^2} & & & \\ & \frac{K_{2-1/2}}{C_2 \Delta x^2} & -\frac{K_{2-1/2} + K_{2+1/2}}{C_2 \Delta x^2} & \frac{K_{2+1/2}}{C_2 \Delta x^2} & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & \frac{K_{M-3/2}}{C_{M-1} \Delta x^2} & -\frac{K_{M-3/2} + K_{M-1/2}}{C_{M-1} \Delta x^2} & \frac{K_{M-1/2}}{C_{M-1} \Delta x^2} \\ & & & & \frac{K_{M-1/2} + K_{M+1/2}}{C_M \Delta x^2} & -\frac{K_{M-1/2} + K_{M+1/2}}{C_M \Delta x^2} \end{pmatrix}$$

$$D(u) = \begin{pmatrix} \frac{-(K_{-1/2} + K_{1/2})}{C_0 \Delta x} \\ \frac{K_{1/2} - K_{-1/2}}{C_1 \Delta x} \\ \frac{K_{3/2} - K_{1/2}}{C_2 \Delta x} \\ \vdots \\ \frac{K_{M-3/2} - K_{M-1/2}}{C_{M-1} \Delta x} \\ \frac{(K_{M-1/2} + K_{M+1/2}) + 2 \frac{K_0}{K_M} \left( \frac{h_1 - h_M}{d} - 1 \right) K_{M+1/2}}{C_M \Delta x} \end{pmatrix}, \quad u = \begin{pmatrix} h_0 \\ h_1 \\ h_2 \\ \vdots \\ h_{M-1} \\ h_M \end{pmatrix}$$

Generally, one approximates the conductivity between two nodes  $K_{j+1/2}$  using a linear interpolation of the pressure head,  $K_{j+1/2} = K((h_j + h_{j+1})/2)$ , or using a linear interpolation of the conductivity,  $K_{j+1/2} = ((K(h_j) + K(h_{j+1}))/2)$ . We will use FD1 and FD2 to refer to the finite difference approaches with the interpolation based on  $h$  and  $K$  respectively. It is noted that the finite element method uses the  $K$ -based interpolation of the conductivity between nodes as can be easily seen by comparing the matrices  $A$  in eq. (2.18) and eq. (2.22).

The errors due to space discretization can be analyzed more easily when the errors due to time discretization are negligible. The discretization errors in the finite element scheme as well as in the finite difference scheme are then of the order,

$$error = O(\Delta x^2). \quad (2.23)$$

One can estimate the spatial discretization errors by calculating the difference between several approximations of the pressure head, each consecutive approximation uses half the step size (in space direction) of the previous one. These differences are proportional to the spatial truncation errors, according to



$$e_{\Delta x}(x,t) \approx 1/3(h_{\Delta x}(x,t) - h_{2\Delta x}(x,t)), \quad (2.24)$$

in case the asymptotic expression (2.23) is valid in sufficient extent.

Equations (2.13-2.16) are first solved for the case that a relatively low pneumatic head,  $h_c = 250$  cm, is applied to the soil sample. The pressure head in the soil is then almost in equilibrium state after 2 hours. Therefore, the experiment is simulated from time  $t = 0$  h up to  $t = 2$  h. In a second experiment a relatively high pneumatic head,  $h_c = 1000$  cm, is used. Here, the experiment is simulated from  $t = 0$  h up to  $t = 4$  h.

Tables 2.2, 2.3 and 2.4 give the approximations of the pressure head, for a series of step sizes  $\Delta x_i$  with  $\Delta x_{i+1} = \Delta x_i/2$ ; Moreover, the differences of consecutive approximations are given, so that the validity of the asymptotic expression (2.23) can be verified. Table 2.2 gives the results of the first experiment and table 2.3 and 2.4 the results of the second experiment. In the first case the largest and smallest errors are found at  $x = 0.0$  cm and  $x = 4.0$  cm respectively. Therefore, the approximations in table 2.2 are given at these nodes. For the second case the situation is more complicated. The smallest error is for all approaches found at  $x = 0.0$  cm. The largest error is found at  $x = 4.0$  cm for the FD1 approach; for the FD2 and the FE approach the largest error is found at the last but one node  $M-1$ . Therefore, the approximations of the pressure head are given at  $x = 0.0$  cm and  $x = 4.0$  cm for the first approach and at  $x = 0.0$  cm and  $x = 3.8$  cm for the two other approaches.

Table 2.2. Pressure head,  $h(x,t)$ , at  $x = 0.0$  cm,  $t=2.0$  h and  $x=4.0$  cm,  $t=2.0$  h for several approximations, with  $h_c=250$  cm. Each consecutive approximation uses half the step size (in space direction) of the previous one. FE stands for the finite element approach, FD1 for the finite difference approach with a  $h$ -based interpolation of  $K_{j+1/2}$  and FD2 for the finite difference approach with a  $K$ -based interpolation.

	$\Delta x$	$h_{\Delta x}(0.0,2.0)$	$h_{\Delta x}-h_{2\Delta x}$	$h_{\Delta x}(4.0,2.0)$	$h_{\Delta x}-h_{2\Delta x}$
FE	0.2	-228.1091		-247.3963	
FE	0.1	-228.0937	0.0154	-247.3917	0.0046
FE	0.05	-228.0898	0.0039	-247.3905	0.0012
FE	0.025	-228.0889	0.0009	-247.3902	0.0003
FD1	0.2	-227.8942		-247.3578	
FD1	0.1	-228.0403	-0.1461	-247.3822	-0.0244
FD1	0.05	-228.0765	-0.0362	-247.3882	-0.0060
FD1	0.025	-228.0855	-0.0090	-247.3897	-0.0015
FD2	0.2	-227.9863		-247.3738	
FD2	0.1	-228.0623	-0.0760	-247.3860	-0.0122
FD2	0.05	-228.0819	-0.0196	-247.3890	-0.0030
FD2	0.025	-228.0868	-0.0049	-247.3898	-0.0008

Table 2.3. Pressure head,  $h(x,t)$ , at  $x = 0.0$  cm,  $t=4.0$  h and  $x=4.0$  cm,  $t=4.0$  h for several approximations, with  $h_c=1000$  cm. Each consecutive approximation uses half the step size (in space direction) of the previous one. FD1 stands for finite difference approach with a  $h$ -based interpolation of  $K_{j+1/2}$ .

	$\Delta x$	$h_{\Delta x}(0.0,4.0)$	$h_{\Delta x}-h_{2\Delta x}$	$h_{\Delta x}(4.0,4.0)$	$h_{\Delta x}-h_{2\Delta x}$
FD1	0.2	-244.513		-249.957	
FD1	0.1	-272.561	-28.048	-291.806	-41.849
FD1	0.05	-293.025	-20.464	-339.898	-48.092
FD1	0.025	-305.368	-12.343	-396.684	-56.786
FD1	0.0125	-311.897	-6.529	-465.039	-68.355
FD1	0.00625	-315.110	-3.213	-549.294	-84.255
FD1	0.003125	-316.634	-1.524	-658.944	-109.650
FD1	0.0015625	-317.356	-0.722	-888.145	-224.201
FD1	0.00078125	-317.654	-0.298	-998.400	-110.255
FD1	0.000390625	-317.665	-0.011	-998.408	-0.009

Table 2.4. Pressure head,  $h(x,t)$ , at  $x = 0.0$  cm,  $t=4.0$  h and  $x=3.8$  cm,  $t=4.0$  h for several approximations, with  $h_c=1000$  cm. Each consecutive approximation uses half the step size (in space direction) of the previous one. FE stands for the finite element approach, FD2 for the finite difference approach with a  $K$ -based interpolation of  $K_{j+1/2}$ .

	$\Delta x$	$h_{\Delta x}(0.0,4.0)$	$h_{\Delta x}-h_{2\Delta x}$	$h_{\Delta x}(3.8,4.0)$	$h_{\Delta x}-h_{2\Delta x}$
FE	0.2	-324.491		-590.942	
FE	0.1	-320.710	3.781	-546.202	44.740
FE	0.05	-318.987	1.732	-526.837	19.364
FE	0.025	-318.205	0.773	-519.071	7.766
FE	0.0125	-317.875	0.330	-516.141	2.930
FE	0.00625	-317.741	0.134	-515.135	1.006
FE	0.003125	-317.691	0.050	-514.823	0.312
FE	0.0015625	-317.674	0.017	-514.734	0.089
FD2	0.2	-324.376		-590.639	
FD2	0.1	-320.673	3.703	-546.084	44.555
FD2	0.05	-318.964	1.709	-526.791	19.293
FD2	0.025	-318.200	0.764	-519.048	7.743
FD2	0.0125	-317.872	0.328	-516.126	2.922
FD2	0.00625	-317.739	0.133	-515.126	1.000
FD2	0.003125	-317.690	0.049	-514.820	0.306
FD2	0.0015625	-317.673	0.017	-514.733	0.087

In the first experiment the difference between consecutive approximations reduces each time with a factor 4. This implies that the asymptotic expression (2.23) is valid in sufficient extent to have an accurate error estimate. For this situation the results are very well for all approaches and they are comparable in precision. In the second experiment the difference between consecutive approximations does not reduce with a factor 4. Therefore, an accurate estimate of the errors cannot be given. However, it is obvious that the errors are very large, especially near the lower boundary. The results of the FD2 and the FE approach (table 2.4), which both use a  $K$ -based interpolation of  $K_{j+1/2}$  are almost identical. On the other hand the results of the FD1 and FD2 approach (tables 2.3 and 2.4), using the  $h$ - and  $K$ -based interpolation of  $K_{j+1/2}$  respectively, are quite different, especially near the lower boundary. This suggests that the accuracy of the solution strongly depends on the definition of the conductivity between nodes. From figure 2.2 it is seen that near the lower boundary the pressure head has a steep gradient (drying front) and that in this region the pressure head lies between approximately -300 cm and -1000 cm. For these  $h$ -values the conductivity is a strongly nonlinear function of  $h$  (e.g. Van Dam, 1990). So, to increase the accuracy of the solution one has to approximate the fluxes near and at the boundary more precisely. This may be done by using a smaller step size in this region, that is, by using a non-uniform grid. The advantage of a non-uniform grid is that the accuracy thus obtained is comparable to the accuracy obtained with the finest grid.

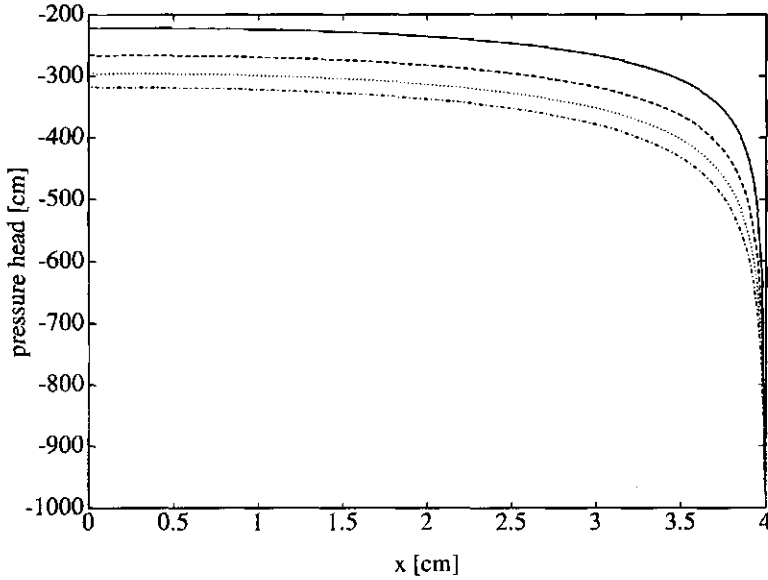


Figure 2.2. Approximated pressure head  $h(x,t)$  for  $t=1.0$  h (solid line),  $t=2.0$  h (dashed line),  $t=3.0$  h (dotted line) and  $t=4.0$  h (dotted-dashed line), with  $h_c=1000$  cm,  $\Delta x=0.025$  cm and using the FE approach.

We have implemented a scheme with a non-uniform grid using the finite element approach; the same result will be probably obtained using the finite difference approach. Between the lower boundary and a distance  $d_1$  from the lower boundary, we use a step-size  $\Delta x_2$ ; in the other region a step-size  $\Delta x_1$ . As a first choice we take  $d_1 = 0.1$  cm. The results of these runs are given in figures 2.3 and 2.4 and table 2.5.

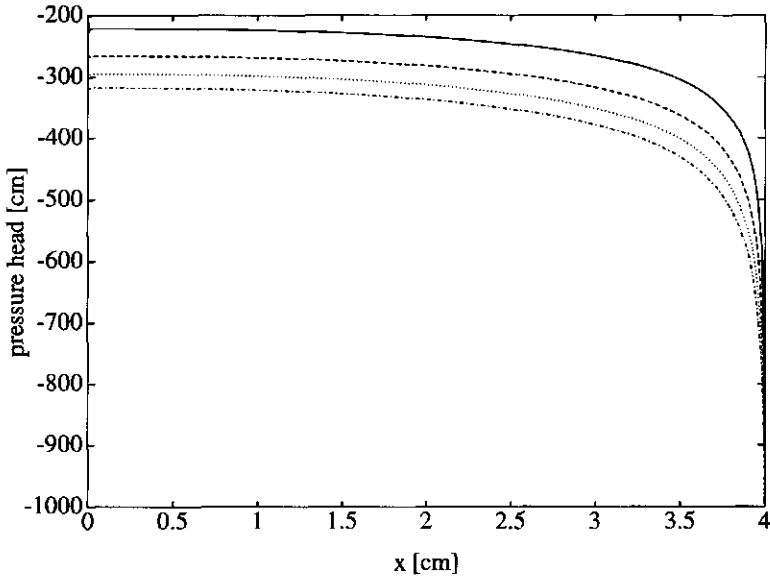


Figure 2.3. Approximated pressure head  $h(x,t)$  for  $t=1.0$  h (solid line),  $t=2.0$  h (dashed line),  $t=3.0$  h (dotted line) and  $t=4.0$  h (dotted-dashed line), with  $h_c=1000$  cm,  $\Delta x_1=0.195$  cm,  $\Delta x_2=0.005$  cm and using the FE approach with a nonuniform grid.

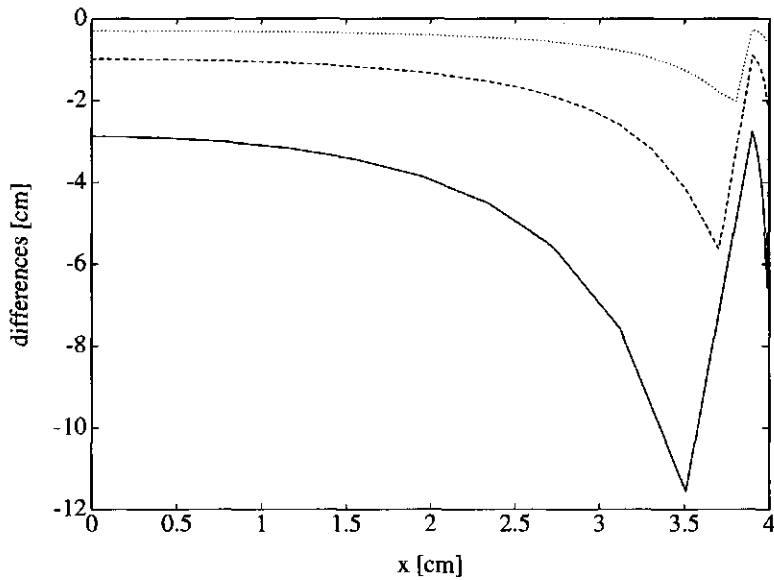


Figure 2.4. Difference between several approximations of the pressure head at  $t=4.0$  h, each consecutive approximation uses half the step sizes  $\Delta x_1$  and  $\Delta x_2$  of the previous one.



Table 2.5. Pressure head,  $h(x,t)$ , at  $x = 0.0$  cm,  $t=4.0$  h and  $x=3.99$  cm,  $t=4.0$  h for several approximations, with  $h_c=1000$  cm. Each consecutive approximation uses half the step sizes  $\Delta x_1$  and  $\Delta x_2$  (in space direction) of the previous one. FEN stands for the finite elements approach with a nonuniform grid.

	$\Delta x_1$	$\Delta x_2$	$h_{\Delta x}(0.0,4.0)$	$h_{\Delta x}-h_{2\Delta x}$	$h_{\Delta x}(3.99,4.0)$	$h_{\Delta x}-h_{2\Delta x}$
FEN	0.39	0.01	-321.933		-855.425	
FEN	0.195	0.005	-319.058	2.875	-848.844	6.581
FEN	0.0975	0.0025	-318.070	0.988	-846.725	2.119
FEN	0.04875	0.00012	-318.774	0.296	-846.129	0.596
FEN	0.024375	0.0000625	-317.694	0.080	-845.972	0.157

The performance of this scheme is very well because the differences reduce again with a factor 4. Also the estimated errors are very small, although we only used 1/4 of the number of nodes compared with the other approaches. From figure 4 it can also be seen that the errors in the area, where a larger step size is used, are comparable with the errors in the area where the step-size is small. This suggests that the choice of the ratio  $\Delta x_1 / \Delta x_2$  and that of  $d_1$  is good. If  $d_1$  is larger the number of nodes in the small region increases, because the same  $\Delta x_2$  is needed to obtain the same accuracy. However, if  $d_1$  is smaller the number of nodes in the small region decreases, but the number of nodes in the large region must increase to obtain the same accuracy.

## 2.4 CONCLUSIONS

In this paper it is shown that numerical approximations may be very poor in regions where the spatial derivatives of the pressure head are large. In outflow experiments this occurs if the pneumatic head imposed on the pressure cell is high,  $h_c = 1000$  cm. To analyze the errors due to the spatial discretization one first has to make the errors due to time discretization comparatively small. For this purpose a high order ODE-solver, adopted from a standard software library, is used. The advantages of a standard ODE-solver are a higher efficiency and reliability and a facility for error control.

Three approximation schemes have been used for the spatial discretization of Richards equation. One approximation scheme was based on a finite element approach and the other two schemes were based on finite differences. All three schemes have a second order spatial discretization error. The difference between the two finite difference schemes is the approximation of the conductivity between nodes. One scheme uses an approximation of  $K_{j+1/2}$  based on a linear interpolation of  $h$ ; the other uses an approximation based on a linear interpolation of  $K$ . The fact that the results of the two finite difference schemes are quite distinct in regions where the solution has a steep front, suggests that the solution strongly depends on the non-linear hydraulic conductivity relation. If the pressure head has a steep front, the approximations of the nodal fluxes may be very poor and result in large discretization errors. In the finite element approach the same approximation of the nodal fluxes is used as in the finite difference approach with the interpolation based on  $K$ . This explains why the results of these two schemes are almost identical. It also confirms our point of concern that the spatial discretization of the flux may be the cause of large errors at and near the lower boundary.

Because steep fronts of the pressure head occur also in related problems, for example infiltration into dry soils (Ceila et al., 1990), the numerical approximation of the fluxes at and near the wetting front will be very poor in

case a uniform grid is used and  $\Delta x$  is not sufficiently small. It is then almost necessary to use non-uniform adaptive grid methods. Although an adaptive grid is very accurate it is not very fast and therefore it is not suitable to be used in inverse modeling problems. It is shown that a simple non-uniform 'fixed' grid is sufficient to obtain very accurate numerical simulations of the ONE-STEP experiment. The reason that this simple grid is sufficient is that the position of the steep front is known beforehand as opposed to the moving wetting front in the infiltration problem. The advantage of a non-uniform grid with respect to a uniform grid is that only 1/4 of the number of nodes is needed to yield the same accuracy.

## REFERENCES

- Ceila, M.A., E.T. Bouloutas and R.L. Zarba, A general mass-conservative numerical solution for the unsaturated flow equation, *Water Resour. Res.*, 26, 1483-1496, 1990.
- Hanks, R.J., A. Klute and E. Bresler, A numeric method for estimating infiltration, redistribution, drainage and evaporation of water from soil, *Water Resour. Res.*, 5, 1064-1069, 1969
- IMSL, *IMSLMATH/LIBRARY* user's manual, 640-652, 1987
- Kool, J.B., J.C. Parker and M.T. van Genuchten, Determining soil hydraulic properties from one-step outflow measurements by parameter estimation, I, theory and numerical studies, *Soil Sci. Soc. Am. J.*, 49, 1348-1354, 1985.
- Kool, J.B. and J.C. Parker, Analysis of the inverse problem for transient unsaturated flow, *Water Resour. Res.*, 24, 513-522, 1988.
- Milly, P.C.D., A mass-conservative procedure for time-stepping in models of unsaturated flow, *Adv. Water Resour.*, 8, 32-36, 1985.

- Mualem, Y, A new model for predicting the hydraulic conductivity of unsaturated porous media, *Water Resour. Res.*, 12, 513-522, 1976.
- Van Dam, J.C., J.N.M. Stricker and P. Droogers, From One-Step to Multi-Step, determination of soil hydraulic functions by outflow experiments, Rep. 7, Dep. of Hydrology, Soil Physics and Hydraulics, Agric. Univ. Wageningen, 1990.
- Van Genuchten, M.T., A closed-form equation for predicting the hydraulic conductivity of unsaturated soils. *Soil Sci. Soc. Am. J.*, 44, 892-898, 1980.
- Van Genuchten, M.T., A comparison of numerical solutions of the one-dimensional unsaturated-saturated flow and transport equation, *Adv. Water Resour.*, 5, 47-55, 1982.
- Wösten, J.H.M. and M.T. van Genuchten, Using texture and other soil properties to predict the unsaturated soil hydraulic functions, *Soil Sci. Soc. Am. J.*, 52, 1762-1770, 1988.

# **Chapter 3**

## **IDENTIFICATION OF THE MOVEMENT OF WATER IN UNSATURATED SOILS; THE PROBLEM OF IDENTIFIABILITY OF THE MODEL**

### **Abstract**

The estimation of model parameters using nonlinear regression techniques is one of the aspects of inverse modeling and is known as the identification problem, the solution of which may be non-unique. The main causes of this non-uniqueness are the structure of the model and the design of the input signal. It will be shown that the parameters can be estimated only if a model with different parameter values yields different output signals. A model that has this feature is called identifiable. As an example, the identification of a model for the movement of water in unsaturated soils is used.

This model appears to be non-identifiable, which results in non-unique solutions.

### 3.1 INTRODUCTION

Computers have made it possible to build complex models of physical dynamical processes. To use such a model, for example for prediction or simulation of a process, the model has to be identified first. A model that can often be used to describe a dynamical process is a partial differential equation of the form

$$\frac{\partial u(x,t)}{\partial t} = f(\nabla_x^2 u(x,t), \nabla_x u(x,t), u(x,t), v(x,t), a), \quad (3.1)$$

where the state variable  $u(x,t)$  is a function of space and time. The state variable often has physical significance, for example temperature or pressure. The input of the dynamical process is denoted by  $v(x,t)$  and is also a function of space and time, for example a source of heat. The model parameters (coefficients) are denoted by the vector  $a$ . In most practical cases the partial differential eq. (3.1) cannot be solved analytically. The solution of eq. (3.1) is then often approximated using a semi-discrete finite element or a finite difference method. This semi-discretisation approach leads to a system of ordinary differential equations for the discretised state vector  $u(t)=(u(x_1,t), u(x_2,t), \dots, u(x_n,t))$ ,

$$\frac{\partial u(t)}{\partial t} = f_1(u(t), v(t), a). \quad (3.2)$$

The observations  $y(kT)$  of the states are often restricted to be sampled at discrete and equidistant instants of time  $kT$ ,  $k=1, \dots, N$ , where  $T$  is the sampling time. It is assumed that the output is disturbed by random measure-

ment errors  $e(kT)$ . The model for the observations is then given by

$$y(kT) = f_2(u(kT), a) + e(kT). \quad (3.3)$$

In the remainder of this paper a model is called deterministic if the disturbances are equal to 0 and stochastic otherwise.

Before the parameter vector  $a$  can be estimated in the case of a stochastic model the input signals  $v(t)$  have to be designed and the initial values of the states have to be known. The set of the input signals is called the experimental design.

In the last decade, studies of the movement of water in unsaturated soils have encountered problems with the uniqueness of the parameter estimates (Hornung, 1983, Kool et al., 1985 and van Dam et al., 1990). The experiment mostly used to identify this process is known as the ONE-STEP experiment of Kool et al. (1985). In this experiment outflow measurements are taken from a soil core. In conjunction with an inverse modeling procedure these outflow data are used to identify the unknown model parameters. In this paper it will be shown that the non-uniqueness of the parameter estimates is not due to a bad choice of the optimization algorithm; it is merely a consequence of the structure of the model and the design of the experiment. Because the latter can be avoided, it is worthwhile to analyze the cause of this non-uniqueness.

In some other papers, where rainfall-runoff models are discussed, non-uniqueness of the parameter estimates is also said to be one of the main problems (e.g. Kleissen et al., 1990, Sorooshian and Gupta, 1985). Sorooshian and Gupta (1985) suggested using the sensitivity ratio to analyze the influence of the structure of the model on the non-uniqueness of the parameter estimates, but this ratio does not take into account the influence of significant numerical errors in the sensitivity matrix. For this purpose the concept of identifiability may be more suitable (Walter, 1992). In section 3.3 this concept is further developed so as to take also numerical errors into

account.

### 3.2 ONE-DIMENSIONAL TRANSIENT WATER FLOW

The ONE-STEP experiment was developed to identify one-dimensional transient water-flow in porous media in a simple and fast way. The usual way to model the movement of water in unsaturated soil is to use the Richards equation with the pressure head as state variable,

$$C(h) \frac{\partial h}{\partial t} = \frac{\partial}{\partial x} \left( K(h) \left( \frac{\partial h}{\partial x} - 1 \right) \right), \quad (3.4)$$

where  $x$  is the vertical coordinate taking positive downward (cm),  $h(x,t)$  is the pressure head at time  $t$  at point  $x$  in the medium (cm),  $C(h)$  is the differential water capacity of the soil at a pressure head  $h$  (cm) and  $K(h)$  is the hydraulic conductivity of the soil at a pressure head  $h$  (cm h<sup>-1</sup>).

The differential water capacity and the hydraulic conductivity in unsaturated soils depend on the pressure head in the soil. The expressions proposed by Mualem (1976) and van Genuchten (van Genuchten, 1980) are often used to model these relationships and are given by

$$C(h) = \alpha m (\theta_s - \theta_r) S_e^{1/m} (1 - S_e^{1/m})^m (1 - m)^{-1}, \quad (3.5)$$

$$K(h) = K_s S_e^{\gamma} \left( 1 - (1 - S_e^{1/m})^m \right)^2 \quad (3.6)$$

with

$$S_e = (1 + |\alpha h|^n)^{-m},$$

$$m = 1 - 1/n.$$



The model parameters  $\alpha$ ,  $n$ ,  $\theta_r$ ,  $\theta_s$ ,  $K_s$ ,  $\gamma$  are the parameters to be estimated. The parameter  $\theta_s$  (-) is the saturated water content ( $\theta$  at  $h = 0$  cm),  $K_s$  (cm/h) is the saturated hydraulic conductivity. The parameters  $\alpha$  (cm<sup>-1</sup>),  $n$  (-),  $\theta_r$  (-) and  $\gamma$  (-) have no clear physical significance.

The initial values of the pressure head and the boundary conditions belonging to eq. (3.4) can be derived from the experimental setup. The experiment consists of a soil sample and a porous plate placed in a pressure cell (Kool et al., 1985). Initially, the water in the soil core is in a state of equilibrium and the volumetric water content in the soil is high (almost saturated). The porous plate is saturated. Then an additional pressure  $h_c(t)$  is applied to the pressure cell and water seeps out. Since no water infiltrates in the sample, there is a zero flux condition,

$$q(0,t) = -K(h) \left( \frac{\partial h}{\partial x} - 1 \right) \Big|_{x=0} = 0,$$

at the top of the soil core, which gives boundary condition (3.8). Assuming that the pressure at the bottom of the plate is atmospheric, the pressure  $h_c(t)$  leads to the boundary condition (3.9) at the bottom of the porous plate:

$$h(x,t_0) = h_0(x), \quad (3.7)$$

$$\frac{\partial h(0,t)}{\partial x} = 1, \quad (3.8)$$

$$h(L,t) = -h_c(t). \quad (3.9)$$

Because the porous plate is saturated in the beginning of the experiment and remains saturated during the experiment the capacity  $C(h)$  and the conductivity  $K(h)$  in the porous plate are assumed to satisfy

$$C(h) = 0,$$

$$K(h) = K_p.$$

The output signal, the cumulative outflow  $Q(t)$ , can be calculated by integrating the volumetric water content over the soil sample, according to

$$Q(t) = A \int_0^L \{ \theta(h_0(x)) - \theta(h(x,t)) \} dx \quad (3.10)$$

$$\theta = \theta_r + S_e(\theta_s - \theta_r)$$

where  $A$  is the core area perpendicular to the flow and  $\theta$  is the volumetric watercontent and. The cumulative outflow is sampled with a constant sampling time from time  $t_0$  up to  $t_e$ . The sampling time is then equal to

$$T = \frac{t_e - t_0}{N}, \quad (3.11)$$

where  $N$  is equal to the number of samples.

By discretizing the space variable in the mathematical description of the process given above, the problem can be reformulated into the general form of the state-space model given by eq. (3.2) (Mous, 1990). It is then clear that the pressure  $h_c(t)$  imposed on the soil sample is the input signal of the model.

### 3.3 IDENTIFIABILITY

After a model structure is specified, the unknown parameters can be

estimated using a nonlinear regression program. However, it is worthwhile to analyze the identifiability of the model first. The purpose of this analysis is to study whether the parameters can be estimated uniquely given the input and output signals.

To study identifiability in a more mathematical sense, some definitions have to be given. In this paper only the continuous-time state-space model, given by eqs. (3.2-3.3), is considered. If the structures of the functions  $f_1$  and  $f_2$  are specified, but the numerical values of the parameter vector  $a$  are unknown, the models with different parameter vector form a class. More precisely, the model class  $M$  is the range of the model, specified by the functions  $f_1$  and  $f_2$ , where the parameter vector  $a$  varies over the parameter space  $A$ .

Alternatively, If the input signal  $v(t)$  of the process and the initial values of the state vector are known, the output signal of the deterministic model ( $e(kT) = 0$ ) is a function of the parameter vector  $a$ . In this way eqs. (3.2-3.3) define a mapping of the parameter space to the response space.

### *Definition 1*

Let a model structure be specified by the functions  $f_1$  and  $f_2$  and let the initial values of the state vector  $u(t_0)$  and the input signal  $v(t)$  be known. The deterministic model is said to be globally identifiable if different  $a$  values yield different response vectors.

This definition of global identifiability is similar to the definition given by several other authors (e.g., Nguyen and Wood, 1982, Distefano and Cobelli, 1980). The importance for a model to be globally identifiable is clear, because it is a requirement to obtain a unique solution of the parameter estimating problem (Walter, 1985). However, it may be very difficult to prove that a model is globally identifiable, because the functions  $f_1$  and  $f_2$  in

physical problems are often nonlinear.

In most practical situations, it may be easy to show that a model is locally identifiable at a specified  $a^* \in A$  (Bellman and Åström, 1970).

*Definition 2*

Let a model structure be specified by the functions  $f_1$  and  $f_2$  and let the initial values of the state vector  $u(t_0)$  and the input signal  $v(t)$  be known. For any  $a^* \in A$  the deterministic model is said to be locally identifiable at  $a^*$  if the function

$$V(a) = \sum_{k=1}^N (y(kT; a^*) - y(kT; a))^2, \quad a \in A, \quad (3.12)$$

has a local minimum, equal to 0, in  $a^*$ .

In literature this definition of local identifiability has become known as least-squares identifiability (Distefano and Cobelli, 1980). Usually one checks only if the model structure is locally identifiable at an initial guess,  $a^* = a_0$ , or at an estimated value,  $a^* = \hat{a}$ , produced by the regression program. Of course, it is clear that in case the function  $V(a)$  has a global minimum in  $a^*$  for every  $a^* \in A$ , the model is also globally identifiable in the sense of the first definition.

Several reasons can make a model non-identifiable. The two important reasons are a bad design of the input signal and overparameterization of the functions  $f_1$  and  $f_2$  (Bellman and Åström, 1970, Ljung, 1987). This latter reason causes the model also to be not globally identifiable. In this case it is impossible to estimate all parameters, and other output signals have to be measured or the model structure has to be reparameterized. If the first reason causes a non-identifiable model, another design of the input signal (experimental design) may be sufficient. However, it may be better to measure other

output signals as well.

Assuming that the first and second order derivatives of the output vector with respect to the parameter vector exist, conditions for these derivatives can be formulated for the nonlinear function  $V(a)$  to have a local minimum at  $a^*$  (e.g., Scales, 1985). The first necessary condition is that the gradient  $g(a)$  is equal to 0 at  $a^*$ . The second condition is that the Hessian matrix  $H(a)$  is positive definite at  $a^*$ . These two conditions are sufficient for a minimum at  $a^*$ . The gradient and Hessian matrix of  $V(a)$  are given by

$$g_i(a) = \sum_{k=1}^N -2(y(kT; a^*) - y(kT; a)) \left( \frac{\partial y(kT; a)}{\partial a_i} \right), \quad (3.13)$$

$$h_{ij}(a) = \sum_{k=1}^N \left\{ -2(y(kT; a^*) - y(kT; a)) \left( \frac{\partial^2 y(kT; a)}{\partial a_i \partial a_j} \right) + 2 \left( \frac{\partial y(kT; a)}{\partial a_i} \right) \left( \frac{\partial y(kT; a)}{\partial a_j} \right) \right\}. \quad (3.14)$$

From eq. (3.13) it follows that  $g(a^*) = 0$ , so the first condition is always fulfilled. The second condition remains but it can be simplified, because the first term of the right-hand side,

$$-2(y(kT; a^*) - y(kT; a)) \left( \frac{\partial^2 y(kT; a)}{\partial a_i \partial a_j} \right)$$

is equal to 0 at  $a = a^*$ . The Hessian matrix at  $a^*$  is thus equal to

$$H(a^*) = 2X(a^*)'X(a^*). \quad (3.15)$$

Here, the  $N \times p$  matrix  $X(a)$  denote the partial derivatives of  $y(kT; a)$ ,  $k=1, \dots, N$ ,

with respect to the parameters  $a_i$ ,  $i=1,\dots,p$ . This matrix is often called the sensitivity matrix. In statistical literature this matrix is also called the design matrix. The condition that the Hessian matrix  $H(a^*)$  is positive definite is equivalent to the condition that the sensitivity matrix  $X(a^*)$  is of full rank. The determination of the rank of a matrix is numerically difficult since the matrix  $X(a^*)$  can only be calculated with finite precision. Especially in highly nonlinear models, such as the model for the movement of water in unsaturated soil, the truncation and roundoff errors in the numerical approximations of the solutions of these models are often significant (Mous, 1990). The best way to analyze the rank of a noisy matrix is to use the concept of  $\epsilon$ -rank (Dongarra et al., 1979).

### *Definition 3*

Let  $X$  be a  $m \times n$  matrix and  $Z_k$  be a set of  $m \times n$  matrices of at least rank  $k$ . The distance  $d_k$  between  $X$  and  $Z_k$  is defined as

$$d_k := \min_{Z \in Z_k} \|X - Z\|,$$

where  $\|\cdot\|$  is the quadratic-norm. The  $\epsilon$ -rank of  $X$  is then defined as the smallest value of  $k$  so that  $d_k \leq \epsilon \|X\|$ ,  $\epsilon \geq 0$

Owing to the numerical errors, the matrix  $X(a^*)$  will in general be of full rank. Let  $\Delta X(a^*)$  be the error in  $X(a^*)$  and let  $E$  be an  $N \times p$  matrix with  $\|E\| \leq \|\Delta X\|$ . The question then is not whether  $X(a^*)$  is of full rank but whether the matrix  $X(a^*) + E$  is rank deficient ( $\text{rank}(X(a^*) + E) \leq p-1$ ), because  $X(a^*) + E$  could be the "error-free" sensitivity matrix.

By choosing  $\epsilon = \|\Delta X(a^*)\| / \|X(a^*)\|$ , the  $\epsilon$ -rank  $r_\epsilon$  of  $X(a^*)$  can be used as a rank test for  $X(a^*)$ . If  $r_\epsilon$  is smaller than the number of unknown parameters of the model then there is a matrix  $E$  with  $\|E\| \leq \|\Delta X\|$  such

that  $X(a^*) + E$  is rank deficient. In that case the model is said to be numerically not identifiable at the parameter vector  $a^*$ .

The  $\epsilon$ -rank  $r_\epsilon$  may be calculated using a singular value decomposition (S.V.D.) of the matrix  $X(a^*)$ .

$$X(a^*) = U\Sigma V', \quad (3.17)$$

where  $U$  is an  $N \times N$  orthonormal matrix and  $V$  is a  $p \times p$  orthogonal matrix. The diagonal matrix  $\Sigma$  contains the singular values  $\sigma_1, \sigma_2, \dots, \sigma_p$ , with  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p$ . The  $\epsilon$ -rank of  $X(a^*)$  is the number of  $\sigma_i$  which are larger than  $\eta = \epsilon \|X(a^*)\|$ , i.e. larger than  $\epsilon \sigma_1$ .

### 3.4 EXPERIMENTAL DESIGN

If the calculations, as described in the previous section, are done for different input signals, not only identifiability can be studied; these experimental designs can also be compared with each other with respect to some optimality criteria. For example, a criterion may be used for designing an experiment that gives richer and more informative output signals.

In this study it is assumed that the disturbances  $e(kT)$  are independently and identically distributed. By definition the least squares estimate  $\hat{a}$  is obtained by minimizing the sum of squares

$$S(a) = \sum_{k=1}^N (z(kT) - y(kT; a))^2. \quad (3.18)$$

Here, the observed responses are denoted by  $z(kT)$  and the model outcome by  $y(kT; a)$ . The gradient  $g_s(a)$  and Hessian matrix  $H_s(a)$  of  $S(a)$  are given by

$$g_{s,i}(a) = \sum_{k=1}^N -2(z(kT) - y(kT;a)) \left( \frac{\partial y(kT;a)}{\partial a_i} \right), \quad (3.19)$$

$$h_{s,ij}(a) = \sum_{k=1}^N \left\{ -2(z(kT) - y(kT;a)) \left( \frac{\partial^2 y(kT;a)}{\partial a_i \partial a_j} \right) + 2 \left( \frac{\partial y(kT;a)}{\partial a_i} \right) \left( \frac{\partial y(kT;a)}{\partial a_j} \right) \right\}. \quad (3.20)$$

It is assumed that the data set is large and that the i.i.d. assumption of the disturbances is valid. The contribution of the first terms on the right-hand side of eq. (3.20) to the Hessian matrix are then negligible, because successive terms

$$-2(z(kT) - y(kT;a)) \left( \frac{\partial^2 y(kT;a)}{\partial a_i \partial a_j} \right)$$

cancel through sign differences. Asymptotically as  $N \rightarrow \infty$  and  $T = (t_e - t_0)/N \rightarrow 0$  the contribution in  $N^1 H$  of the first terms on the right-hand side will tend to zero and the second to a constant. The contribution of the first terms is also negligible if the output  $y(kT;a)$  depends almost linearly on the parameter vector  $a$  so that the second order derivatives are almost equal to 0. Therefore, the model must be at least locally identifiable at  $\hat{a}$  to obtain unique parameter estimates. However, because  $\hat{a}$  is not known beforehand, the best one can do is to test whether the model, given the input design  $v(t)$ , is locally identifiable at an initial guess  $a_0$ . This condition is used as a first criterion. Since this criterion does not select the "best" input design, a second criterion is needed. The D-optimal criterion is often used to select an input design and is attained by the design that minimizes the determinant



$|X(a)'X(a)|^{-1}$  (e.g. Box and Lucas, 1950).

If the Hessian matrix can be approximated by eq. (3.15), with  $a^* = \hat{a}$ , the asymptotical covariance matrix may be estimated by (Bates & Watts, 1988)

$$C = \frac{S(\hat{a})}{N-p} (X(\hat{a})'X(\hat{a}))^{-1}, \quad (3.21)$$

The determinant used to find the D-optimal design is proportional to the determinant of this estimated asymptotical covariance matrix, with  $a = \hat{a}$ , so it can be seen as a measure of the uncertainty in the parameter values. Although the objective is to evaluate the D-optimal criterion at the "true" parameter vector, again an initial guess  $a_0$  has to be used instead.

### 3.5 RESULTS

To decide whether it is meaningful to carry out the ONE-STEP experiment, we will use the tools presented in the previous section. To do this analysis, some prior information albeit hypothetical about the parameters in the Mualem-van Genuchten model is needed (table 3.1).

In the experiments the soil core is almost saturated soil. Therefore, in this numerical experiment the initial state is set to  $h(x, t_0) = -50$  cm. As input signal  $h_c(t)$  a step function is used. We will compare several experiments with the level of the step function between 100 cm and 1000 cm. In addition, a multiple step function is also used as input signal (Van Dam et al., 1990). The levels of this step function are set to 75 cm, 150 cm, 250 cm and 1000 cm at times 0.0 h, 2.0 h, 5.5 h and 7.0 h respectively. A sampling time  $T$  equal to  $t_e/100$  is used, with  $t_e$  is equal to the duration of the experiment. The time  $t_e$  is chosen in such a way that the water content in the soil sample at time  $t_e$  is almost in a state of equilibrium.

Table 3.1. Parameter values used in the Mualem-van Genuchten model

Parameter	Value(unit)
$\alpha$	0.01 (cm <sup>-1</sup> )
$n$	2.0 (-)
$\theta_r$	0.17 (-)
$\theta_s$	0.47 (-)
$K_s$	3.0 (cm/h)
$\gamma$	2.0 (-)

The sensitivity matrix  $X(a_0)$  is approximated simultaneously with the cumulative outflow  $Q(kT; a_0)$  (Caracotsios and Stewart, 1985), using a finite element method, with variable step size, for the spatial discretization (Mous, 1990) and the Crank-Nicolson method for the time integration.

To estimate the error matrix  $\Delta X(a_0)$ , the errors due to time discretization are made small in comparison with the errors due to space discretization, by using a very small time-step. The errors due to space discretization are estimated by calculating the difference of several approximations, each consecutive approximation uses half the step-size (in space direction) of the previous one. Because the errors in the columns of the sensitivity matrix are not of the same order of magnitude, the columns are scaled (Dongarra et al., 1979). The scaling matrix  $M$  is given by

$$m_{ij} = \begin{cases} \left( \frac{1}{N} \sum_{k=1}^N |\Delta x_{k,i}| \right)^{-1} & i = j \\ 0 & i \neq j \end{cases} \quad (3.22)$$

Table 3.2. Rank analysis of the sensitivity matrix  $X(a_0)$ 

Level $h_c$	$t_e$	$\  \Delta X(a^*) \ $	$\  \Delta \tilde{X}(a^*) \ $	$\epsilon$ -rank
100	4.0	6.34	0.024	4
200	4.5	10.66	0.150	4
300	5.0	21.37	0.214	3
400	5.5	47.33	0.206	3
500	6.0	50.72	0.207	3
600	6.5	43.64	0.184	3
700	7.0	34.89	0.148	3
800	7.5	27.16	0.116	3
900	8.0	22.38	0.098	3
1000	8.5	19.61	0.088	3
MS	10.0	10.31	0.111	5

The levels of the multiple step function, MS, were set to 75 cm, 150 cm, 250 cm and 1000 cm at times 0.0 h, 2.0 h, 5.5 h and 7.0 h respectively.

The scaled sensitivity matrix  $\tilde{X}(a_0) = X(a_0)M$  is used to analyze the rank of the sensitivity matrix.

Table 3.2 shows that the  $\epsilon$ -rank for all the input designs is smaller than the number of parameters and it may be concluded that the model is numerically not identifiable at  $a_0$ . In figures 3.1(a)-(f) the columns of the sensitivity matrix of the experiment with  $h_c = 200$  cm are plotted. The shape of figures 3.1(b), (c) and (d) are almost equal or are mirror images. This indicates a linear dependence between two columns of the sensitivity matrix and consequently a rank deficient sensitivity matrix. Linear dependency is also found between the fifth and sixth column of the sensitivity matrix (see figures

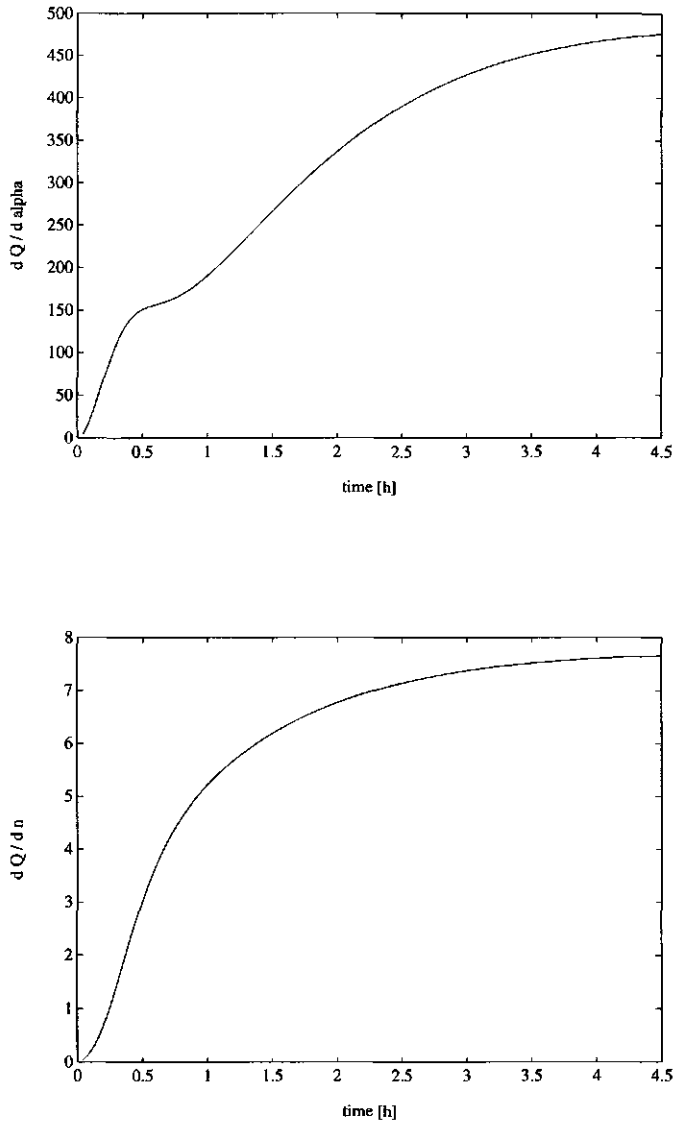


Figure 3.1 Sensitivity analysis of the ONE-STEP experiment with the level of the input signal  $h_c = 200$  cm: (a)  $\partial Q/\partial\alpha$ ; (b)  $\partial Q/\partial n$ ; (c)  $\partial Q/\partial\theta_s$ ; (d)  $\partial Q/\partial\theta_r$ ; (e)  $\partial Q/\partial K_s$ ; (f)  $\partial Q/\partial l$ .

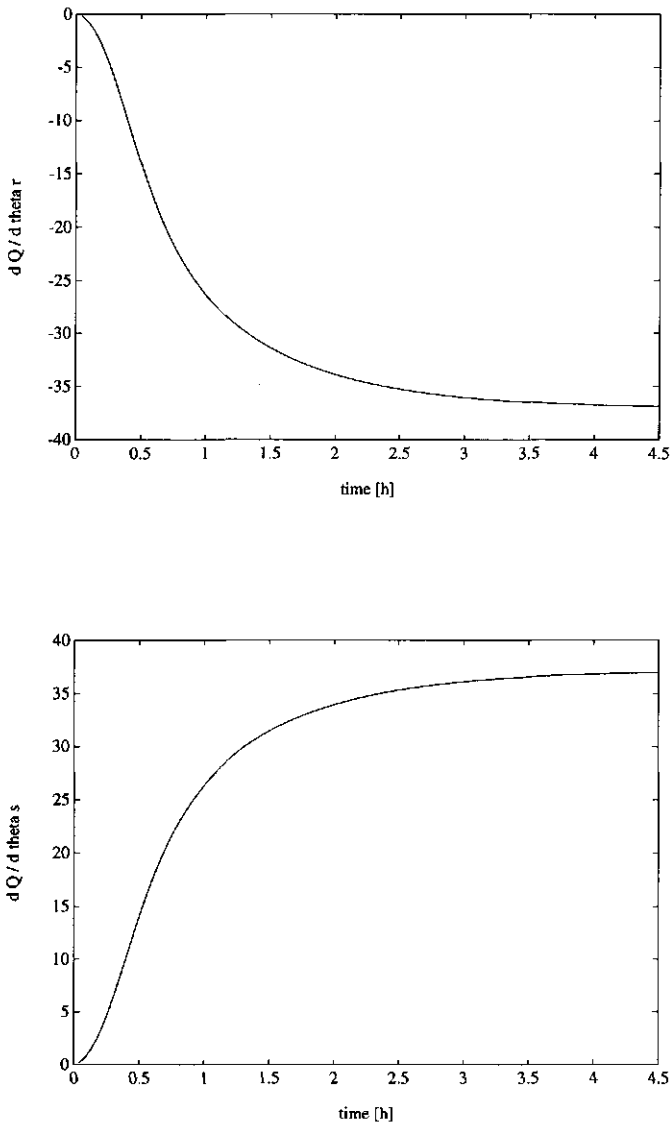


Figure 3.1 Continued.

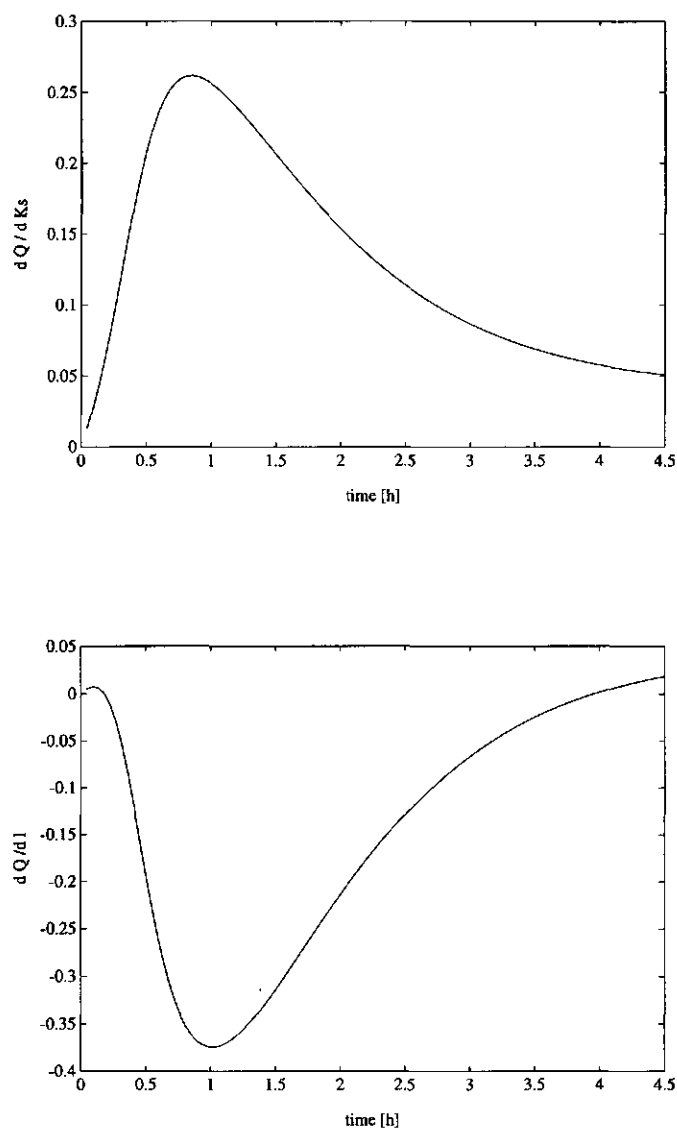


Figure 3.1 Continued.

3.1(e) and (f)). Because of these 2 groups, the sensitivity matrix will have at least two large singular values.

This non-identifiability is caused partly by the structure of the model. Equations (3.4-3.6) and (3.10) show that the output  $Q(kT)$  depends on the difference of  $\theta_s$  and  $\theta_r$  only. Therefore, at least one singular value will be equal to 0. This does not explain, however, why the  $\varepsilon$ -rank is equal to 3 or 4 for most designs, which mean that only 3 or 4 parameters can be estimated. Because the  $\varepsilon$ -rank with the multiple step function as input signal is equal to 5, non-identifiability of the model is caused here also by the design of the experiment. Suppose that the difference of  $\theta_s$  and  $\theta_r$  is of interest, then the parameter vector can be reduced by one parameter. The model, given the multiple step function as input signal, is then numerically identifiable at the reduced parameter vector  $a_0$ .

To analyze the linear dependency between the columns of the sensitivity matrix using the multiple step function as input signal, the inner products of the column-vectors of the sensitivity matrix, which are first normalized to length 1, are calculated. Table 3.3 shows that there are almost linear dependencies between these columns, because some inner products are almost equal to  $\pm 1$ . This means that, for example, the variation in the output due to a variation in  $K_s$  can be almost completely compensated by a variation in  $(\theta_s - \theta_r)$  or by  $l$ . For this example the singular values analysis shows that only three singular values ( $\sigma_1 = 27.18$ ,  $\sigma_2 = 12.66$ ,  $\sigma_3 = 2.34$ ) are substantial larger than  $\|\Delta\tilde{X}(a_0)\|$ . The other two singular values ( $\sigma_4 = 0.95$ ,  $\sigma_5 = 0.30$ ) are only one order of magnitude larger than  $\|\Delta\tilde{X}(a_0)\|$ . This means that for some parameters the marginal confidence intervals are very large. These large confidence intervals emphasize the importance of checking the approximation of the Hessian matrix by  $H(\hat{a}) = 2X(\hat{a})'X(\hat{a})$ . If, in the validation phase, it turns out that the approximation of the Hessian matrix is not valid, one has to repeat the experiment with a smaller sampling time so that more samples are obtained. The approximation of the Hessian matrix will then probably be

Table 3.3. Inner products between the column vectors, scaled to length 1, of the sensitivity matrix.

	$\alpha$	$n$	$\theta_r - \theta_s$	$K_s$	$l$
$\alpha$	1.00	0.18	0.16	-0.20	0.53
$n$	0.18	1.00	-0.93	0.89	0.70
$\theta_r - \theta_s$	0.16	-0.93	1.00	-0.99	0.91
$K_s$	-0.20	0.89	-0.99	1.00	-0.93
$l$	0.53	0.70	0.91	-0.93	1.00

better, because the first terms on the right-hand side of eq. (3.20) are more likely to be negligible. Another way to overcome this problem is to use higher order approximations of the confidence region, as shown in Bates and Watts (1988). Although the calculation of higher approximations is cumbersome, it might be helpful to find an accurate approximation of the confidence region. To estimate all six model parameters additional measurements such as the pressure head in the soil sample are necessary, because with only outflow measurements the model is globally not identifiable, regardless of the design of the input signal.

### 3.6 SUMMARY AND DISCUSSION

It is reported by other authors (e.g., Kool et al., 1985) that the identification of transient one-dimensional water flow in unsaturated soils, using outflow measurements on a soil core, will give parameter estimates that are not unique. A model that is often used to describe this process is based on the Richards equation and the Mualem-van Genuchten relations. We have analyzed this model and have compared several input designs. Because the



$\epsilon$ -rank of the sensitivity matrix is smaller than the number of parameters to be estimated, we conclude that this model is numerically not identifiable. The sufficient conditions for a unique minimum in the parameter estimation problem are not fulfilled. Numerical errors can dominate the sum of squares and there may well be many local minima. One may be misled because the Hessian matrix in this minimum can be positive definite. In such a situation the location of the minimum that is found by an optimization program, such as the Levenberg-Marquard algorithm, strongly depends on the initial guess. The location of such a minimum is, however, meaningless.

The multiple step function, proposed by van Dam et al. (1990), is the most promising input signal because the parameters  $\alpha$ ,  $n$ ,  $(\theta_s - \theta_r)$ ,  $K_s$  and  $l$  are numerically identifiable at an initial guess. Since some singular values of the sensitivity matrix are rather small, the assumption that the Hessian matrix may be approximated by  $H(\hat{a}) = 2X(\hat{a})'X(\hat{a})$  has to be validated. However, to estimate all model parameters it is clear that other output signals are necessary.

## REFERENCES

- Bates, D.M. and D.G. Watts, *Non-linear regression analysis and its applications*, Wiley, New York, 1988.
- Bellman, R. and K.J. Åström, On structural identifiability, *Math. Biosci.*, 7, 329-339, 1970.
- Caracotsios, M and W.E. Stewart, Sensivity analysis of initial value problems with mixed ODES and algebraic equations, *Comp. & Chem. Engin.*, 9(4), 359-365, 1985.

- DiStefano, J.J. and C. Cobelli, On parameter and structural identifiability: nonunique observability/reconstructibility for identifiable systems, other ambiguities, and new definitions, *IEEE Trans. Autom. Control.*, 25, 830-833, 1980.
- Dongarra, J.J., J.R. Bunch, C.B. Moler and G.W. Steward, *LINPACK User's Guide*, SIAM, Philadelphia, 1979.
- Hornung, U., Identification of non-linear soil physical parameters from an input-output experiment, in: *P. Deuflhard and E. Hairer (ed.), Workshop on numerical treatments of inverse problems in differential and integral equations*, Birkhäuser, Boston, 227-237, 1983
- Kleissen, F.M., M.B. Beck and H.S. Wheeler, The identifiability of conceptual hydrochemical models, *Water Resour. Res.*, 26, 2979-2992, 1990.
- Kool, J.B., J.C. Parker and M.T. van Genuchten, Determining soil hydraulic properties from one-step outflow measurements by parameter estimation, I, theory and numerical studies, *Soil Sci. Soc. Am. J.*, 49, 1348-1354, 1985.
- Kool, J.B. and J.C. Parker, Analysis of the inverse problem for transient unsaturated flow, *Water Resour. Res.*, 24, 513-522, 1988.
- Ljung, L., *System identification: theory for the user*, Prentice-Hall, Englewood Cliffs, 1987.
- Mous, S.L.J., Numerical solutions of the one-dimensional transient flow equation. Technical note 90-05, Dep. of Mathematics, Agric. Univ. Wageningen, 1990.
- Mualem, Y., A new model for predicting the hydraulic conductivity of unsaturated porous media, *Water Resour. Res.*, 12 513-522, 1976.
- Nguyen, V.V. and E.F. Wood, Review and unification of linear indentifiability concepts, *SIAM Rev.*, 24, 34-51, 1982.

- Sorooshian, S. and V.K. Gupta, The analysis of structural identifiability: Theory and application to conceptual rainfall-runoff models, *Water Resour. Res.*, 21, 478-495, 1985.
- Van Dam, J.C, J.N.M. Stricker and P. Droogers, From One-Step to Multi-Step, determination of soil hydraulic functions by outflow experiments, Rep. 7, Dep. of Hydrology, Soil Physics and Hydraulics, Agric. Univ. Wageningen, 1990.
- Walter E., *Identifiability of state space models*, Springer-Verlag, New York, 1982.

# Chapter 4

## TWO METHODS FOR ASSESSING THE SIZE OF EXTERNAL PERTURBATIONS IN CHAOTIC PROCESSES

### Abstract

This paper deals with the assessment of an external perturbation in nonlinear chaotic dynamical processes using either a modified sentinel function or an extended Kalman filter treatment. We consider processes that can be modeled by a system of nonlinear ordinary differential equations. The sentinel function is used to detect an external perturbation that is not included in the model of the process. In cases where the time dependency of the external perturbation is known but the size of the perturbation is unknown, the sentinel function is also used to estimate the size of this perturbation. We

have compared the sentinel method with an extended Kalman filter treatment. As an example to illustrate these two approaches we have analyzed a low-order spectral model of the atmospheric circulation with a perturbed equator-pole temperature gradient.

#### 4.1 INTRODUCTION

The irregular behavior of many dynamical processes can be modeled by systems of nonlinear differential equations that manifest a chaotic behavior. In this paper we analyze the difference between the dynamical process as it is observed and the solution of a differential equation for this process. Because this type of dynamical model has a sensitive dependence on the initial state, the observed values and the model values will diverge due to the limited accuracy of the estimated initial values. This divergence may get worse in case the physical process is not correctly modeled or certain external perturbations are neglected. The purpose of this study is to analyze methods to detect external perturbations as well as methods to estimate the size of such a perturbation.

The classical least squares method cannot be applied to estimate the size of an external perturbation. The reason for this is the exponential growth of an error in the initial values, making accurate predictions of the output of the model impossible. Therefore, we have to use other methods to analyze the perturbed chaotic dynamical process. Such methods must account for the propagation of the error in the initial state (Baake et. al., 1992). In this paper we compare two approaches. The first approach uses the sentinel function, introduced by Lions (1988, 1990), and the second approach uses the extended Kalman filter (e.g., Jazwinski, 1970, Brammer and Siffling, 1975).

We will first use the sentinel function to detect an external perturbation (Grasman and Houtekamer, 1992). The sentinel function is a weighted

average of the state vector of the process with the weighting chosen in such a way that the effect of an initial error upon the sentinel function is reduced. The external perturbation is then detected by comparing the sentinel as function of the observation with the sentinel as function of the model values. The method of Lions (1988) for constructing sentinel functions was introduced to analyze models of distributed processes. These models are based on a system of partial differential equations. In this study we analyze models that are based on a system of ordinary differential equations and therefore we have modified Lions method in order to make it applicable to this type of systems. In the next section a description of this modified sentinel method is given. We will elaborate as an example the Rössler attractor to illustrate the method.

In section 3 we consider the case where we have more information about the external perturbation and we assume that the time dependency of the perturbation is known. The sentinel function is then used to estimate the size of the perturbation. We have applied this method to a low-order spectral model of the atmospheric circulation with a perturbed equator-pole temperature gradient. The idea is that one detects a change in the heat flux (e.g., greenhouse effect) from the observation of systematic deviations in the circulation.

As an alternative method to estimate the size of an external perturbation we describe in section 4 the extended Kalman filter. This method is frequently used to estimate states and/or parameters of a stochastic dynamical system, so we have to reinterpret our chaotic dynamical system as a stochastic process. To estimate the size of the external perturbation we have regarded the size of the perturbation as a random variable and we have augmented the state vector with this variable as described in Jazwinski (1970).

In the final section we compare the results of these two approaches and make some final remarks.

## 4.2 SENTINELS FOR DETECTING PERTURBATIONS

We consider an  $n$ -dimensional system of nonlinear ordinary differential equations of the form:

$$\frac{dx}{dt} = f(x) \quad (4.1)$$

with initial values

$$x(t_0) = x_0. \quad (4.2)$$

In this study this system constitutes an approximating model of a physical process with a strange attractor. Inaccurate knowledge of parameters or an external perturbation is modeled in the equations by an additive perturbation term. Although this perturbation term is often small, it can still be relevant for the dynamics of the chaotic process. Including the external perturbation and the inaccuracy in the initial state, we have in reality for the state vector a perturbed dynamical system of the form:

$$\frac{dz}{dt} = f(z) + \lambda g(t), \quad t_0 < t < t_1, \quad (4.3)$$

$$z(t_0) = x_0 + \tau \xi_0 \quad (4.4)$$

in which we normalize the vector function  $g(t)$  and the vector  $\xi_0$  according to

$$\frac{1}{t_1 - t_0} \int_{t_0}^{t_1} (g(t), g(t)) dt = 1, \quad (\xi_0, \xi_0) = 1,$$

with  $(\cdot, \cdot)$  the inner product in  $\mathbb{R}^n$ . In case  $g(t)$  or  $\xi_0$  is a stochastic variable another normalization is chosen. In eqs. (4.3-4.4) the parameter  $\lambda$  represents the size of the (unknown) perturbation and the parameter  $\tau$  the inaccuracy in the initial state. We assume that all states can be observed and that the observation errors can be neglected. The observations of the system are then given by the solution of eqs.(4.3-4.4).

Due to sensitive dependence on the initial state the solution of eqs. (4.3-4.4) without external perturbation ( $\lambda=0$ ) and the one of the model system (4.1-4.2) will diverge. Thus, at forehand it is not clear that the discrepancy between the observations  $z(t)$  and the prediction  $x(t)$  stems from shortcomings of the model or from the inaccuracy in the initial values being amplified by sensitive dependence. The aim of the sentinel method is to reduce the latter error so that an assessment of the quality of the model can be made.

In order to keep the notations simple we will write the solution of eqs. (4.3-4.4) as  $z(t) = x(t; \lambda, \tau)$ . To compare the solution of the model eqs. (4.1-4.2), given now by  $x(t; 0, 0)$ , with the observations we introduce the average

$$A(\lambda, \tau) = \int_{t_0}^{t_1} (h(t), x(t; \lambda, \tau)) dt, \quad (4.5)$$

where the vector function  $h(t)$ , satisfying

$$\int_{t_0}^{t_1} |h(t)| dt = 1,$$

is a weight function. It is noted that for  $\lambda$  and  $\tau$  small the difference between the observation average and the model average is approximated by



$$A(\lambda, \tau) - A(0,0) \approx A_\lambda(0,0)\lambda + A_\tau(0,0)\tau. \quad (4.6)$$

Due to sensitive dependence on the initial state the coefficient  $A_\tau(0,0)$  may be large. In that case the second term dominates the difference. Therefore, the average cannot be used to detect a perturbation. A reduction of the contribution coming from the initial error can be achieved by taking the generalized average or sentinel function, introduced first by Lions (1988),

$$S(\lambda, \tau) = \int_{t_0}^{t_1} (h(t) + w(t), x(t; \lambda, \tau)) dt \quad (4.7)$$

with  $w(t)$  a vector function such that

$$S_\tau(\lambda, \tau) = \int_{t_0}^{t_1} (h(t) + w(t), x_\tau(t; 0, 0)) dt = 0. \quad (4.8)$$

Moreover,  $w(t)$  should be as small as possible so that the weight function  $h(t)$  is the least affected. Thus,  $w(t)$  must be such that

$$I(w) = \frac{1}{2} \int_{t_0}^{t_1} (w(t), w(t)) dt \text{ is minimal.} \quad (4.9)$$

Before we analyze this optimal control problem, we first study for the system (4.1-4.2) the initial value problem for the first variation  $x_\tau(t; 0, 0)$ :

$$\frac{d}{dt}x_\tau = f'(t)x_\tau, \quad x_\tau(0) = \xi_{\tau 0}, \quad (4.10)$$

with

$$f'(t) = \left[ \frac{\partial f_i}{\partial x_j}(x(t)) \right]_{n \times n}.$$

The solution of this local tangent linear equation, being a system of linear differential equations, is given by

$$x_\tau = R(t, t_0) \xi_{\tau_0}, \quad (4.11)$$

where  $R(t, t_0)$  is the state transition matrix. It is remarked that the initial error  $\xi_{\tau_0}$  is unknown. Therefore,  $S_\tau$  cannot be computed directly using eq. (4.8). The minimization problem (4.9), with constraint (4.8), can be solved by switching to the nonhomogeneous adjoint local tangent linear equation

$$-\frac{dq}{dt} = [f'(t)]^* q + h(t) + w(t), \quad q(t_1) = 0. \quad (4.12)$$

The solution of this initial value problem is given by

$$q(t) = \int_t^{t_1} P(t, s) (h(s) + w(s)) ds, \quad (4.13)$$

where the state transition matrix  $P(t, s)$  of the adjoint local tangent linear equation is the formal adjoint of  $R(s, t)$ . The following relation holds:  $P(t, s) = R^*(s, t)$  (Talagrand and Courtier, 1987). We then have

$$\begin{aligned} S_\tau(0, 0) &= \int_{t_0}^{t_1} ((h(t) + w(t)), R(t, t_0) \xi_{\tau_0}) dt = \int_{t_0}^{t_1} (R^*(t, t_0) (h(t) + w(t)), \xi_{\tau_0}) dt \\ &= \int_{t_0}^{t_1} (P(t_0, t) (h(t) + w(t)), \xi_{\tau_0}) dt = (q(t_0), \xi_{\tau_0}). \end{aligned} \quad (4.14)$$

Thus in order to reduce the effect of the initial error in the generalized

average  $S(\lambda, \tau)$ , we must choose  $w(t)$  such that  $q(t_0) = 0$ . This optimal control problem for  $w(t)$  satisfying the minimization problem (4.9) with constraint  $q(t_0) = 0$  is solved as follows. We introduce the Lagrange multipliers  $\mu = (\mu_1, \dots, \mu_n)'$  and minimize

$$J(w) = \frac{1}{2} \int_{t_0}^{t_1} (w(t), w(t)) dt + (\mu, q(t_0)). \quad (4.15)$$

Using eq. (4.13) we have

$$\begin{aligned} J(w) &= \int_{t_0}^{t_1} \frac{1}{2} (w(t), w(t)) + (\mu, P(t_0, t)(h(t) + w(t))) dt = \\ &= \int_{t_0}^{t_1} \frac{1}{2} (w(t), w(t)) + (P^*(t_0, t)\mu, (h(t) + w(t))) dt, \end{aligned} \quad (4.16)$$

so that  $\partial J / \partial w = 0$  for

$$w(t) = -P^*(t_0, t)\mu = -R(t, t_0)\mu. \quad (4.17)$$

Rewriting the constraint (4.6) gives the condition

$$q(t_0) = \int_{t_0}^{t_1} P(t_0, t)(h(t) - R(t, t_0)\mu) dt = 0. \quad (4.18)$$

Consequently, the Lagrange multipliers satisfy

$$\mu = \left[ \int_{t_0}^{t_1} P(t_0, t)R(t, t_0) dt \right]^{-1} \int_{t_0}^{t_1} P(t_0, t)h(t) dt. \quad (4.19)$$

Assuming that  $S_\lambda(0,0) \neq 0$ , the difference between  $S(\lambda,\tau)$  and  $S(0,0)$  yields an order estimate of the average error in  $x(t)$  on the interval  $(t_0, t_1)$  caused by the external perturbation:

$$S(\lambda,\tau) - S(0,0) \approx \lambda S_\lambda(0,0) \quad (4.20)$$

with

$$S_\lambda(0,0) = \int_{t_0}^{t_1} ((h(t)+w(t)), x_\lambda(t;0,0)) dt \quad (4.21)$$

and

$$x_\lambda(t;0,0) = \int_{t_0}^t R(t,s)g(s) ds. \quad (4.22)$$

The estimate of the perturbation strongly depends on the type of forcing function  $g(t)$ . An oscillating function tends to produce a vanishing contribution to the integral. These are "stealthy" perturbations (Lions, 1990). Therefore the estimate gives a lower bound for the average error in  $x(t)$  on the interval  $(t_0, t_1)$  due to the perturbation  $\lambda g(t)$ .

#### 4.2.1 THE RÖSSLER ATTRACTOR

We apply the sentinel method to a third order system of ordinary differential equations, first formulated by Rössler (1976):

$$\frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -(x_1 + x_3) \\ x_1 + ax_2 \\ b + x_1x_3 - cx_3 \end{bmatrix}. \quad (4.23)$$

We take the model values  $a=0.2$ ,  $b=0.2$  and  $c=5.7$ , so that the system has a strange attractor as solution. We assume that for the "true" process the parameter  $b$  is perturbed by some external force,

$$b_r = b + \lambda g_3(t).$$

To analyze the potential use of the sentinel function for detecting an external perturbation, we test the method by simulating the "true" process with several types of test functions  $g(t)$  and with the initial conditions of the state vector contaminated with white noise with variance  $\tau^2=10^{-4}$ . We use three test-functions: a fixed perturbation:  $g(t)=(0,0,1)'$ , a standard white Gaussian process with zero mean:  $g(t)=(0,0,\xi(t))'$  and a white Gaussian process with mean equal to one:  $g(t)=(0,0,0.5(\xi(t)+1))'$ . We have taken the size of these perturbations as  $\lambda=0.01$  and  $\lambda=0.05$ . The processes were simulated over a time period  $T=t_1-t_0=1$  and  $T=5$ . For the sentinel function we have taken a uniform weighing  $h(t)$ . The results of these tests are shown in tables 4.1 and 4.2. Here, the average error in the state vector is compared with the sentinel function. It is seen that for the white Gaussian processes the sentinel function is smaller than for the constant perturbation. Because the sentinel function for these oscillating processes is an order of magnitude larger than  $S(0,\tau)$  these processes are not completely "stealthy", however, they are hard to detect. It is also seen that the result does not always improve by calculating the sentinel with observations taken over a longer time period. The reason for this is that the error in the state vector caused by the perturbation  $x_\lambda(t)$  oscillates, causing a cancellation in the sentinel function.

The average error in the state vector caused by  $\lambda g(t)$  is shown in the

last columns of tables 4.1 and 4.2. In most cases this error is of the same order as the difference  $S(\lambda, \tau) - S(0,0)$ , so we indeed can use this difference to estimate the average error in the state vector from external perturbations.

Table 4.1. Detection of external perturbations using the sentinel method. The white Gaussian process is represented by  $\xi(t)$ . In the last column we have

$$D(\lambda g(t)) = \int_{t_0}^{t_1} (h(t), |x(t; \lambda, 0) - x(t; 0, 0)|) dt,$$

denoting the average error in  $x(t)$  during the time period  $T$ . We have taken a time period  $T=1.0$ .

$\lambda$	$g(t)$	$S(\lambda, \tau) - S(0,0)$	$D(\lambda^* g(t))$
0.0		$2.76 \cdot 10^{-6}$	0
0.01	1	$1.12 \cdot 10^{-3}$	$2.66 \cdot 10^{-3}$
0.05	1	$5.59 \cdot 10^{-3}$	$1.33 \cdot 10^{-2}$
0.01	$\xi(t)$	$4.17 \cdot 10^{-5}$	$1.37 \cdot 10^{-4}$
0.05	$\xi(t)$	$1.98 \cdot 10^{-4}$	$6.85 \cdot 10^{-4}$
0.01	$0.5(1+\xi(t))$	$5.02 \cdot 10^{-4}$	$1.38 \cdot 10^{-3}$
0.05	$0.5(1+\xi(t))$	$2.90 \cdot 10^{-3}$	$6.92 \cdot 10^{-3}$

---

Table 4.2. See table 4.1 for the explanation. Now for the time period is taken  $T=5.0$ .

$\lambda$	$g(t)$	$S(\lambda, \tau) - S(0,0)$	$D(\lambda * g(t))$
0.0		$9.32 \cdot 10^{-7}$	0
0.01	1	$4.32 \cdot 10^{-4}$	$6.17 \cdot 10^{-3}$
0.05	1	$2.15 \cdot 10^{-3}$	$3.08 \cdot 10^{-2}$
0.01	$\xi(t)$	$1.37 \cdot 10^{-5}$	$1.01 \cdot 10^{-4}$
0.05	$\xi(t)$	$6.48 \cdot 10^{-5}$	$5.05 \cdot 10^{-4}$
0.01	$0.5(1+\xi(t))$	$2.29 \cdot 10^{-4}$	$3.11 \cdot 10^{-3}$
0.05	$0.5(1+\xi(t))$	$1.14 \cdot 10^{-3}$	$1.56 \cdot 10^{-2}$

### 4.3 SENTINELS FOR ESTIMATING UNCERTAIN PARAMETERS

In the foregoing it was assumed that the size  $\lambda$  and the time dependency  $g(t)$  of the perturbation were unknown. The sentinel method can then be used to estimate the average error in the state vector, caused by this perturbation. Now we suppose that  $g(t)$  is a known function. This extra information is used to estimate the size  $\lambda$  of the perturbation. Anew we consider a process that is described by a system of ordinary differential equations given by eqs. (4.1-4.2). The formal analysis of the preceding section again applies. The only difference is that  $S_\lambda(0,0)$  can be computed by using eqs. (4.21) and (4.22). Assuming that  $S_\lambda(0,0) \neq 0$ , we can estimate the size  $\lambda$  of the perturbation using the Taylor expansion (4.20). We then arrive at the estimator:

$$\hat{\lambda} = \frac{S(\lambda, \tau) - S(0,0)}{S_\lambda(0,0)}. \quad (4.24)$$

For most time intervals this estimate is sufficiently accurate, however, for certain time intervals  $S_\lambda(0,0)$  can be small. Then higher order terms in the

Taylor expansion are not always negligible with respect to  $S_\lambda(0,0)$  and produce inaccurate and biased estimates.

In section 4.2 we neglected the errors in the observations. We will now include them in our analysis. These observation errors appear not to be negligible in cases where the difference  $S(\lambda, \tau) - S(0,0)$  is small, because then the contribution of the observation errors to  $S(\lambda, \tau)$  will be of the same order as the difference  $S(\lambda, \tau) - S(0,0)$ . This also yields inaccurate estimates.

#### 4.3.1 A LOW-ORDER SPECTRAL MODEL OF THE ATMOSPHERIC CIRCULATION WITH A PERTURBED EQUATOR-POLE TEMPERATURE GRADIENT

We study the use of the sentinel function for a 10-component spectral model of the atmospheric circulation that has been analyzed by De Swart (1988a, 1988b). For barotropic flow the atmospheric circulation is described by a streamfunction. This streamfunction,  $\psi(x,y,t)$ , satisfies the so-called quasi-geostrophic barotropic potential vorticity equation:

$$\frac{\partial}{\partial t} \nabla^2 \psi + J(\psi, \nabla^2 \psi + f) + \gamma J(\psi, m) + C \nabla^2 (\psi - \psi^*) = 0, \quad (4.25)$$

where  $J$  is the Jacobian operator given by

$$J(a,b) = a_2 b_1 - a_1 b_2, \quad a = (a_1, a_2)', \quad b = (b_1, b_2)'. \quad (4.26)$$

The Coriolis parameter  $f$  is taken to be fixed, meaning that the flow is restricted to a channel in the tangent plane at a given latitude. The term  $m(x,y)$  describes the topography (mountains) of the domain and the coefficient  $\gamma$  accounts for the effect of this topography. The coefficient  $C$  is a



measure for frictional effects and finally the term  $\psi^*(x,y,t)$  models the driving force of the atmospheric flow: the equator-pole temperature gradient.

The solution  $\psi(x,y,t)$  of eq. (4.25) is approximated by expanding  $\psi$ ,  $\psi^*$  and  $m$  in a series of eigenfunctions  $\{\phi_j\}$  of the Laplace operator for the channel being a domain in the  $x,y$ -plane, periodic in  $x$  and bounded by lower and upper values of  $y$ :

$$\psi(x,y,t) = \sum_{j=1}^{\infty} x_j(t) \phi_j(x,y). \quad (4.27)$$

In De Swart (1988a) the system of nonlinear differential equations for the coefficients  $x_j(t), j=1, \dots, 10$  of the truncated series  $\psi(x,y,t)$  is derived. It is of the form

$$\frac{dx}{dt} = f(x) + Cx^*, \quad (4.28)$$

(see appendix A). This spectral model has a vacillating solution and is a system of the lowest possible order that still exhibits this behavior (De Swart, 1988a). Vacillation means that the solution visits in an irregular manner domains in state space where it remains for some time. These domains correspond to so-called preferent weather regimes.

The effect of the equator-pole temperature gradient is represented in eq. (4.28) by the term  $Cx^*$ , with  $x^*=(x_1^*, 0, 0, x_4^*, 0, \dots, 0)$ . We want to study the influence of a perturbation in this temperature gradient on the evolution of this spectral model. Since this system has a sensitive dependence on the initial state, the sentinel method is used.

As in most situations in practise the observations are taken at discrete times  $t_k, k=1, \dots, N$ . Therefore we use the discrete analogon of the sentinel function and we define the discrete sentinel function as

$$S(\lambda, \tau) = \sum_{k=1}^N (h_k + w_k x(t_k)) \quad (4.29)$$

with  $h_k, w_k \in \mathbb{R}^{10}$  representing the weighing. Following the same computations as in the previous section we derive for the weighing  $w_k$

$$w_k = -R_k \mu \quad (4.30)$$

with

$$\mu = \left( \sum_{k=1}^N R_k^T R_k \right)^{-1} \sum_{k=1}^N R_k^T h_k, \quad R_k = R(t_k, t_0). \quad (4.31)$$

In this test example the "true" process is simulated by integrating eq. (4.28), with the initial conditions contaminated with white noise with variance  $\tau^2 = 10^{-6}$  and with the vector  $x^*$  perturbed with  $\lambda g(t)$ . We have taken  $\lambda=0.05$  and  $g(t)=(1,0,\dots,0)'$ . The observations are sampled with a sampling period  $T_s = 0.25$ . We have also contaminated the observations with white observation noise, with variance  $\tau^2 = 10^{-6}$ .

The sentinel function  $S(\lambda, \tau)$  is calculated using a batch of  $N=12$  consecutive observations from the data collection of the "true" process, and with a uniform weighing  $h_k$ . The model data for a batch are obtained by integrating eq. (4.28) with the initial conditions equal to the first observations of the batch of the "true" process and with  $\lambda=0$ . These data are then used to calculate  $S(0,0)$ . We also calculate  $S_\lambda(0,0)$  for this batch and estimate  $\lambda$  according to eq. (4.24). Because the data set is much larger than a single batch, we can monitor  $\hat{\lambda}$  in a "moving sentinel" approach. By this we mean that the first batch of  $N$  samples is used to calculate  $\hat{\lambda}(1)$ . Shifting one place through the data set we get a second batch of data, from which  $\hat{\lambda}(2)$  is

obtained, and so on.

In figure 4.1(a) we show  $\hat{\lambda}$  as function of the counting index  $k$  of the batches. The large outliers that appear in this picture are caused by small values of  $S_{\lambda}$ , as we see from figure 4.1(b); the outliers occur when the corresponding value of  $|S_{\lambda}(k)|$  gets small. This makes these estimates unreliable for reasons discussed before. To improve the result we take the mean of all estimates  $\hat{\lambda}(k)$ , except those estimates  $\hat{\lambda}(k)$  for which the absolute value of the accessory  $S_{\lambda}$  is smaller than a threshold value. As threshold value we choose the mean of the absolute values of  $S_{\lambda}$ . The results of this procedure for various values of  $\lambda$  and  $\tau$  are presented in table 4.3. It is seen that the size of  $\lambda$  with respect to  $\tau$  determines the quality of the estimates. This procedure yields less accurate estimates only for small values of  $\lambda$ . In all other cases the estimate  $\hat{\lambda}$  is sufficiently accurate.

Table 4.3. Estimates of  $\hat{\lambda}$  using the sentinel method, as function of the "true" parameters  $\lambda$  and  $\tau$  that are used to simulate the "true" process.

$\tau$	0.0	0.005	0.001	0.0015	0.002
$\lambda$					
0.0	$6.9 \cdot 10^{-7}$	-0.0039	-0.0037	-0.0059	-0.0166
0.01	0.0099	0.0076	0.0034	$-2.9 \cdot 10^{-4}$	0.0058
0.02	0.020	0.018	0.019	0.013	0.030
0.03	0.030	0.029	0.027	0.025	0.046
0.04	0.039	0.037	0.036	0.046	0.046
0.05	0.054	0.051	0.046	0.047	0.050

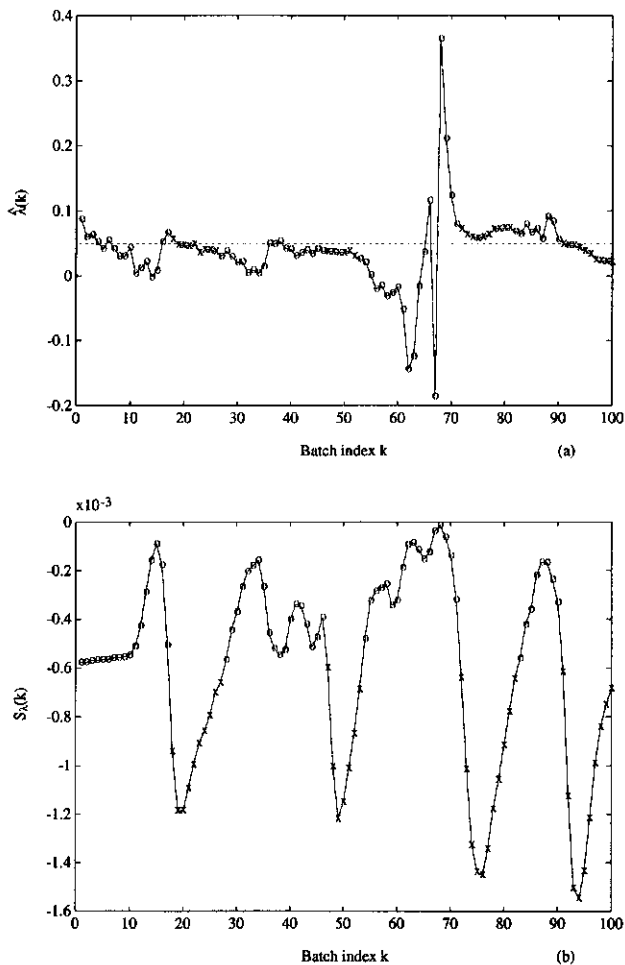


Figure 4.1. The sentinel method for the atmospheric circulation problem. The "true" process was simulated with  $\lambda=0.05$  and  $\tau^2=10^{-6}$ . Each batch consists of 12 consecutive observations; batch  $k$  starts with observation  $k$ . (a) Estimated parameters  $\hat{\lambda}(k)$  and (b) first derivatives  $S_{\lambda}(k)$  as function of the counting index  $k$  of the batches; cases where  $|S_{\lambda}|$  is below the threshold value are indicated by "o".

#### 4.4 EXTENDED KALMAN FILTERING FOR ESTIMATING UNCERTAIN PARAMETERS

A more traditional method to estimate parameters and states of a dynamical system is the Kalman filtering technique. This filter is developed in the early sixties by Kalman and Bucy. It gives, for linear dynamical systems, the optimal estimator for the state of the dynamical system, in the sense of minimum variance. For nonlinear dynamical systems there are many approximated optimal filters developed (e.g., Jazwinski, 1970, Sorenson, 1988). In this study we will use the extended Kalman filter and we adapt a procedure described in Jazwinski (1970, pp. 281-282) to estimate an uncertain parameter. First we will give a description of the extended Kalman filter. Consider a nonlinear dynamical system of the form

$$\frac{dx}{dt} = f(x) + \xi_1(t), \quad x(0) = x_0 + \xi_0, \quad t > 0, \quad (4.32)$$

$$y(t_k) = Mx(t_k) + \xi_2(t_k), \quad (4.33)$$

where  $x(t)$  is the state vector and  $y(t_k)$  are the observations taken at time  $t_k$ . We assume that the noise term  $\xi_2(t_k)$ , caused by observations errors, can be modeled by a white Gaussian process in discrete time. Errors in the state equation are modeled by a continuous white Gaussian process  $\xi_1(t)$ . The autocovariance functions of these processes are of the form

$$E(\xi_1(t)\xi_1(t+\tau)') = Q_1(t)\delta_D(\tau) \quad \text{and} \quad E(\xi_2(t_k)\xi_2(t_l)') = R(t_k)\delta_K(t_k - t_l),$$

where  $\delta_D(\cdot)$  denotes the Dirac delta-function and  $\delta_K(\cdot)$  denotes the Kronecker delta. Furthermore we assume that the error in the initial state  $\xi_0$  and these noise processes,  $\xi_1(t)$  and  $\xi_2(t_k)$ , are mutually uncorrelated.

Starting from estimates of the mean  $\hat{x}(t_k|t_k)$  and the error covariance matrix  $P(t_k|t_k)$  of the state vector at time  $t_k$ , given the observations  $y(t_1), \dots, y(t_k)$ , we predict the state  $\hat{x}(t_{k+1}|t_k)$  and the error covariance matrix  $P(t_{k+1}|t_k)$  at time  $t_{k+1}$  with the first two equations of the extended Kalman filter

$$\hat{x}(t_{k+1}|t_k) = \hat{x}(t_k|t_k) + \int_{t_k}^{t_{k+1}} f(\hat{x}(t|t_k)) dt, \quad (4.34)$$

$$P(t_{k+1}|t_k) = \Phi(t_{k+1}, t_k; \hat{x}(t_k|t_k)) P(t_k|t_k) \Phi^T(t_{k+1}, t_k; \hat{x}(t_k|t_k)) + Q(t_{k+1}), \quad (4.34)$$

where  $\Phi(t_{k+1}, t_k; \hat{x}(t_k|t_k))$  is the state transition matrix of the linearized state equation

$$\frac{d(\delta x)}{dt} = f'(\hat{x}(t_k|t_k)) \delta x \quad (4.36)$$

with

$$f'(\hat{x}(t_k|t_k)) = \left[ \frac{\partial f(\hat{x}(t_k|t_k))}{\partial x} \right]_{n \times n}$$

and where

$$Q(t_{k+1}) = \int_{t_k}^{t_{k+1}} \Phi(t_{k+1}, t_k) Q_1(t) \Phi^T(t_{k+1}, t_k) dt.$$

At time  $t_{k+1}$  a new observation comes available. This observation is used with the remaining three equations of the Kalman filter to update the estimates of the mean and the error covariance matrix of the state vector,

$$\hat{x}(t_{k+1}|t_{k+1}) = \hat{x}(t_{k+1}|t_k) + K(t_{k+1})(y(t_{k+1}) - M\hat{x}(t_{k+1}|t_k)), \quad (4.37)$$

$$P(t_{k+1}|t_{k+1}) = (I - K(t_{k+1})M)P(t_{k+1}|t_k), \quad (4.38)$$

$$K(t_{k+1}) = P(t_{k+1}|t_k)M^T(MP(t_{k+1}|t_k)M^T + R(t_k))^{-1}. \quad (4.39)$$

The filter is initialized with estimates  $\hat{x}(t_0|t_0)$  and  $P(t_0|t_0)$ . These estimates of the initial state are assumed to be known.

We will use this filter to estimate the state and the size of the external perturbation of a dynamical system given by eqs. (4.3-4.4). As in the previous section we assume that the shape  $g(t)$  of the external perturbation is known and that the size  $\lambda$  is unknown. The size of the external perturbation can be seen as an uncertain parameter in this model. Therefore, we augment the model (4.3-4.4) with an additional state variable by regarding this parameter as a random variable with mean and variance

$$E(\lambda) = \lambda_r, \quad \text{var}(\lambda) = 0,$$

meaning that this variable is a constant and thus its variance is equal to zero. It satisfies

$$\frac{d\lambda}{dt} = 0 \quad (4.40)$$

with an initial mean  $\hat{\lambda}(t_0|t_0)$  and initial variance  $P_\lambda(t_0|t_0)$ . We now combine eq. (4.3) and eq. (4.40) into the augmented model

$$\frac{dX}{dt} = \frac{d}{dt} \begin{bmatrix} x \\ \lambda \end{bmatrix} = \begin{bmatrix} f(x) \\ 0 \end{bmatrix} + \begin{bmatrix} \lambda g(t) \\ 0 \end{bmatrix} \quad (4.41)$$

with state vector  $X(t) = (x(t), \lambda(t))'$ . Since the state equation of the original model (4.3) and the state equation (4.40) are noise free, the state equation of the augmented model (4.41) is also noise free. Therefore, the autocovariance function  $Q(t)$  can be set equal to 0 in the extended Kalman filter for this augmented model.

We will now apply the extended Kalman filter to the spectral model of the atmospheric circulation, so that we can estimate the size  $\lambda$  of the perturbation in the equator-pole temperature gradient. We then can make a comparison between the extended Kalman filter approach and the sentinel method. We will regard  $\lambda$  as a random variable with mean  $\lambda$ , and variance equal to 0, so that the spectral model (4.28) changes into

$$\frac{dX}{dt} = \frac{d}{dt} \begin{bmatrix} x \\ \lambda \end{bmatrix} = \begin{bmatrix} f(x) \\ 0 \end{bmatrix} + \begin{bmatrix} C(x^* + \lambda) \\ 0 \end{bmatrix}, \quad (4.42)$$

$$y(t_k) = \begin{bmatrix} I & 0 \end{bmatrix} \begin{bmatrix} x(t_k) \\ \lambda(t_k) \end{bmatrix} + \xi_2(t_k), \quad (4.43)$$

where  $y(t_k)$  are the observations of the state vector at time  $t_k$ . For the observations of the "true" process we will use the data sets that also have been used by the sentinel method. We take as the initial estimate for  $\hat{x}(t_0|t_0)$  the first observation of the data set. As the initial estimate for  $P(t_0|t_0)$  we take

$$\begin{aligned} P_{i,j}(t_0|t_0) &= \tau^2 \text{ if } i=j, \\ &= 0 \text{ if } i \neq j, \end{aligned}$$



because the observations are contaminated with white noise with variance  $\tau^2$ . Accordingly the error covariance function of  $\xi_2(t_k)$  is therefore also equal to  $R\delta(t_k-t_l)$  with

$$\begin{aligned} R_{i,j} &= \tau^2 \text{ if } i=j, \\ &= 0 \text{ if } i \neq j. \end{aligned}$$

We set  $\hat{\lambda}(t_0|t_0)=0$  and  $P_\lambda(t_0|t_0)=100$  as the initial statistics for  $\lambda$ . This large variance  $P_\lambda(t_0|t_0)$  stands for the initial uncertainty in the parameter  $\lambda$ .

We will filter two data sets, where we take as "true" parameter values  $\lambda=0.01$ ,  $\tau^2=10^{-6}$  and  $\lambda=0.05$ ,  $\tau^2=10^{-6}$  respectively. We saw in the previous section that the result of the sentinel method was only sufficiently accurate for the second data set, where we found  $\hat{\lambda}=0.046$ . The filtering results are shown in figures 4.2(a) and (b). In these figures we plotted the estimated values  $\hat{\lambda}(t_k|t_k)$ . The dotted line represents the error standard deviation,  $\lambda \pm \sqrt{P_\lambda(t_k|t_k)}$  (root mean square error curve). We have found as final estimates  $\hat{\lambda}(t_{100}|t_{100})=0.010$  for the first data set and  $\hat{\lambda}(t_{100}|t_{100})=0.051$  for the second data set.

Figures 4.2(a) and 4.2(b) show that the error standard deviation decreases fast, which means that the data contain much information on  $\lambda$ . Especially for the data set with  $\lambda=0.01$  we see that the standard deviation decreases very fast for samples taken at the interval  $(t_{30}, t_{40})$ . Apparently, the process is very sensitive to this parameter on this part of the attractor.

It is noted that for the data set with  $\lambda=0.05$ , the estimates  $\hat{\lambda}(t_k|t_k)$ , for large  $k$  fall somewhat outside the rms-error curve. This may be due to the particular realizations of the observation errors and the choice of the initial estimate  $\hat{\lambda}(t_0|t_0)$ . Since the filter is an approximated optimal filter, it may also be that the higher order moments of the probability density functions and numerical errors are not negligible for this perturbation. Sometimes the result can be improved by including a small fictitious error term in the state

equation (Jazwinski, 1970, pp. 301-307).

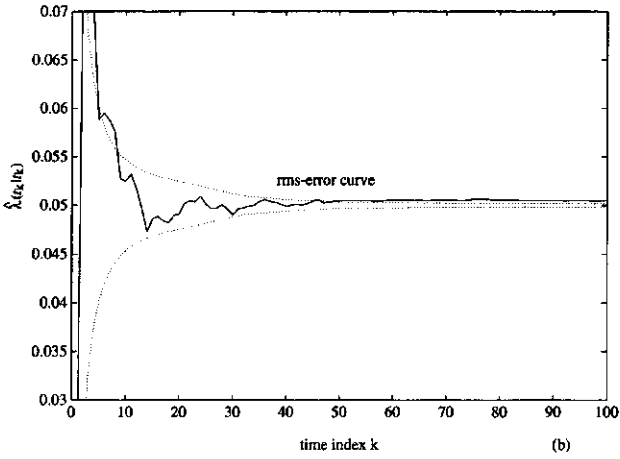
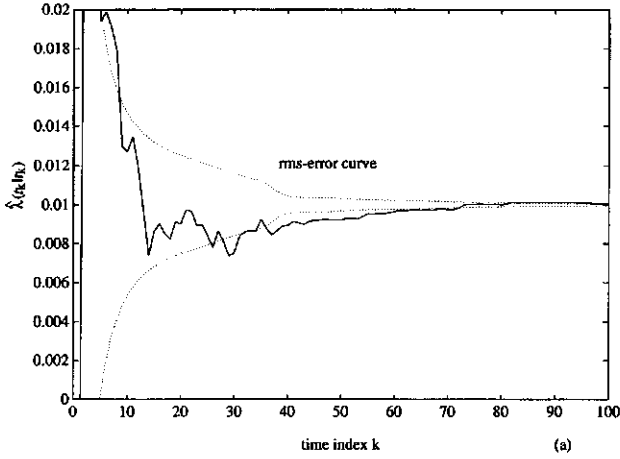


Figure 4.2. The extended Kalman filter for the atmospheric circulation problem. The solid line denotes the estimated  $\hat{\lambda}(t_k|t_k)$  and the dotted line denotes the root mean square error curve,  $\lambda \pm \sqrt{P_{\lambda}(t_k|t_k)}$ . (a) The "real" process was simulated with  $\lambda=0.01$  and  $\tau^2=10^{-6}$ , (b) the "real" process was simulated with  $\lambda=0.05$  and  $\tau^2=10^{-6}$ .

## 4.5 CONCLUSIONS

In this paper we have modified Lions' sentinel function so that it can be applied to systems that are described by ordinary differential equations. We also have presented an analog of the modified sentinel function for situations where the observations of the system are discrete. This modified sentinel function is applied to systems that are perturbed by some external force. Two situations can be distinguished: the external perturbation is completely unknown, and the case that the time dependency of the perturbation is known, but the size of the perturbation is unknown.

For the first situation, we have shown that the sentinel can be interpreted as the averaged error in the state vector caused by the unknown perturbation. It can be used to detect an external perturbation by comparing the sentinel as function of the observations with the sentinel as function of the model values. Because the sentinel is a weighted average of the state vector over a time interval  $(t_0, t_1)$ , oscillating perturbations are hard to detect. Such so-called "stealthy" perturbations can be made visible by an appropriate choice of the time interval and the weighing  $h(t)$  (Lions, 1990).

For the second situation, we regarded the size of the external perturbation as an uncertain parameter in the model. We have constructed an estimator for this uncertain parameter based on the sentinel function. As an alternative we also have constructed an estimator based on the external Kalman filter. These two estimators have been compared in a case study, where we have analyzed a spectral model of the atmospheric circulation with a perturbed equator-pole temperature gradient. In this case study we compared simulations of the perturbed model, representing the "true" process, with simulations of the model. For large values of the uncertain parameter  $\lambda$ , compared with the accuracies in the initial state and the observations, the error in the estimator based on the sentinel function  $\varepsilon = \hat{\lambda} - \lambda$  is of the order of 5-10%. For smaller values of  $\lambda$  the inaccuracy even increases. In the exten-

ded Kalman filter treatment the error in the estimator of the uncertain parameter is considerably smaller,  $\epsilon < 2\%$ , in the two cases that we have analyzed.

The less accurate results of the sentinel method are caused by the relatively low sampling frequency. The derivative  $S_\lambda(0,0)$  can then be small for certain time intervals, so that the effect of neglecting higher order terms of the Taylor expansion is felt. Moreover, the influence of observation errors becomes significant.

From figure 4.2 it is seen that the root mean square error curve calculated by the extended Kalman filter drops considerably after a few observations are taken. This implies that the data contains much information about the size of the perturbation. This observation is also reported in a paper of Baake et al. (1992), where they concluded that an observed trajectory of a chaotic process may be expected to contain a large amount of information about uncertain parameters.

The extended Kalman filter that we have used for estimating an uncertain parameter, is less suitable for the first situation. We then have to augment the state vector with more variables, because it is not known if the system is affected by one or more perturbations and which elements of the state vector are affected by these perturbations. For the special situation where the extended perturbation can be modeled by an additional noise term, we refer to Mehra (1970) and Iglehart and Leondes (1974). They present some methods to estimate the covariance matrix of such a noise term using methods that are based on the Kalman filter.

*Acknowledgement.* We are grateful to Huib de Swart for providing the code of his 10-component spectral model of the atmospheric circulation and to Maarten de Gee and Albert Otten for their advices during the preparation of the manuscript.

## REFERENCES

- Baake, E., M. Baake, H.G. Bock and K.M. Briggs, Fitting ordinary differential equations to chaotic data, *Phys. Rev. A*, vol. 45, 8, 5524-5529, 1992.
- Brammer, K., and G. Siffling, *Kalman-Bucy filter: deterministische beobachtung und stochastische filterung*, Oldenbourg, München, 1975.
- De Swart, H.E., Vacillation and predictability properties of low-order atmospheric spectral models, Ph.D.Thesis, Rijksuniversiteit Utrecht, 1988a.
- De Swart, H.E., Low-order spectral models of the atmospheric circulation: a survey. *Acta Appl. Math.*, 11, 49-96, 1988b.
- Grasman, J., and P. Houtekamer, Methods for improving the prediction of dynamical processes with special reference to the atmospheric circulation, *Proc. of IUGG Symp. "Nonlinear Dynamics and Predictability of critical Geophysical Phenomena*, W.I. Newman and A.M. Gabriellov (eds.), Am. Geophysical Union, 1992.
- Iglehart, S.C., and C.T. Leondes, Estimation of a dispersion parameter in discrete Kalman filtering, *IEEE Trans. Autom. Control*, AC-19, 262-263, 1974.
- Lions, J.L., *Sur les sentinelles de systèmes distribuées*, C.R.A.S., Paris, 1988.
- Lions, J.L., Sentinels and stealthy perturbations, *Proc. Int. Symp. on Assimilation of Observations in Meteorology and Oceanography*, Clermont-Ferrand, World Meteorological Organization, 13-18, 1990.
- Jazwinski, A.H., *Stochastic Processes and Filtering Theory*, Academic Press, Paris, 1970.
- Mehra, R.K., On the identification of variances and adaptive Kalman filtering, *IEEE Trans. Autom. Control*, AC-15, 2, 175-184, 1970.
- Rössler, O.E., An equation for continuous chaos, *Phys. Lett. A*, 57, 397-398, 1976.

- Sorenson, H.W., Recursive estimation for nonlinear dynamical systems, in *Bayesian analysis of time series and dynamic models*, J.C. Spall ed., Dekker, New York, 1988.
- Talagrand, O., and P. Courtier, Variational assimilation of meteorological observations with the adjoint vorticity equation. Part. 1: Theory, *Quart. J. Roy. Meteo. Soc.*, 113, 1311-1328, 1987.

# APPENDIX A.

In De Swart (1988) a 10-component spectral model for the barotropic flow (4.25) in a rectangular channel is derived. The channel has length  $2\pi$  in the zonal direction and width  $\pi b$  in the meridional direction. The 10-component model is given by

$$\frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \\ x_{10} \end{bmatrix} = \begin{bmatrix} +\gamma_{11}^* x_3 & -C(x_1 - x_1^*) & & & & & & & & \\ -(\alpha_{11} x_1 - \beta_{11}) x_3 & -C x_2 & -\delta_{11} x_4 x_6 & -\rho_{11} (x_5 x_8 - x_6 x_7) & & & & & & \\ +(\alpha_{11} x_1 - \beta_{11}) x_2 & -\gamma_{11} x_1 & -C x_3 & +\delta_{11} x_4 x_5 & +\rho_{11} (x_5 x_7 + x_6 x_8) & & & & & \\ +\gamma_{12}^* x_6 & -C(x_4 - x_4^*) & +\epsilon_1 (x_2 x_6 - x_3 x_5) & +\epsilon_2 (x_7 x_{10} - x_8 x_9) & & & & & & \\ -(\alpha_{12} x_1 - \beta_{12}) x_6 & -C x_5 & -\delta_{12} x_3 x_4 & +\rho_{12} (x_2 x_8 - x_3 x_7) & +\gamma_{12}^* x_8 & & & & & \\ +(\alpha_{12} x_1 - \beta_{12}) x_5 & -C x_6 & +\delta_{12} x_2 x_4 & -\rho_{12} (x_2 x_7 + x_3 x_8) & -\gamma_{12}^* x_7 & & & & & \\ -(\alpha_{21} x_1 - \beta_{21}) x_8 & -C x_7 & -\delta_{21} x_4 x_{10} & -\rho_{21} (x_2 x_6 + x_3 x_5) & +\gamma_{21}^* x_6 & & & & & \\ +(\alpha_{21} x_1 - \beta_{21}) x_7 & -C x_8 & +\delta_{21} x_4 x_9 & +\rho_{21} (x_2 x_5 - x_3 x_6) & -\gamma_{21}^* x_5 & & & & & \\ -(\alpha_{22} x_1 - \beta_{22}) x_{10} & -C x_9 & -\delta_{22} x_4 x_8 & & & & & & & \\ +(\alpha_{22} x_1 - \beta_{22}) x_9 & -C x_{10} & +\delta_{22} x_4 x_7 & & & & & & & \end{bmatrix}$$

advection    topo-    forcing/    advection    wave    topo-  
                  graphy    advec-                   triad    graphy  
                                  tion

where

$$\begin{aligned}
\alpha_{nm} &= \frac{8\sqrt{2}n}{\pi} \frac{m^2}{4m^2-1} \frac{n^2b^2+m^2-1}{n^2b^2+m^2}, \quad \beta_{nm} = \frac{\beta nb^2}{n^2b^2+m^2} \\
\delta_{nm} &= \frac{64\sqrt{2}n}{15\pi} \frac{n^2b^2-(m^2-1)}{n^2b^2+m^2}, \quad \gamma_{nm}^* = \frac{4m}{4m^2-1} \frac{\sqrt{2}nb\gamma}{\pi} \\
\varepsilon_n &= \frac{16\sqrt{2}n}{5\pi}, \quad \gamma_{nm} = \frac{4m^3}{4m^2-1} \frac{\sqrt{2}nb\gamma}{\pi(n^2b^2+m^2)} \\
\rho_{nm} &= \frac{9}{2} \frac{(n-2)^2-(m-2)^2}{n^2b^2+m^2}, \quad \gamma'_{nm} = \frac{3b\gamma}{4(n^2b^2+m^2)}
\end{aligned}$$

The  $\beta_{nm}$ -contributions represent planetary vorticity advection, the  $\gamma_{nm}^*$ ,  $\gamma_{nm}$  and  $\gamma'_{nm}$ -terms the various couplings between flow and topography and  $Cx_n^*$  the equator pole-temperature gradient. We have taken  $b=1.6$ ,  $\beta=1.25$   $\gamma=1$ ,  $C=0.1$ ,  $x_1^*=4$  and  $x_4^*=8$ .

# **Chapter 5**

## **ESTIMATING UNCERTAIN MODEL AND NOISE PARAMETERS IN CHAOTIC MODELS WITH APPLICATIONS IN METEOROLOGY**

### **Abstract**

For a well operating extended Kalman filter it is necessary to have an accurate description of the physical process and knowledge of the noise statistics. In this paper, adaptive extended Kalman filter techniques are analyzed for nonlinear chaotic models that involve unknown parameters of two different types. First, there are parameters describing unknown systematic perturbations. Second, the error covariance matrix of the model is unknown; this matrix is assumed to be described in terms of a few (unknown) param-



ters. The values of both sets of parameters are estimated using an approximated Maximum Likelihood method. To this end, the probability distributions of the forecast errors are approximated by Gaussian distributions. The value of the log-likelihood function is then easily obtained using the extended Kalman filter. However, in case of divergence of the extended Kalman filter, the approximate log-likelihood function may have many local maxima and consequently, commonly used optimization algorithms cannot locate the optimal solution. Here, we mean with divergence of the extended Kalman filter that the actual forecast errors are larger than expected from the computed error statistics. In our approach to this problem we add an artificial noise term to the state equations of the model. This artificial noise term is used to control the accuracy of the estimates, so that the filter does not learn the wrong state too well. The values of the unknown model and noise parameters are approximated by an optimization procedure that is designed in such a way that it does not tend to settle in a local optimum caused by divergence of the filter. The above mentioned problems arise in particular in chaotic systems. To illustrate how these problems can be overcome a spectral model of the atmospheric circulation is chosen as an example of the application of the method.

## 5.1 INTRODUCTION

During the last decades the quality of forecasts of numerical weather models has been substantially improved. In general, this quality depends on at least three separate factors. First, there is the accuracy with which the mathematical model describes the *true* atmospheric circulation. Second, there is the accuracy with which the numerical model approximates the mathematical model. Third, there is the problem of estimating the initial state of the model sufficiently accurate. For weather models this last problem is promi-

nently present, because in these models there is sensitive dependence on the initial state. Starting from almost identical initial conditions, model solutions may diverge from each other so that integrating only a few days, little resemblance between the final states remains (Lorenz, 1963).

The rapid development of computer technology creates the possibility of improving both the mathematical and the numerical modeling: we can handle more complex models with a larger number of variables. The knowledge of the initial state has improved using more dense observation networks and more sophisticated data assimilation methods. Yet, in spite of all these improvements we are not able to give accurate forecasts beyond a range of about five days.

In this study we assume that observations are contaminated by discrete white Gaussian noise and that the model state equations are perturbed by a continuous white Gaussian noise process. In case the *true* process is described by this system and the statistics of these noise processes are known, there exists an optimal filter solution for the estimator of the state vector in the sense of minimum variance (Jazwinski, 1970). In case the system is linear and the initial state is Gaussian distributed, the probability distribution of the state vector remains Gaussian. Therefore, for an optimal filter for linear systems it is sufficient to predict and update the first and second moment of the probability distribution of the state vector; the resulting filter is the well-known Kalman-Bucy filter. On the other hand, for nonlinear systems the probability distribution of the state vector does not remain Gaussian. Therefore, this probability distribution is not described completely by its first and second moment, and consequently, an optimal filter cannot depend on these first two moments only. However, in many cases we may approximate the probability distributions of the state vector by a Gaussian distribution. A nonlinear filter that predicts and updates these first two moments generally works quite well (Anderson and Moore, 1979, Sorenson, 1988).

For linear models, a variety of adaptive filters has been developed that

deal with the uncertainty in the model. For meteorological practice, Dee et al. (1985) developed an efficient method to estimate the statistics of the model errors of a linearized weather model. Dee's method is based on a result of Bélanger (1974) for linear models: if the error covariance matrix of the model is a linear combination of a set of parameters, then the forecast error covariance matrix is also a linear combination of these parameters. Using this result it is possible to estimate these unknown parameters in a secondary filter.

Mous and Grasman (1993) described two methods to estimate uncertain model parameters. The first method is based on Lions' sentinel function (Lions, 1990). It compares a weighted average of the state vector, the so-called sentinel function, with a weighted average of the observations. The weights are chosen such that the sentinel function is, to first order, insensitive to errors in the initial state and can be easily calculated using the adjoint equations. A disadvantage of the method is that the sentinel function may be less sensitive to certain perturbations. These perturbations are called "stealthy" (Lions, 1988). The second method they described is based on an adaptive extended Kalman filter. In this method the uncertain parameters are regarded as random variables and the state vector is augmented with these variables. Since the extended Kalman filter yields estimates of the state vector, it also provides an estimate of the uncertain parameters. A disadvantage of this on-line estimating procedure is that initially the filter uses the wrong parameters, which lowers the performance of the filter. This disadvantage may be overcome by repeating the whole process using the estimated values for the parameters from the first run as initial estimates for the parameters for a second run. This off-line approach may have a better performance.

In this study an approximated Maximum Likelihood method is used to estimate the uncertain model parameters and the noise statistics. To this end, it is assumed that nonsystematic perturbations, numerical errors, the effect of

deleting higher order moments in filtering the data, etc., may be described by an additive white Gaussian noise process. In our method the data are filtered with the extended Kalman filter, because then we can easily obtain an approximation of the value of the log-likelihood function. We are aware of the computational costs of the extended Kalman filter. However, because the unknown parameters are estimated off-line, this is not our major concern.

In section 2 we show how the extended Kalman filter can be used in conjunction with an optimization algorithm to estimate the uncertain parameters. A complication in using the extended Kalman filter is that it may diverge, that is to say, the forecast errors may become inconsistent with their computed error statistics. To study this divergence we will analyze as an example the Lorenz attractor in section 3. We will show that divergence of the extended Kalman filter may lead to many local optima. A way to control this divergence is to add an artificial noise term to the state equations (Jazwinski, 1970). The statistics of this artificial noise term then have to be estimated. In section 4 we propose an optimization procedure for estimating the unknown model parameters and the noise statistics. This optimization procedure is designed in such a way that it does not tend to settle in a local optimum caused by divergence of the filter. We believe that this method may work well to estimate the uncertain parameters in nonlinear weather models, having a vacillating solution (de Swart, 1988). Vacillation means that the solution irregularly visits domains in state space where it remains for some time. These domains correspond to so-called preferent weather regimes. This behavior is very similar to that of the Lorenz equation. The Lorenz equation has a strange attractor, which consists of two "sheets" that are pinched together into a cantor book near the origin (Sparrow, 1982); the state vector switches irregularly between these two sheets. This kind of chaotic behavior is also characteristic for weather systems. The switching between the preferent domains of the system forms a problem for the extended Kalman filter; it easily leads to filter divergence. In section 5 the results of some numerical

experiments with the barotropic vorticity equation are presented. Here, the extended Kalman filter together with the proposed optimization procedure is used to estimate the parameters that describe a small systematic perturbation and the statistics of an artificial noise term.

## 5.2 EXTENDED KALMAN FILTERING

Starting point of our analysis is a system of  $(\hat{t}_0)$  stochastic nonlinear differential equations, which describes the time evolution of a state-vector  $x(t)=(x_1(t),\dots,x_n(t))'$ ,

$$\frac{dx}{dt} = f(x,t) + g(x,t;\lambda) + \xi_1(t), \quad t > t_0. \quad (5.1)$$

We assume that this system describes the *true* evolution of the physical process under study. The first term on the right-hand side of this system, the vector function  $f(x,t)$ , denotes the model of the physical processes. The second term on the right-hand side, the vector function  $g(x,t;\lambda)$  with unknown parameter vector  $\lambda$ , stands for small systematic perturbations that are neglected in the model. The last term stands for processes that inherently cannot be modeled, because they cannot be idealized. We think of nonsystematic perturbations, numerical errors, deletion of higher order moments of the probability distribution of the state vector, etc. As we will show later, it may also be seen as an artificial noise term that is used to solve the problem of filter divergence. We assume that these unknown noise processes are described by a white Gaussian noise process with autocovariance function

$$E(\xi_1(t)\xi_1(\tau)^T) = \sigma^2 W(t)\delta_p(t-\tau) \quad (5.2)$$

with  $\delta_p(t-\tau)$  the Dirac delta function. We further assume that the state vector is observed at discrete times  $t_k$  and that these observations

$y(t_k) = (y_1(t_k), \dots, y_m(t_k))'$  satisfy the so-called observation equation,

$$y(t_k) = m(x(t_k)) + \xi_2(t_k), \quad k=1, \dots, N. \quad (5.3)$$

Here,  $\xi_2(t_k)$  is the observation error. We assume that it can be described by a discrete white Gaussian noise process with autocovariance function

$$E(\xi_2(t_k)\xi_2(t_l)^T) = \sigma^2 R(t_k) \delta_K(t_k - t_l) \quad (5.4)$$

with  $\delta_K(t_k - t_l)$  the Kronecker delta. We also assume that the noise processes  $\xi_1(t)$  and  $\xi_2(t_k)$  are independent. For more complicated systems, e.g., with correlation between the noise processes, we refer to Jazwinski (1970). In this reference a thorough description of nonlinear filtering theory is presented. We will use the extended Kalman filter to calculate the forecast and the filtering solution and for completeness we will give a brief description of this filter.

Starting from the filtering solution  $\hat{x}(t_k|t_k)$ , given the observations  $y(t_1), \dots, y(t_k)$ , the forecast  $\hat{x}(t_{k+1}|t_k)$  is simply given by integrating eq. (5.1), with  $\xi_1(t)$  set to its mean value  $E(\xi_1(t))=0$  and with an initial guess for the parameter vector  $\lambda$ . In meteorology the filtering solution is often called the analysis. For the calculation of the error covariance matrix  $P(t_{k+1}|t_k)$  at time  $t_{k+1}$  the state transition matrix or resolvent of the linearized system is needed. This calculation is very laborious and is often considered as the major weakness of the extended Kalman filter. Dee (1991) presented a simplified Kalman filter, which is based on a simple model for the propagation of the model error. He showed that this filter may work quite well for weather forecasting. However, we believe that a good approximation of this error covariance matrix is necessary to make a good judgement of the quality of the model. Therefore, we will not use a simplified form of the state transition matrix to calculate the propagation of the model error, but we will use the complete state transition matrix of the linearized model equations. The propagation of the forecast  $x(t_{k+1}|t_k)$  and error covariance matrix  $P(t_{k+1}|t_k)$  are then calculated by

$$\hat{x}(t_{k+1}|t_k) = \hat{x}(t_k|t_k) + \int_{t_k}^{t_{k+1}} f(\hat{x}(t|t_k)) + g(\hat{x}(t|t_k); \lambda) dt, \quad (5.5)$$

$$P(t_{k+1}|t_k) = \Phi(t_{k+1}, t_k; \hat{x}(t_k|t_k)) P(t_k|t_k) \Phi^T(t_{k+1}, t_k; \hat{x}(t_k|t_k)) + Q(t_{k+1}), \quad (5.6)$$

where  $\Phi(t_{k+1}, t_k; \hat{x}(t_k|t_k))$  is the state transition matrix of the linearized state equation

$$\frac{d(\delta x)}{dt} = (f'(\hat{x}(t_k|t_k)) + g'(\hat{x}(t_k|t_k); \lambda)) \delta x \quad (5.7)$$

with

$$f'(\hat{x}(t_k|t_k)) = \left[ \frac{\partial f(\hat{x}(t_k|t_k))}{\partial x} \right]_{n \times n}, \quad g'(\hat{x}(t_k|t_k); \lambda) = \left[ \frac{\partial g(\hat{x}(t_k|t_k); \lambda)}{\partial x} \right]_{n \times n}. \quad (5.8)$$

The matrix  $Q(t_{k+1})$  in eq. (5.6) is related to the autocovariance function of the white Gaussian process of the model,  $\xi_1(t)$ , by

$$Q(t_{k+1}) = \int_{t_k}^{t_{k+1}} \Phi(t_{k+1}, \tau) W(\tau) \Phi^T(t_{k+1}, \tau) d\tau. \quad (5.9)$$

Observation  $y(t_{k+1})$  are used to update the forecast and the error covariance matrix,

$$\hat{x}(t_{k+1}|t_{k+1}) = \hat{x}(t_{k+1}|t_k) + K(t_{k+1})(y(t_{k+1}) - m(\hat{x}(t_{k+1}|t_k))), \quad (5.10)$$

$$P(t_{k+1}|t_{k+1}) = (I - K(t_{k+1})M(t_{k+1}))P(t_{k+1}|t_k) \quad (5.11)$$

with

$$K(t_{k+1}) = P(t_{k+1}|t_k)M(t_{k+1})^T(M(t_{k+1})P(t_{k+1}|t_k)M(t_{k+1})^T + R(t_{k+1}))^{-1}, \quad (5.12)$$

$$M(t_{k+1}; \hat{x}(t_{k+1}|t_k)) = \left[ \frac{\partial m(\hat{x}(t_{k+1}|t_k))}{\partial x} \right]_{m \times n}. \quad (5.13)$$

The new analysis  $\hat{x}(t_{k+1}|t_{k+1})$  and error covariance matrix  $P(t_{k+1}|t_{k+1})$  are used to calculate a new forecast. Of course it is unlikely that the initial state  $x_0$  belonging to the state eq. (5.1), which is needed to initialize the filter, is known exactly. Therefore, we assume that the initial state is a random vector with mean  $x_0$  and with covariance matrix  $P_0$ . Furthermore, we assume that the initial state is independent of the noise processes  $\xi_1(t)$  and  $\xi_2(t)$ . The recursive extended Kalman filter is then started with mean  $\hat{x}(t_0|t_0)=x_0$  and covariance matrix  $P(t_0|t_0)=P_0$ .

The extended Kalman filter works well if the probability distribution of the state vector  $x(t)$  can be approximated by a Gaussian distribution. This approximation is valid in case the linearized state equation describes the process between  $t_k$  and  $t_{k+1}$  adequately. In practice this means that the observations are made with a high sampling frequency, and that the model is an accurate description of the *true* system. To check whether the extended Kalman filter is doing well, we calculate the forecast errors,  $e_i(t_{k+1})=y_i(t_{k+1})-m(\hat{x}_i(t_{k+1}|t_k))$ . For an optimal filter these forecast errors have to be white. Therefore, the more white they are the better the filter works (Anderson and Moore, 1977). A simple test for the whiteness of the forecast errors is based on the autocorrelation functions. First the forecast error time series are normalized, so that the forecast errors are approximately identically distributed. The autocovariance functions are then estimated by

$$\hat{C}_{i,i}(k) = \frac{1}{N} \sum_{l=1}^N e_i(t_l) e_i(t_{l+k}), \quad (5.14)$$



and the autocorrelation by

$$\hat{\Gamma}_{i,i}(k) = \frac{\hat{C}_{i,i}(k)}{\hat{C}_{i,i}(0)}. \quad (5.15)$$

The variances of the estimated autocorrelation of a time series with independent, identically distributed normal errors may be approximated by the formula of Moran (1947),

$$\text{var}(\hat{\Gamma}_{i,i}(k)) = \frac{N-k}{N(N+2)}. \quad (5.16)$$

The 95% critical values for  $\hat{\Gamma}_{i,i}(k)$  are then  $\pm(1.96 \text{ var}(\hat{\Gamma}_{i,i}(k))^{1/2})$ . In case less than 5% of  $\hat{\Gamma}_{i,i}(k)$  lies outside this band, we say that the sequence is white. In this approach, the forecast errors are tested by component. An alternative may be found in testing them simultaneously, as indicated by Mehra (1970).

In many cases the model contains a number of unknowns, such as the model parameters  $\lambda$  and the matrices  $Q(t_k)$  and  $R(t_k)$ . The importance of having a good approximation of the matrix  $Q(t_k)$  can be understood from eq. (5.6); a large matrix  $Q(t_k)$  may dominate the error covariance matrix  $P(t_{k+1}|t_k)$  and the Kalman gain  $K(t_{k+1})$  making the new analysis  $\hat{x}(t_{k+1}|t_{k+1})$  unreliable. Therefore one needs a good approximation of this matrix. It should be estimated from the data. However, it is impossible to estimate all elements of this matrix, even if we assume that this matrix is time-invariant. Instead of estimating all elements of  $Q(t_k)$  we rather estimate certain linear combination of the elements, by approximating  $Q(t_k)$  by a linear combination of a fixed set of time independent matrices (Cohn and Parrish, 1991):

$$Q(t_k) = Q = \sum \mu_i Q_i. \quad (5.17)$$

The matrices  $Q_i$  are chosen a priori. The parameters  $\mu_i$ , which in the following are called noise parameters, are yet to be determined. The other matrix  $R(t_k)$  is a measure for the accuracy of the observations. It is important for the

update equations of the filter. However, because in practice one often knows the accuracy of the observations quite accurately, we assume that this matrix is known.

Using the extended Kalman filter an approximation of the log-likelihood function can be obtained. To do so the log-likelihood function is written as

$$L(y_N, y_{N-1}, \dots, y_1; \sigma^2, \lambda, \mu) = L(y_N | y_{N-1}, \dots, y_1; \sigma^2, \lambda, \mu) + L(y_{N-1}, \dots, y_1; \sigma^2, \lambda, \mu) \quad (5.18)$$

(cf., e.g., Harvey, 1981). As mentioned above the probability distributions of the forecast errors are approximated by Gaussian distributions. The mean and covariance matrices of these distributions are given by

$$E(e(t_{k+1})) = E(y(t_{k+1}) - m(\hat{x}(t_{k+1} | t_k))) = 0, \quad (5.19)$$

$$\begin{aligned} E(e(t_{k+1})e(t_{k+1})^T) &= \sigma^2(R(t_{k+1}) + M(t_{k+1})P(t_{k+1} | t_k)M(t_{k+1})^T), \\ &= \sigma^2 H(t_{k+1}). \end{aligned} \quad (5.20)$$

The second term on the right-hand side of eq. (5.18) is then equal to

$$\begin{aligned} L(y_N | y_{N-1}, \dots, y_1; \sigma^2, \lambda, \mu) &= \\ &= -\frac{m}{2} \ln(2\pi) - \frac{m}{2} \ln(\sigma^2) - \frac{1}{2} \ln(|H(t_k)|) - \frac{1}{2} \sigma^{-2} e^T(t_k) H(t_k) e(t_k). \end{aligned} \quad (5.21)$$

Substitution of eq. (5.21) into eq. (5.18) shows that the approximate log-likelihood differs by no more than an additive constant from

$$\begin{aligned} L(y_N, \dots, y_1; \sigma^2, \lambda, \mu) &= \\ &= -\frac{1}{2} N m \ln \sigma^2 - \frac{1}{2} \sum_{k=1}^N \ln(|H(t_k)|) - \frac{1}{2} \sigma^{-2} \sum_{k=1}^N e^T(t_k) H(t_k)^{-1} e(t_k). \end{aligned} \quad (5.22)$$

We can simply reduce the dimension of the optimization problem by one, since  $\sigma^2 = SS/Nm$ , with

$$SS = \sum_{k=1}^N e^T(t_k) H(t_k)^{-1} e(t_k), \quad (5.23)$$

optimizes  $L(\sigma^2, \lambda, \mu | y_N, \dots, y_1)$  for all  $\lambda$  and  $\mu$ . The object function for the

optimization of  $\lambda$  and  $\mu$  is therefore equal to

$$L^*(y_N, \dots, y_1; \lambda, \mu) = -\frac{1}{2}Nm \ln \left( \frac{\sum_{k=1}^N e(t_k)^T H(t_k)^{-1} e(t_k)}{Nm} \right) - \frac{1}{2} \sum_{k=1}^N \ln(|H_k|). \quad (5.24)$$

Since the system is nonlinear, one cannot give an explicit expression for the approximated Maximum Likelihood estimates of  $\lambda$  and  $\mu$ . We may obtain these estimates using the extended Kalman filter in conjunction with an optimization algorithm, e.g., the Quasi-Newton algorithm or the Conjugated Gradients algorithm. However, after getting some practical experience with a chaotic system, we concluded that these algorithms rarely converge. The reason is that the log-likelihood as function of  $\lambda$  and  $\mu$  has many local optima. These local optima occur due to divergence of the extended Kalman filter. It is important to have a good understanding of this phenomenon, because the cause of this phenomenon is important for constructing a robust optimization procedure and for finding good initial guesses. Therefore, we will study a simple chaotic system as an example.

### 5.3 NUMERICAL EXPERIMENTS WITH THE LORENZ ATTRACTOR

The purpose of these experiments is twofold. First, we will study the use of the extended Kalman filter for chaotic models. In particular we will study divergence of the filter. Second, we will analyze the sensitivity of the log-likelihood function as function of the unknown parameters. Due to filter divergence the log-likelihood function will have many local optima. The hierarchical optimization procedure, that is presented in the next section, is designed to deal with this problem. We have chosen the Lorenz attractor as example because the small dimension of the state vector reduces the compu-

tational cost of the filter. The state equations of the Lorenz attractor are given by

$$\begin{aligned}\frac{dx_1}{dt} &= c(x_3 - x_1) + \xi_{1,1}, \\ \frac{dx_2}{dt} &= -x_1x_3 + rx_1 - x_2 + \xi_{1,2}, \\ \frac{dx_3}{dt} &= x_1x_2 - bx_3 + \xi_{1,3}\end{aligned}\tag{5.25}$$

(Lorenz, 1963). We take  $c=10$ ,  $r=48$  and  $b=8/3$ . It is remarked that we have added small stochastic perturbation terms to the equations. Usually these stochastic terms are not present in this system. Here, these stochastic terms are used as an artificial method to account for deleting higher order moments, numerical errors, etcetera.

For the numerical experiments we have generated a dataset of the deterministic system by integrating the state equation using a high-order Runge-Kutta method with variable time-step. All state variables are observed with a sampling period  $T=t_k-t_{k-1}=0.1$ . We have contaminated the observations with generated pseudo-random numbers from a normal distribution  $N(0,1)$ .

In the first experiment we take

$$P_0 = \sigma_p^2 I, \quad Q(t_k) = 0, \quad R(t_k) = \sigma_R^2 I$$

with  $I$  the identity matrix and with  $\sigma_p^2=1$  and  $\sigma_R^2=1$ . Figure 5.1(a) shows the forecast errors,  $e(t_{k+1})=y_i(t_{k+1})-\hat{x}_i(t_{k+1}|t_k)$ . After some time the forecast errors suddenly increase, which indicates filter divergence. The reason for this sudden increase of the forecast errors is that the filter "assumes" that the state is winding around one side of the attractor, while the *true* system is winding around the other side. The filter is not able to track a sudden change in the trajectory, because the increased accuracy of the estimate of the state vector causes the Kalman gain to become small and consequently the next consecutive observations have almost no effect. Therefore, for vacillating systems we

have to modify the state equations by adding a small artificial noise term to them, so that the extended Kalman filter is able to follow the observations.

In the second experiment we study the effect of such an artificial noise term. We have chosen for a simple parameterization of the matrices  $Q(t_k)$ ,

$$Q(t_k) = Q = \mu I$$

and we take  $\mu=0.001$ . The forecast errors in figure 5.1(b), indicate no filter divergence. Apparently, the filter does not learn the wrong state too well. An other way to analyze the working of the filter is to calculate the autocorrelation function of the normalized forecast errors. In figure 5.2(b) the straight lines denote the critical values for the autocorrelation function. Since the autocorrelation function lies for more than 95% between these lines, we conclude that the filter works well.

Filter divergence affects the calculation of the approximated log-likelihood: the approximated log-likelihood will have many local optima as function of  $\lambda$  and  $\mu$ . Therefore, the usefulness of the approximated log-likelihood as object function seems questionable. The previous two experiments showed that the filter is not divergent in case we add an artificial noise term to the state equations. The response surface of the log-likelihood function will be smooth in the region of the parameter space where  $\mu$  and thus  $Q$  is large. In many cases the log-likelihood function will have an optimum in this region that locates the *true* value of  $\lambda$ . To illustrate this, we have perturbed the state equation by adding a small external force

$$g(x,t;\lambda) = (g_1(x,t;\lambda), g_2(x,t;\lambda), g_3(x,t;\lambda)) = (\lambda, 0, 0)$$

to the state equations, with an actual value  $\lambda=5.0$ . This external force does not cause a dramatic change in the qualitative behavior of the Lorenz attractor. We have studied the sensitivity of the approximated log-likelihood function as a tool to estimate the value of  $\lambda$ , for various values of  $N$  and of

$\mu$ . The results of the sensitivity analysis are summarized in figure 5.3. For short datasets,  $N=20$ , the approximated log-likelihood has a well-defined minimum. However, the accuracy of  $\hat{\lambda}$  is low and the log-likelihood is not very sensitive to  $\lambda$ . Comparison of figures 5.3(a),(b),(c),(d) shows that the log-likelihood is not very sensitive to  $\mu$  either. For larger datasets, the log-likelihood has many local minima. The reason that larger datasets give rise to many local minima is that the filter starts to diverge after the processing of some data. The accuracies of estimates obtained with larger datasets is, however, substantially higher. In figure 5.3 this trend is seen: a steeper minimum corresponds to a more accurate estimate (note the different scaling). Accurate estimates of  $\lambda$  are obtained even with initially overestimated noise parameters  $\mu$ . Accurate estimates of  $\mu$  are harder to obtain, because with an underestimation of  $\mu$  the filter easily diverges (see figures 5.3(k) and 5.3(l)). Although accurate estimates of  $\lambda$  can be obtained with an overestimation of  $\mu$ , these noise parameters are important because they identify the optimal extended Kalman filter.

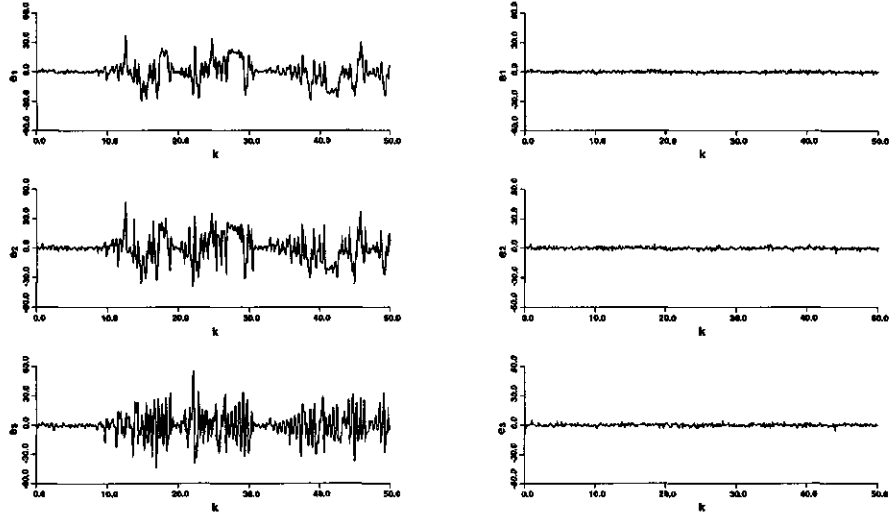


Figure 5.1 Estimated forecast errors, (a) with model error covariance matrix  $Q=0$ , (b) with model error covariance matrix  $Q=0.001I$ .

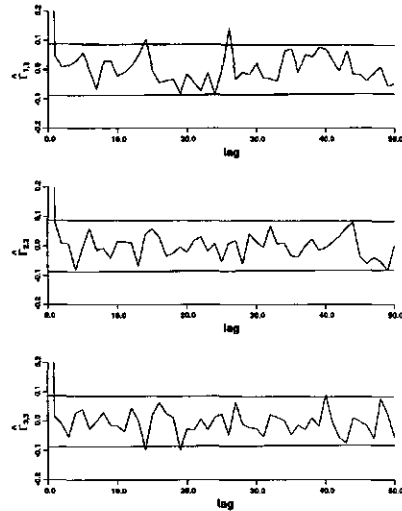


Figure 5.2 Autocorrelation functions of the forecast errors.

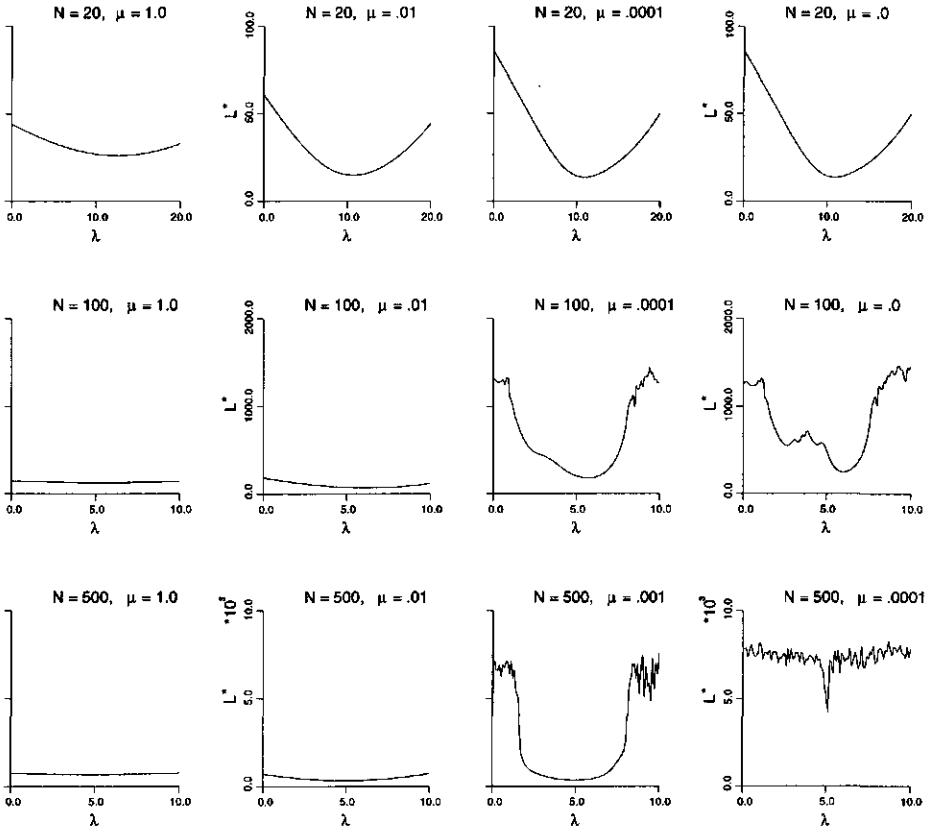


Figure 5.3 Sensitivity of the log-likelihood.

## 5.4 HIERARCHICAL OPTIMIZATION

The problems and experiments discussed in the previous section lead to the following requirements for the optimization procedure:

- a partial decoupling of the optimization of the object function with respect to the model parameters  $\lambda$  and the noise parameters  $\mu$ ,
- a simple method for checking filter divergence.

The idea of a partial decoupling is based on the experience that  $\lambda$  is estima-



ted rather well when we initially overestimate the noise parameters  $\mu$ . A second argument for partial decoupling is that filter divergence is strongly related to the estimate of  $\mu$ . This divergence makes it difficult to estimate the parameters  $\mu$ . For chaotic models, filter divergence can easily be checked with a visual inspection of the time series of the forecast error, because a sudden increase of this time series indicates filter divergence.

The procedure for hierarchical optimization is schematically presented in figure 5.4. Initially we take  $\mu$  large and we use some initial guess for  $\lambda$ . The forecast errors are calculated using the extended Kalman filter and are visually checked on divergence. If necessary we adjust the parameters  $\mu$ . The matrices  $K(t_k)$  and  $H^1(t_k)$  that are computed at each moment  $t_k$  during this run are stored.

The model parameters  $\lambda$  are estimated, with the noise parameters  $\mu$  fixed, using a Quasi-Newton algorithm. This is called the inner optimization. We remark that initially the large matrix  $Q$  dominates the error covariance matrix  $P(t_{k+1}|t_k)$  and therefore also the Kalman gain. Since the matrix  $Q$  does not change within the inner optimization, we may use the stored matrices  $K(t_k)$  and  $H^1(t_k)$  to approximate the actual matrices. This reduces the computational costs of the inner optimization considerably.

In the other optimization the parameters  $\mu$  are estimated, with the parameters  $\lambda$  fixed. This is called the outer optimization. As mentioned before, filter divergence occurs in case the matrix  $Q$  is too small. Therefore, we adjust the estimate of  $\mu$  with small steps. In the outer optimization only one step is made in the steepest descent direction. The new matrices  $K(t_k)$  and  $H^1(t_k)$  are again stored so that they can be used in the inner optimization. We also check whether the filter diverges and if necessary we adjust the update for the parameters  $\mu$ .

In the inner optimization the dominance of the matrix  $Q$  decreases after a few iterations. Our experience is that it is not necessary to calculate the matrices  $K(t_k)$  and  $H^1(t_k)$  every time the extended Kalman filter algorithm is

called, because the forecast errors are more sensitive to a small change in  $\lambda$  than the matrix  $H^{-1}(t_k)$ .

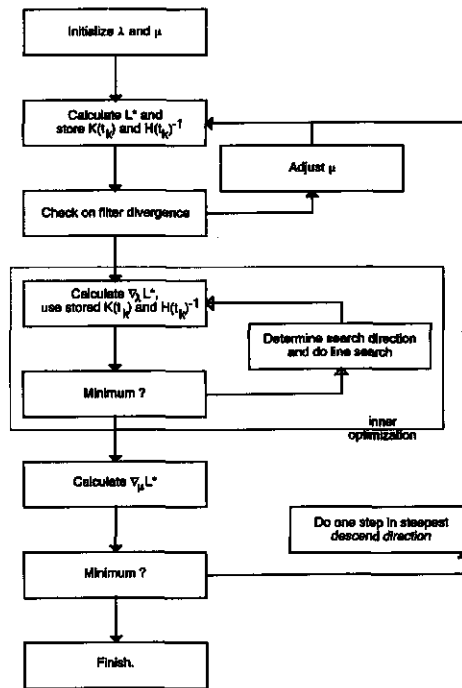


Figure 5.4 Optimization procedure

## 5.5 NUMERICAL EXPERIMENTS WITH THE BAROTROPIC VORTICITY EQUATION

The barotropic vorticity equation describes the dynamics of a two-dimensional non-divergent and inviscid flow at the surface of a rotating sphere. It can be used as an approximating model of the atmospheric circulation at 500 mbar. The forecast of the flow field is not very accurate, because of the intrinsic growth of errors in the initial state and because several physical processes are neglected in the model. Some of these processes, for example the effect of large mountains, can be modeled by small systematic perturbation terms (de Swart, 1988). On the other hand there are also processes that cannot be modeled easily by systematic perturbation terms, like the effect of cultivation of the land. Assuming that we may describe these processes by adding noise to the state-equation, the method described in the previous sections can be used to assess the different sources of model deficiencies.

The barotropic vorticity equation we consider in this section is given by

$$\frac{\partial \zeta}{\partial t} = J(\zeta + f, \psi), \quad (5.26)$$

where  $\zeta(\alpha, \beta, t)$  is the relative vorticity,  $\psi(\alpha, \beta, t)$  is the streamfunction,  $f$  is the Coriolis parameter and  $J$  is the Jacobi operator. This operator is defined by

$$J(g, h) = \frac{\partial g}{\partial \alpha} \frac{\partial h}{\partial \beta} - \frac{\partial g}{\partial \beta} \frac{\partial h}{\partial \alpha}, \quad (5.27)$$

where  $\alpha$  is the geographic longitude and  $\beta$  is the sine of the geographic latitude. The radius of the earth and the inverse of the angular speed of rotation of the earth are used as unit of length and time respectively. The relative vorticity is related to the streamfunction by  $\zeta = \Delta \psi$ . As state variable

we may choose either  $\zeta(\alpha, \beta, t)$  or  $\psi(\alpha, \beta, t)$ . In this study, we have chosen for the streamfunction  $\psi(\alpha, \beta, t)$ .

We approximate the solution  $\psi(\alpha, \beta, t)$  of eq. (5.26) by expanding it in spherical harmonics and using a triangular truncation at T11,

$$\psi(\alpha, \beta, t) = \sum_{m=-11}^{m=11} \sum_{n=|m|}^{n=11} x_{mn}(t) Y_{mn}(\alpha, \beta) \quad (5.28)$$

with

$$Y_{mn}(\alpha, \beta) = P_{mn}(\beta) \exp(im\alpha).$$

The functions  $P_{mn}(\beta)$  denote the associated Legendre functions of the first kind and of order  $m$  and degree  $n$ . Since the spherical harmonics are eigenfunctions of the Laplace operator, substitution of eq. (5.28) into eq. (5.26) leads to system of ordinary differential equations for  $x(t)$  which is of the form of eq. (5.1) (Machenhauer, 1979).

The state transition matrix of the linearized system can be obtained by linearizing the barotropic vorticity equation around  $\psi(\alpha, \beta, t)$ . We then obtain the so-called tangent equation:

$$\frac{\partial \varepsilon}{\partial t} = \Delta^{-1} J(\Delta \varepsilon, \psi) + \Delta^{-1} J(\Delta \psi + f, \varepsilon) \quad (5.29)$$

(Barkmeyer, 1992). Expanding the solution  $\varepsilon(\alpha, \beta, t)$  of this equation also in spherical harmonics using a T11 truncation we obtain the state transition matrix by integrating this equation one row at a time.

Numerical experiments have been carried out to test the performance of the extended Kalman filter and the hierarchical optimization procedure for this application. In the same way as in the example of the Lorenz equation, we perturbed some of the state equations with a small constant perturbation. Here we have perturbed the state equations for  $x_{0,i}$ ,  $i=1, \dots, 4$  with

$$\begin{aligned} g_{0,1}(x,t;\lambda) &= \lambda_1 = 1.0 \times 10^{-5}, & g_{0,2}(x,t;\lambda) &= \lambda_2 = 2.0 \times 10^{-5}, \\ g_{0,3}(x,t;\lambda) &= \lambda_3 = 3.0 \times 10^{-5}, & g_{0,4}(x,t;\lambda) &= \lambda_4 = 4.0 \times 10^{-5}, \end{aligned}$$

The reference trajectory is then calculated by integrating the state equations, starting from an arbitrary initial state, with the stochastic noise term set to zero. All the state variables are observed with a sampling period of  $T=12$  h. and we have contaminated the observations with generated pseudo-random numbers from a normal distribution  $N(0, \sigma_R^2)$ .

In the first experiment we take  $\sigma_R^2=1.0 \times 10^{-10}$ . Since the magnitude of the state variables lies in the order of  $1.0 \times 10^{-3}$ , this agrees with an accuracy of the observations of 0.1-1%. We set the covariance matrices in the extended Kalman filter to

$$P_0 = 1.0 \times 10^{-10} I, \quad Q(t_k) = \mu I, \quad R(t_k) = 1.0 \times 10^{-10} I$$

and we start the optimization procedure with  $\mu=1.0 \times 10^{-3}$  and  $\lambda_i=0$ ,  $i=1, \dots, 4$ . The optimum of  $L^*(\lambda, \mu | y_N, \dots, y_1)$  is then found at

$$\begin{aligned} \hat{\lambda}_1 &= 0.972 \times 10^{-5}, & \hat{\lambda}_2 &= 1.987 \times 10^{-5}, \\ \hat{\lambda}_3 &= 3.007 \times 10^{-5}, & \hat{\lambda}_4 &= 3.867 \times 10^{-5}, \\ \hat{\mu} &= 0.0. \end{aligned}$$

Since  $\hat{\mu}=0.0$ , the likelihood criterion indicates that the state equations are deterministic. Therefore the covariance matrix of  $\lambda_1, \lambda_2, \lambda_3$  and  $\lambda_4$  is approximated by the inverse of the Hessian matrix of  $L^*(\lambda, 0 | y_N, \dots, y_1)_{\lambda=\hat{\lambda}}$ . This matrix is given in table 5.1. It shows that the accuracy of the estimates strongly varies, with  $\hat{\lambda}_1$  is the most accurate and  $\hat{\lambda}_4$  is the less accurate one.

The performance of the extended Kalman filter, using the optimal parameters, is analyzed by calculating the autocorrelation functions of the forecast errors. For illustration, we show in figure 5.5(a) the autocorrelations of the forecast errors  $e_{0,1}(t_{k+1})=y_{0,1}(t_{k+1})-\hat{x}_{0,1}(t_{k+1}|t_k)$ ,  $e_{-5,5}(t_{k+1})=y_{-5,5}(t_{k+1})-\hat{x}_{-5,5}(t_{k+1}|t_k)$

Table 5.1 Inverse Hessian matrix of  $L^*(\lambda, 0|y_N, \dots, y_1)_{\lambda=\hat{\lambda}}$ 

	0.0918	0.0140	-0.0360	0.0167
$1.0 \times 10^{-14} *$	0.0140	1.2537	-0.2226	0.9685
	-0.0360	-0.2226	4.6671	-0.4348
	0.0167	0.9685	-0.4348	22.4646

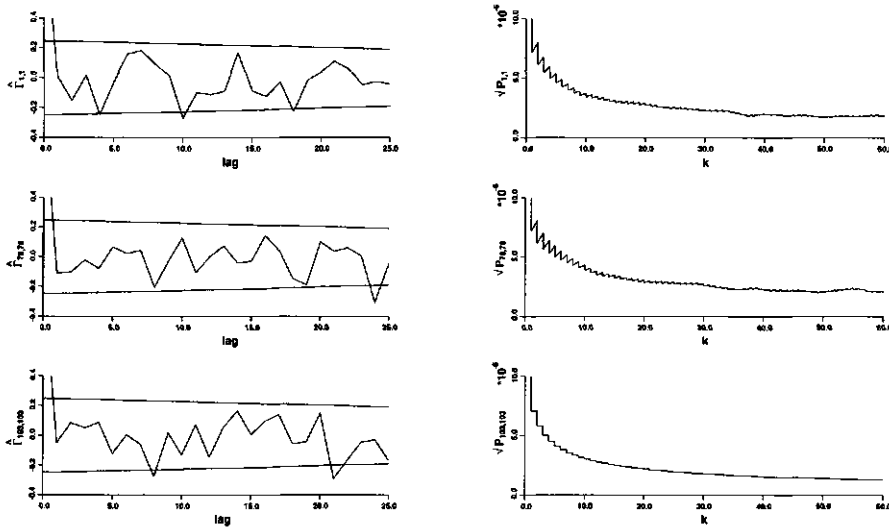


Figure 5.5 (a) Autocorrelation functions of the forecast errors;  $\hat{\Gamma}_{1,1}$  is the autocorrelation function of the forecast errors time series  $e_{0,1}(t_{k+1}) = y_{0,1}(t_{k+1}) - \hat{x}_{0,1}(t_{k+1}|t_k)$ , and  $\hat{\Gamma}_{78,78}$  and  $\hat{\Gamma}_{103,103}$  are the autocorrelations of the forecast errors time series  $e_{-5,5}(t_{k+1}) = y_{-5,5}(t_{k+1}) - \hat{x}_{-5,5}(t_{k+1}|t_k)$  and  $e_{5,5}(t_{k+1}) = y_{5,5}(t_{k+1}) - \hat{x}_{5,5}(t_{k+1}|t_k)$  respectively. (b) Accuracy of these forecast errors time series.

and  $e_{5,5}(t_{k+1}) = y_{5,5}(t_{k+1}) - \hat{x}_{5,5}(t_{k+1}|t_k)$ . Since more than 5% lies outside the confidence band, the forecast errors are not white. Figure 5.5(b) shows the evolution of the square root of the variances  $\sqrt{P(t_k)}$  or root mean square error curve that go with these forecast error time series. From this figure it is clear

that the performance of the filter is still dominated by the initial variances. The long settling time of the filter and the nonwhite forecast errors indicate that in this case more data are required to obtain an accurate estimate of the model error covariance matrix.

In the second experiment we analyze the case with less accurate observations. We take  $\sigma_R^2 = 1.0 \times 10^{-10}$  and we set the covariance matrices in the extended Kalman filter to

$$P_0 = 2.5 \times 10^{-7} I, \quad Q(t_k) = \mu I, \quad R(t_k) = 2.5 \times 10^{-7} I$$

We use the same starting values for  $\mu$ ,  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$  and  $\lambda_4$ . The optimization procedure then yields an optimum of  $L^*(\lambda, \mu | y_N, \dots, y_1)$  at

$$\begin{aligned} \hat{\lambda}_1 &= 1.242 \times 10^{-5}, & \hat{\lambda}_2 &= 1.505 \times 10^{-5}, \\ \hat{\lambda}_3 &= -9.65 \times 10^{-5}, & \hat{\lambda}_4 &= 6.566 \times 10^{-5}, \\ \hat{\mu} &= 1.0 \times 10^{-9}. \end{aligned}$$

The approximated covariance matrix shows that the estimates for  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$  and  $\lambda_4$  are far less reliable than in the previous experiment. In this second experiment we have increased the observation error causing an increased effect of deleting higher order moments. This is reflected by the large estimated artificial noise terms in the state equations. The autocorrelations, shown in figure 5.6(a), are more white than in the previous experiment. However, the filter is still not working optimally. The evolution of the square root of the variances of the forecast error time series show that the settling time of the filter is much shorter, but it also shows that the accuracies of the forecast and analysis are less accurate.

Although the forecast errors are in both cases not completely white the extended Kalman filter did not diverge. The advantage of a non-diverging filter is that we use all the available information to produce an forecast. For a small artificial noise term the improvement of the quality of the forecast is significant as can be seen from the root mean square error curves.

Table 5.2 Inverse Hessian matrix of  $L^*(\lambda, \mu | y_N, \dots, y_1)_{\lambda=\hat{\lambda}, \mu=\hat{\mu}}$ 

	0.0712	0.0031	0.0650	0.0403	0.0
$1.0 \times 10^{-10} *$	0.0031	3.8369	-2.1553	0.2416	0.0
	0.0650	-2.1553	9.8911	5.0725	0.0
	0.0403	0.2416	5.0725	8.6555	0.0
	0.0	0.0	0.0	0.0	$1.1 \times 10^{-10}$

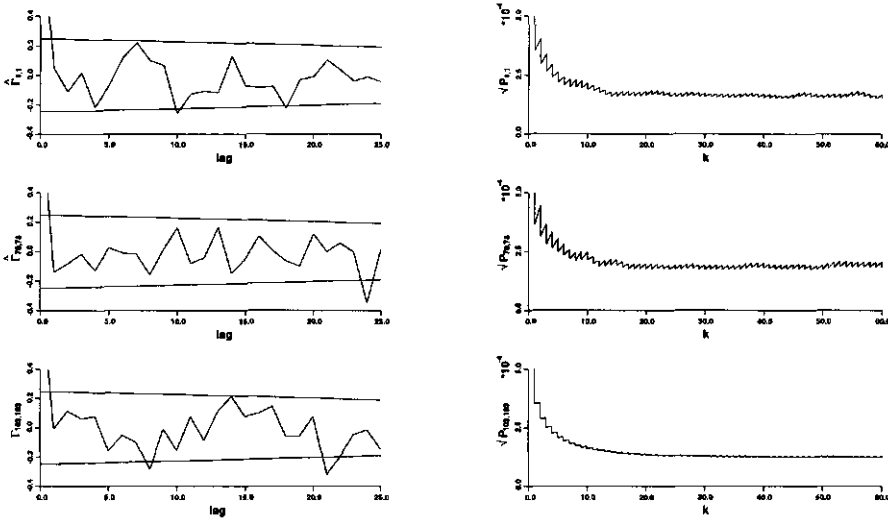


Figure 5.6 (a) Autocorrelation functions of the forecast errors;  $\hat{\Gamma}_{1,1}$  is the autocorrelation function of the forecast errors time series  $e_{0,1}(t_{k+1}) = y_{0,1}(t_{k+1}) - \hat{x}_{0,1}(t_{k+1}|t_k)$ , and  $\hat{\Gamma}_{78,78}$  and  $\hat{\Gamma}_{103,103}$  are the autocorrelations of the forecast errors time series  $e_{5,5}(t_{k+1}) = y_{5,5}(t_{k+1}) - \hat{x}_{5,5}(t_{k+1}|t_k)$  and  $e_{5,5}(t_{k+1}) = y_{5,5}(t_{k+1}) - \hat{x}_{5,5}(t_{k+1}|t_k)$  respectively. (b) Accuracy of these forecast errors time series.



## 5.6 CONCLUSIONS

If a non-linear filter is operating over a large range of time (large datasets), filter divergence often occurs when determining the state of vacillating systems. A sudden increase of the forecast errors shows the beginning of divergence. The reason for this sudden increase is that the filter "assumes" that the state is in one preferent domain whereas the *true* state is in an other preferent domain. Imperfect knowledge of the forecast model, uncertain noise parameters and deletion of higher order moments of the probability distributions may give rise to filter divergence. A reasonable approach to this problem is to add a noise term to the state equation of the model. In this way we can control the accuracy of the state estimates and thus prevent divergence. The optimal working extended Kalman filter for such a system reduces the error in tracking this system and in this sense it is related to the shadowing theory of chaotic systems (Farmer and Sidorowich, 1992)

We have used the Maximum Likelihood method to estimate uncertain model parameters and noise parameters. An approximation of the value of the log-likelihood function can easily be obtained with the extended Kalman filter. However, filter divergence causes this function to have many local optima. A disadvantage of standard optimization algorithms is that they may converge to such a local optimum. To overcome this problem we have developed a new procedure that does not tend to settle in a local optimum caused by divergence of the filter. This procedure is based on a partial decoupling of the optimization with respect to the model and noise parameters. By checking on filter divergence at strategic moments during the optimization and by adjusting the noise parameters if necessary, the procedure will not converge to an incorrect local optimum. An additional advantage of this procedure is that we may approximate the Kalman gain in the inner optimization, which reduces the computational costs considerably.

The procedure was applied to estimate uncertain model and noise

parameters of the barotropic vorticity equation. The uncertain model parameters belong to a small perturbation term and the noise parameters are part of the parameterization of the model error covariance matrix. This model error covariance matrix is used in the extended Kalman filter to account for the effect of deleting higher order moments and numerical errors. In the first experiment we have used very accurate observations. The Maximum Likelihood method suggests that the noise term in the state equation may be neglected since we found the optimum at  $Q=0$ . Nevertheless we think that the filter works best with a small model error, because the estimated autocorrelation functions of the forecast errors indicate that the forecast errors are not completely white. Since we have used a rather short dataset, which causes the filter to be dominated by the initial error covariance matrix, we may not conclude that the process is best described by a "deterministic" system. A comparison with the example of the Lorenz attractor suggests that an analysis with more data may give a decisive answer to this question.

In the second example we have used inaccurate observations. Here we found that the performance of the extended Kalman filter is nearly optimal in case we add a small artificial noise term to the state equation. The autocorrelation functions show that the forecast errors are more white than in the first experiment. Although the extended Kalman filter seems to work well, the forecasts are not very accurate, because the estimated model error covariance matrix  $Q$  is large. On the other hand we considerably reduced the divergence of the filter.

## REFERENCES

- Anderson, B.D.O. and J.B. Moore, *Optimal Filtering*, Englewood Cliffs, Prentice-Hall, 1979.
- Barkmeyer, J., Local error growth in a barotropic model, *Tellus*, 44A, 314-323, 1992.
- Bélanger, P.R., Estimation of noise covariance matrices for a linear time-varying stochastic process, *Automatica*, 10, 267-275, 1974.
- Cohn, S.E. and D.F. Parrish, The behaviour of forecast error covariances for a Kalman filter in two dimensions, *Mon. Wea. Rev.*, 119, 1757-1785, 1991.
- Dee, D.P., S.E. Cohn, A. Dalcher and M. Ghil, An efficient algorithm for estimating noise covariances in Distributed Systems, *IEEE Trans. Aut. Control*, AC-30, 11, 1057 - 1065, 1985.
- Dee, D.P., Simplification of the Kalman filter for meteorological data assimilation, *Quart. J. Roy. Meteor. Soc.*, 9-16, 1991.
- De Swart, H.E., *Vacillation and predictability properties of low-order atmospheric spectral models*, Ph.D.Thesis, Rijksuniversiteit Utrecht, 1988.
- Farmer, J.D. and J.J. Sidorowich, Optimal shadowing and noise reduction, *Physica D*, 47, 373-392, 1992.
- Harvey, A.C., *Time Series Models*, Wiley, New York, 1981.
- Houtekamer, P., *Predictability in models of the atmospheric circulation*, Ph.D.Thesis, Wagenigen Agricultural University, 1992.
- Jazwinski, A.H., *Stochastic Processes and Filtering Theory*, Academic Press, Paris, 1970.
- Lions, J.L., Sentinels and stealthy perturbations, *Proc. Int Symp. on Assimilation of Observations in Meteorology and Oceanography*, Clermont-Ferrand, World Meteorological Organization, 13-18, 1990.

- Lorenz, E.N., Deterministic nonperiodic flow, *J. Atmos. Sci.*, 20, 130-144, 1963.
- Machenhauer, B., *The spectral method*, Garp Publ. series no. 17 II, 124-277, 1979.
- Mehra, R.K., On the identification of variances and adaptive Kalman filtering, *IEEE Trans. Autom. Control*, AC-15, 2, 175-184, 1970.
- Moran, P.A.P., Some theorems on time series I, *Biometrika*, 34, 281-291, 1947.
- Mous, S.L.J. and J. Grasman, Two methods for assessing the size of external perturbations in chaotic processes, *Math. Meth. and Models in Appl. Sc.*, 3, 4, 577-593, 1993.
- Sparrow, C., *The lorenz equations: bifurcations, chaos, and strange attractors*, Springer Verlag, New York, 1982.
- Sorenson, H.W., Recursive estimation for nonlinear dynamical systems, in: *J.C. Spall (ed.), Bayesian analysis of time series and dynamic models*, Dekker, New York, 1988.

## SUMMARY

The development of accurate models is very important for analyzing problems concerning simulation, prediction, control, etc. Therefore it is not astonishing that many studies in applied science are about the modeling of these processes. In this thesis we will focus on the building of models that are used to describe some nonlinear processes in hydrology and meteorology; the first process is the movement of water in porous media and the second process is the large-scale atmospheric circulations.

The process of model development can be divided in three essential subprocesses: selection of a model structure, determination of a "best fit" criterion and experimental design. In literature, there are several examples of "case-studies" known, where the specific combination of model structure, criterion and experimental design did not lead to unique estimates of the unknown parameters of the model. This situation is designated by the term: "the model is not identifiable".

A model may not be identifiable (given a certain choice of the experimental design) because the chosen object function is insensitive to some linear combinations of the parameters. In this case the identification problem will not have a unique solution. On the other hand, due to noise in the

---

system, the optimization problem may have many local optima. One can then easily be misled because an optimization algorithm may converge to such a local optimum. It will be studied how such a situation can be recognized. Furthermore, it will be studied how the identifiability can be improved by an appropriate choice of the experimental design.

There are also other situations where the chosen combination of model structure, "best fit" criterion and experimental design will not lead to a unique solution. Such a case occurs when we are dealing with chaotic systems. For chaotic systems the optimization problem, using the output-error criterion as "best-fit" criterion, is ill-posed, because the model's solution depends sensitively on its initial state. The observed values and the model values will then diverge due to the limited accuracy of the initial state. Several criteria are analyzed on their capability for detecting small perturbations in the system and for estimating unknown parameters in the system.

In chapter 2 of this thesis the ONE-STEP method is described. This method is developed to identify the parameters in a model for the movement of water in the unsaturated soils. The motivation to analyze the identifiability of this model comes from the statement made by several authors that not all model parameters can be estimated uniquely. In this chapter we will analyze first some numerical schemes to solve the mathematical model, because the efficiency and the accuracy of a numerical scheme are very important for applicability of the ONE-STEP method.

In chapter 3 the concept of "structural identifiability" is further developed. The term "numerical identifiable" is introduced, so that we can take into account the accuracy of the sensitivity matrix. The identifiability analysis of the ONE-STEP method shows that not all parameters can be estimated uniquely. In the best case, where the pressure in the pressure cell is increased during the experiment at certain time instants, only 5 of the 6 model parameters can be estimated uniquely. Analyzing the structure of the model, we can derive that the object function depends on 5 independent parameters only,

which explains the identifiability problem. Only by adding some other measurements, for example the pressure head at a certain position in the soil core, one may expect better results of this method.

As already mentioned above, the output-error criterion in combination with chaotic systems, leads to ill-posed problems. In chapter 4 it is analyzed whether a criterion, based on a modified sentinel function, can be used to detect an external perturbation in a chaotic system. We found that fast varying perturbations are often "stealthy" for this function. Therefore this criterion can only be used to detect slowly varying perturbations.

The sentinel function can also be used to estimate uncertain parameters that are used to describe such a small perturbation term. We have compared the performance of the sentinel approach with an adaptive extended Kalman filter in a test-case. In the example that is presented, the size of a perturbation in the equator-pole temperature gradient is estimated. The equator-pole temperature gradient characterizes the driving force in a low-order spectral model of the atmospheric circulation and therefore a change in the equator-pole temperature gradient may be important in studying the greenhouse effect. In this test-case the performance of the adaptive extended Kalman filter was better than the performance of the sentinel approach. The less accurate results of the sentinel method are caused by the relative slow sampling frequency. The effect of neglecting higher order terms in the Taylor expansion and the influence of observation errors is then felt.

A disadvantage of extended Kalman filtering is that the filter easily diverges. In chapter 5 this problem is studied for chaotic systems. A reasonable approach to solve the divergence problem is to add an artificial noise term to the state equations. This noise term is used to control the accuracy of the state estimates and so preventing that the filter learns the wrong state too well. With the extended Kalman filter one can easily obtain an approximation of the value of the loglikelihood function. For this problem we have developed an optimization procedure that can be used together with the extended

---

Kalman filter to estimate the unknown parameters in the model description as well as the parameters that are used to describe the covariance matrix of the artificial noise term. This method is successfully applied to determine the optimal extended Kalman filter for a T11-spectral model of the atmospheric circulation.



## SAMENVATTING

### Over identificatie van niet-lineaire systemen

Het ontwikkelen van nauwkeurige modellen is erg belangrijk voor het analyseren van problemen betreffende simulatie, predictie, etc. Het is dus niet verwonderlijk dat veel toegepast onderzoek betrekking heeft op het ontwikkelen van modellen. In dit proefschrift spelen een tweetal onderzoeksgebieden een belangrijke rol: stroming van water in poreuze media en grootschalige atmosferische circulaties.

Bij het ontwikkelen van modellen kan men een drietal essentiële deelprocessen onderscheiden: selectie van een model structuur, bepalen van een "best fit" criterium en experiment-ontwerp. Er zijn echter diverse voorbeelden van "case-studies" waarbij de gekozen combinatie van model structuur, "best fit" criterium en experiment-ontwerp niet leidde tot een-eenduidige schattingen van de model parameters. Deze situatie wordt veelal aangeduid met "het model is niet identificeerbaar".

Een model is vaak niet-identificeerbaar omdat de "best fit" functie ongevoelig is voor bepaalde lineaire combinaties van de model parameters. Het identificatie probleem heeft dan geen een-eenduidige oplossing. Door ruis

---

in het systeem heeft het optimaliseringsprobleem wel vele lokale optima. Men kan dus eenvoudig misleid worden omdat gangbare optimaliseringsalgoritmes naar een lokaal optimum convergeren. Aan de hand van een voorbeeld wordt beschreven hoe de identificeerbaarheid van het model geanalyseerd kan worden. Tevens wordt onderzocht of de identificeerbaarheid verbeterd kan worden door een geschikte keuze van het experiment ontwerp.

Er zijn ook nog andere situaties waarbij de gekozen combinatie van modelstructuur, "best fit" criterium en experiment-ontwerp niet tot een-eenduidige oplossing leidt. Een dergelijke situatie doet zich voor bij het identificeren van chaotische systemen. Het optimaliseringsprobleem voor deze systemen, waarbij de kwadraatsom van de output errors als "best fit" functie gebruikt wordt, is slecht geconditioneerd omdat chaotische systemen gevoelig zijn voor kleine verstoringen in de beginvoorwaarden. Ten gevolge van de onnauwkeurigheid in de beginschattingen zullen dan observaties en model uitkomsten divergeren, waardoor het optimaliseringsprobleem vele lokale optima zal hebben. Verschillende criteria zijn onderzocht op hun geschiktheid om een kleine verstoring in het systeem te detecteren en om de onbekende parameters in het systeem te schatten.

In hoofdstuk 2 van dit proefschrift wordt de ONE-STEP methode beschreven. Deze methode is ontwikkeld om de parameters in een model voor stroming van water in de onverzadigde gronden te schatten. De motivatie om de identificeerbaarheid van dit model te analyseren kwam voort uit verschillende onderzoeken waaruit bleek dat de parameters in dit model niet uniek te schatten zijn. In dit hoofdstuk worden echter allereerst enkele numerieke schema's nader onderzocht omdat de efficiëntie en de nauwkeurigheid van deze schema's van groot belang zijn voor de praktische toepasbaarheid van de ONE-STEP methode.

In hoofdstuk 3 wordt het concept van "structural identifiability" verder ontwikkeld. De term "numeriek identificeerbaar" wordt ingevoerd, zodat ook rekening gehouden kan worden met numerieke onnauwkeurigheden in de

gevoeligheidsmatrix. Uit de identificeerbaarheidsanalyse van de ONE-STEP methode blijkt dat niet alle parameters uniek te schatten zijn. In het beste geval, waarbij de druk in de drukcel gedurende het experiment op gezette tijden wordt verhoogd, zijn slechts 5 van de 6 parameters uniek te schatten. De oorzaak van het niet-identificeerbaarheid is de structuur van het model. Alleen door het toevoegen van andersoortige metingen, bijv. de pressure head op een bepaalde plaats in het monster, kunnen we goede resultaten van deze methoden verwachten.

Zoals reeds boven vermeld leidt het "output error" criterium bij chaotische modellen tot een slecht geconditioneerd optimaliseringsprobleem. In hoofdstuk 4 wordt geanalyseerd of een criterium gebaseerd op de sentinel functie gebruikt kan worden om een perturbatie in een chaotisch systeem te detecteren. Het blijkt dat langzaam variërende perturbaties goed gedetecteerd kunnen worden maar dat snel variërende perturbaties vaak "stealthy/onzichtbaar" zijn voor de sentinel functie.

De sentinelfunctie kan ook gebruikt worden om onzekere parameters, die een perturbatie-term beschrijven, te schatten. In een test-case wordt deze methode vergeleken met een adaptief extended Kalman filter. In het gekozen voorbeeld wordt de grootte van een verstoring in de equator-pool temperatuur-gradiënt geschat in een lage-orde spectraal model van de atmosferische circulatie. De equator-pool temperatuur-gradiënt kan geïnterpreteerd worden als de aandrijvende kracht in het systeem. Een nauwkeurige schatting van deze temperatuur-gradiënt kan dus van belang zijn voor het onderzoek naar het broeikas effect. Uit de analyse blijkt de performance van de sentinel methode het af te leggen tegen een adaptive extended Kalman filter. De reden hiervoor is een lage sampling frequentie, waardoor de effecten van het verwaarlozen van hogere orde termen in de Taylor expansie en van de invloed van observatie fouten merkbaar worden.

Een nadeel van het extended Kalman filter is dat het eenvoudig kan gaan divergeren bij chaotische modellen. In hoofdstuk 5 is dit nader onder-

---

zoekt. Een redelijke aanpak om het probleem van divergentie op te lossen is een artificiële ruisterterm op te tellen bij de systeem-vergelijking. Deze ruisterterm kan gebruikt worden om de nauwkeurigheid van de schattingen te regelen, waardoor het filter niet meer divergeert. Met behulp van het extended Kalman filter kan een benadering van de waarde van de log-likelihood berekend worden, zodat met een speciaal ontwikkeld optimaliseringsprocedure zowel de onbekende model parameters alsook de onbekende parameters die de covariantie matrix van de artificiële ruisterterm beschrijven geschat kunnen worden. Deze methode is met succes toegepast om het optimale extended Kalman filter voor een T11-spectraal model voor de atmosferische circulatie te bepalen.

## Dankwoord

Dit proefschrift is het resultaat van vier jaar onderzoek, verricht bij de vakgroep Wiskunde en Toegepaste Statistiek van de Landbouwniversiteit Wageningen. Maarten de Gee en Albert Otten wil ik bedanken voor het initiëren van dit project en voor de prettige samenwerking gedurende deze vier jaar. Jullie stimulerende kritieken en jullie streven naar een nauwkeurige formulering hebben er toe geleid dat de eerste pennevruchten geworden zijn tot dit proefschrift. Ik wil ook Paul van der Laan bedanken die in eerste instantie zijn medewerking aan dit project heeft verleend.

Bij de aanvang van dit project zijn we druk op zoek geweest naar mogelijke toepassingsgebieden. Pieter Raats noemde de problematiek bij het identificeren van stromingsmodellen in de onverzadigde zone. Met Jos van Dam, die werkte aan dit probleem bij de vakgroep hydrologie, bodemnatuurkunde en hydraulica, heb ik veel stimulerende discussies over dit onderwerp gehad. Pieter en Jos, ik wil jullie bedanken voor de medewerking bij het tot stand komen van de eerste twee hoofdstukken van dit proefschrift.

Ik wil Johan Grasman bedanken voor zijn stimulerende begeleiding bij het tot stand komen van het tweede gedeelte van dit proefschrift. Chaotische dynamica was tot voor drie jaar een volledig onbekend gebied voor mij. Jouw

---

goede kennis van dit gebied en jouw kunde om deze kennis over te dragen hebben er toe bijgedragen dat chaotisch dynamica ook voor mij een fascinerende wereld is geworden.

Dieter Rasch en Eligius Hendrix, die mij geholpen hebben om mijn kennis te verbreden op het gebied van resp. niet-lineaire regressie en niet-lineaire optimaliseringstechnieken, wil ik bedanken voor de prettige samenwerking.

De zaalvoetbal ploeg wil ik bedanken voor de vele uurtjes afleiding. Dat het gezegde "Übung macht den Meister" opgaat blijkt uit de stijgende lijn die de prestaties van ons team vertoont. Ook de mensen die meegedaan hebben aan de Tour-Toto wil ik bedanken, met name omdat ze mij de laatste keer lieten winnen.

Met veel plezier denk ik terug aan de vele gesprekken en discussies die ik met mijn vroegere kamergenoot Johan Dourleijn heb gevoerd. Ondanks mijn gebrabbel in quasi fries en limburgs hebben wij ons gedurende deze vier jaren uitstekend kunnen verstaan.

Sipko Mous  
Wageningen, 11 april 1994

## Curriculum Vitae

Sipko Mous werd op 17 januari 1965 geboren in Heerenveen. In 1983 deed hij eindexamen Atheneum-B aan het Bischoppelijk College Schöndeln te Roermond en begon hij technische natuurkunde te studeren aan de technische universiteit van Eindhoven. Zijn afstudeeronderzoek heeft hij verricht bij de vakgroep Systeem- en Regeltechniek. In augustus 1988 slaagde hij voor het ingenieursexamen. Nadat hij enige maanden in militaire dienst is geweest begon hij in augustus 1989 als assistent in opleiding (AIO) bij de vakgroep Wiskunde van de Landbouwniversiteit Wageningen. Hier verrichte hij onderzoek naar identificatie problemen bij niet-lineaire systemen.

Sinds september 1993 is hij als project-medewerker in dienst bij de vakgroep Bodemkunde en Plantevoeding van de Landbouwniversiteit Wageningen.