

STOCHASTIC SYSTEMS WITH EMBEDDED PARAMETER VARIATIONS

APPLICATIONS TO DISTRICT HEATING

Henning T. Sogaard

**LYNGBY 1993
Ph.D. THESIS
NO. 66**

imsot

ISSN 0908-3456

Copyright © 1993 Henning Tangen Sjøgaard.

Typeset in L^AT_EX.

Trykt af  - DTH
Bogbinder Hans Meyer

Some of the work in this thesis has previously been published in:

Søgaard, H. T. and H. Madsen, 1989: Methods for Tracking Time Varying Delays in Dynamic Systems. Preprints No. 3/89, IMSOR, The Technical University of Denmark (presented at ISF 89 - The Ninth International Symposium on Forecasting, 1989).

Madsen, H., O. P. Palsson, K. Sejling and H. T. Søgaard, 1990: *Models and Methods for Optimization of District Heating Systems. Part I: Models and Identification Methods*. IMSOR, The Technical University of Denmark.

Søgaard, H. T. and H. Madsen, 1991: On-line Estimation of Time-Varying Delays in District Heating Systems. *Proceedings of the 1991 European Simulation Multiconference*, 619-624 (Research report No. 2/1991, IMSOR, The Technical University of Denmark).

Søgaard, H. T. and H. Madsen, 1992: A Grey Box Model of the Variations of Air Temperature. *Proceedings of the 5th International Meeting on Statistical Climatology*, 177-180.

Madsen, H., O. P. Palsson, K. Sejling and H. T. Søgaard, 1992: *Models and Methods for Optimization of District Heating Systems. Part II: Models and Control Methods*. IMSOR, The Technical University of Denmark.

Edlund, P. O. and H. T. Søgaard, 1993: Fixed versus Time-varying Transfer Functions for Modelling Business Cycles. *Journal of Forecasting*, **12**, 345-364.

Preface

This thesis has been prepared at the Institute of Mathematical Statistics and Operations Research (IMSOR), the Technical University of Denmark, as one of the requirements which should be fulfilled to earn the Ph.D. degree in Engineering.

The thesis deals with various aspects of embedded variation of the parameters in dynamic systems. This includes linear stochastic state-space models in continuous time. Most results have some relation to forecasting, control and optimization of district heating systems which is a field that is subject to extensive research at IMSOR.

My first acknowledgement is to my supervisors Assoc. Prof., Ph.D. Henrik Madsen and Assoc. Prof., Techn. Dr. Jan Holst (Lund Institute of Technology, Lund, Sweden) for their valuable suggestions, support and critical comments in the course of this project.

I also wish to thank Ph.D. Ken Sejling and M.Sc. Olafur Petur Palsson for numerous interesting and fruitful discussions on various statistical topics in connection with the project. Furthermore, I would like to thank them for comments and corrections to parts of the manuscript.

Several other of my colleagues at IMSOR have in one way or another contributed to the results of the project, and I wish to thank them for that.

I wish to thank stud. ling. merc. Anette Nørgaard Jappe for her linguistic corrections and suggestions to the manuscript.

Lot of the writing had to be done on overtime, and I am indebted to my wife Tina and our daughter Kirstine for their endless patience and helpfulness.

Lyngby, August 1993

Henning T. Sogaard
Henning T. Sogaard

Summary

This thesis is about statistical models and methods for stochastic systems with embedded parameter variations. This includes a discussion of dynamic state-space models in continuous time with a view to description, forecasting and control. The embedded parameters are the physical parameters being embedded in the system, and embedded parameter variations refer to variation in time of the embedded parameters.

Most of the treated statistical models and methods are to some extent motivated from a need for automatic tools for forecasting and control in connection with district heating systems. However, approaches which do not assume very specific technical or physical description of the considered system are emphasized. Thus the methods are applicable to a wide class of dynamic systems. Furthermore, models and algorithms which are operational and easy to implement are emphasized.

In Chapter 2 some methods for adaptive estimation of time-varying time-delays and dynamics are discussed. First three algorithms based upon recursive estimation with exponential forgetting and next two methods based upon explicit models of the embedded parameter variations are proposed. All the methods assume that the basic model

belongs to an ARIMAX model structure. The methods have been tested with simulated data, data from a district heating system in Ishøj and Swedish business cycle data. The simulated data represents a slowly varying system, and the variations of the time-delay are tracked with close accuracy. Concerning the district heating data, the results from two of the methods indicate that the time-delay is exposed to large diurnal variations. However, the presence of these large variations results in bad performance of a third methods. The results of the experiments with Swedish business cycle data show certain variations of the time-delays between various economic indicators. Furthermore, the experiments illustrate that the proposed methods can be applied in technically very different fields.

Chapter 3 is about modelling and forecasting of ambient air temperature. In the first part forecast procedures based on exponential smoothing are described and applied. The procedures are applied for prediction of the air temperature up to 24 hours ahead. Various aspects of Winters' seasonal forecast procedure are reviewed, and four alternative forecast procedures which can be considered as modified versions of Winters' forecast procedure are proposed. Two of the them are non-linear procedures which turn out to give very good prediction results. However, as expected a traditional ARIMA model provides closer prediction accuracy for short prediction horizons (here up to 9 hours ahead). Investigations of a method which combines predictions from two simple exponential smoothing procedures through an adaptively estimated regression model have also been done. This method gives even better results than all the previously mentioned methods for prediction horizons longer than 5 hours.

In the second part of Chapter 3 a local climatic system is modelled by means of linear state-space models in continuous time. Discrete time data is used and the embedded parameters of the continuous time model are estimated by a maximum likelihood method. The state-space models are obtained by approximating the distributed system

by heat capacities and thermal resistances. It is assumed that the primary input of the system is the net radiation. The deviations from the real system and the measurement errors are described through stochastic terms in the model. The data consists of hourly measurements of climatic variables from a period of almost eight years. The estimates of the physically interpretable parameters are discussed, and it is concluded that some of the estimates are in accordance with the expected values while others deviate significantly from the expected values.

Chapter 4 is about multi-step predictive control with special emphasis laid upon controllers for systems with embedded parameter variations. First the use of multi-step predictive control in connection with district heating systems is motivated. Next a weighted predictive control strategy is described. This strategy can be used in case the dynamic relationship between input and output is not sufficiently described for application of generalized predictive control (GPC). The strategy is implemented and applied in the district heating system in Esbjerg/Varde, and results from this implementation are discussed. After this the traditional GPC strategy which does not allow embedded parameter variations in the model is reviewed. Then an extended strategy which permits multi-step predictive control of systems with embedded parameter variations is proposed. This extended strategy implies a generalization of the loss function which defines the optimality of the control. Furthermore, it is possible to expose the control signal to general equality constraints. It is shown that the traditional GPC is a special case of the extended control strategy. Finally, simulation experiments with multi-step generalized predictive control are performed. The results show that the proposed strategy leads to significantly better control than minimal variance control.

Resumé (Summary in Danish)

Denne afhandling omhandler statistiske modeller og metoder for stokastiske systemer med indlejrede parametervariationer. Herunder diskuteres dynamiske tilstandsmodeller i kontinuert tid med henblik på beskrivelse, forudsigelser og regulering. De indlejrede parametre er de fysiske parametre, der er indlejret i det betragtede system, og ved indlejret parametervariation forstås tidsmæssig variation af de indlejrede parametre.

Den overvejende del af de behandlede statistiske modeller og metoder er til en vis grad motiveret ud fra et behov for automatiske prognose og styringsværktøjer i forbindelse med fjernvarmesystemer. Der er dog i behandlingen lagt vægt på at benytte fremgangsmåder, der ikke forudsætter en meget specifik teknisk eller fysisk konfiguration af det betragtede system. Metoderne vil således have generel anvendelighed for en bred klasse af dynamiske systemer. Der er ligeledes lagt vægt på, at de identificerede modeller og udviklede algoritmer er operationelle og enkle at implementere.

I kapitel 2 diskuteres nogle metoder til adaptiv estimation af tidsvarierende tidsforsinkelser og dynamik i dynamiske input-output systemer. Interessen samler sig dog specielt om tidsvarierende tidsforsinkelser. Først foreslås tre algoritmer, der bygger på rekursiv estimation

med eksponentiel glemsel, og dernæst to metoder, der er baseret på eksplicite modelbeskrivelser af de indlejrede parametervariationer. I alle tilfælde antages den overordnede model at have en ARMAX-modelstruktur. Metoderne afprøves på simulerede data, data fra fjernvarmesystemet i Ishøj samt på data fra svensk økonomi. For de simulerede data, som repræsenterer et langsomt varierende system, opnås gode resultater med at følge variationerne. Et par af metoderne afslører, at fjernvarmesystemet i Ishøj er meget tidsvarierende, hvilket giver sig tydeligst til kende i en stærk døgnmæssig variation af tidsforsinkelsen. De store variationer betyder dog, at en tredje af metoderne giver dårlige resultater. Afprøvning af én af metoderne på de svenske økonomidata viser, at svensk økonomi repræsenterer et tidsvarierende dynamisk system med en vis variation af tidsforsinkelserne mellem forskellige økonomiske indikatorer. Afprøvningen på disse data illustrerer desuden, at de foreslåede metoder kan finde anvendelse på fagligt vidt forskellige områder.

Kapitel 3 omhandler modellering og prognostisering af udendørs lufttemperatur. I første del af kapitlet beskrives og anvendes prædiktionsprocedurer baseret på eksponentiel udjævning. Procedurerne anvendes til at prædiktere lufttemperaturen op til 24 timer frem. Der gennemgås aspekter af Winters' sæsonmæssige prædiktionsprocedure, og lignende alternative procedurer foreslås – bl.a. to ulineære procedurer, som viser sig at give gode prædiktionsresultater. Som forventet giver en traditionel ARIMA model dog nøjagtigere prædiktioner for de korte prædiktionshorisonter (her op til 9 timer frem). Der er yderligere foretaget undersøgelse af en alternativ metode, der kombinerer prædiktionerne fra to simple eksponentielle udjævningsmetoder igennem en adaptivt estimeret regressionsmodel. Denne metode giver endnu bedre prædiktioner end alle de tidligere nævnte metoder for horisonter over 5 timer.

I den anden del af kapitel 3 modelleres et lokalt klimasystem ved hjælp af lineære tilstandsmodeller i kontinuert tid. Ved at tilnærme

systemet med et mindre antal varmekapaciteter, der er indbyrdes forbundet af termiske modstande, opnås tilstandsmodellerne via en opstilling af ordinære lineære differentiaalligninger. Det antages, at det drivende input for systemet hovedsageligt er nettostrålingen. Modelfejl og målefejl beskrives ved stokastiske led i modellen, og parametrene estimeres ved maximum likelihood estimation. De benyttede data er timevise målinger af klimavariabel fra en næsten otte år lang periode. Estimererne af de fysisk fortolkelige parametre diskuteres, og det konkluderes at nogle stemmer overens med det forventede, mens der for andre er betydelige afvigelser fra det forventede.

Kapitel 4 handler om flertrins prædiktiv regulering, og der er specielt lagt vægt på regulatorer for systemer med indlejret parametervariation. Der gives først en motivation for benyttelse af flertrins prædiktiv regulering i forbindelse med fjernvarmesystemer. Dernæst foreslås en flertrins prædiktiv reguleringsstrategi, der kan benyttes, hvis den dynamiske sammenhæng mellem input og output er for utilstrækkeligt beskrevet til en anvendelse af generaliseret prædiktiv regulering (GPC). Strategien, som kaldes vægtet prædiktiv regulering, er implementeret og anvendes i fjernvarmesystemet i Esbjerg/Varde. Resultater herfra vises. Af "rigtige" GPC strategier gennemgås først den traditionelle, som er for systemer uden indlejret parametervariation. Herefter foreslås en udvidet strategi, der bl.a. muliggør flertrins prædiktiv regulering af systemer med indlejret parametervariation. Desuden medfører den udvidede strategi en generalisering af den tabsfunktion, der minimeres ved reguleringen, samt en mulighed for at underkaste styresignalet nogle generelle lineære lighedsrestriktioner. Det vises, at traditionel GPC er et specialtilfælde af den udvidede strategi. Til slut udføres simulationseksperimenter med flertrins generaliseret prædiktiv regulering. Resultaterne heraf viser, at den foreslåede strategi fører til en betydelig bedre regulering end minimalvariansregulering.

Contents

Preface	v
Summary	vii
Resumé (Summary in Danish)	xi
Contents	xv
1 Introduction	1
1.1 Minimization of Supply Temperature in District Heating Systems	3
1.2 Outline of the Thesis	5
2 Time-Varying Time-Delay in Dynamic Systems	11
2.1 Transfer Function Models for District Heating Systems	12
2.2 Methods for Tracking Time-Delay	14
2.2.1 Forgetting Factor Methods	15
Formulation and Recursive Estimation of the ARMAX Model	16
Estimating Models in Parallel	20
Bányász and Keviczky’s Method	23
A Continuous Time-Delay Approach	27
2.2.2 Modelling the Variations of the Time-Delay . .	32
Deterministic Variations	33

	Stochastic State-Space Model	40
2.3	Experimental Results	44
2.3.1	Simulation Results	45
	Tracking a Pure Time-Delay	46
	Tracking the Time-Delay of an ARX Model	50
2.3.2	Tracking a Time-Varying Time-Delay in a District Heating System	54
	Estimating Models in Parallel	58
	Continuous Time-Delay Approach	69
	A Stochastic State-Space Model	71
2.3.3	Time-Delays in Models for Business Cycle Forecasting	78
	The Business Cycle Data Set	79
	Time-Invariant Models	82
	Time-Varying Models	86
	Comparison of Time-Invariant and Time-Varying Models	90
2.4	Conclusion	95
3	Dynamic Models for Air Temperature	103
3.1	Exponential Smoothing	106
3.1.1	Winters' Seasonal Forecast Procedure	106
3.1.2	Alternative Procedures for Exponential Smoothing	121
	Equivalent Embedded ARIMA Models	126
	Choice of the Smoothing Constants	127
3.1.3	Prediction of Air Temperature – Results and Discussion	127
	Combining Forecasts	147
3.1.4	Conclusion	150
3.2	Models Relating the Variation in the Air Temperature to Other Climate Variables	152
3.2.1	Formulation of Models	153

	Modelling Heat Fluxes Using the Resistance Method	154
	Dynamic Relation between Air Temperature and Net Radiation	157
	A Simple Linear, Dynamic Model for the Air Temperature	160
	Linear Stochastic Models in Continuous Time .	162
	Estimation of the Parameters	171
3.2.2	Experimental Results	174
	The Data Set	174
	Estimation Results	176
3.2.3	Discussion and Conclusion	179
4	Multi-Step and Embedded Model Based Control	189
4.1	Control Strategy for a District Heating System	191
4.2	Weighted Predictive Control	196
4.2.1	A Simple Example	197
4.2.2	The General Case	198
4.2.3	Weighted Predictive Control of Supply Tem- perature	200
4.2.4	Results Obtained at Vestkraft in Esbjerg . . .	207
4.3	Ordinary Generalized Predictive Control	211
4.3.1	The Model and the Diophantine Equation . . .	214
4.3.2	The GPC Cost Function	216
4.3.3	Optimization of the GPC Cost Function	220
4.4	Generalized Predictive Control for Embedded Models	224
4.4.1	Model Structure and Output Prediction	224
	The ARMAX Model	224
	The j -Step-Ahead Predictor	225
	The Predictor as a Linear Function of Future Controls	226
4.4.2	The Cost Function	232

Choice of the Design Parameters of the Con-	
troller	233
4.4.3 The Resulting Controller	236
4.4.4 Linear Equality Constraints	238
The Constraints used by Clarke <i>et al.</i>	240
4.5 Simulation Experiments with XGPC	242
A Simple Model of the Mass Flow	243
Reference Values	245
Simulation of Suboptimal Minimum Variance	
Control	246
Simulation of XGPC	249
4.6 Conclusion	255
5 Conclusions	259
A The Kalman Filter in Discrete Time	265
B The Computer Program PRESS	269
B.1 PRESS	269
B.2 Results from Using PRESS in Esbjerg	272
References	275
IMSOR Ph.D. Theses	281

Chapter 1

Introduction

THE modelling of dynamic processes is frequently based upon the assumption that time-invariant linear models constitute appropriate approximations of the real processes. In practice most dynamic processes are neither linear nor time-invariant, and depending on the applications it is then often necessary to take the time-variation and the non-linearities into account. This thesis mainly deals with the aspect of time-variation and embedded parameter variation, but through this study, non-linear aspects are also encountered.

Most of the presented models and methods are in some way related to control and prediction of quantities as temperature and heat load in district heating systems. The diurnal and annual variations of the heat consumption in the district heating system induce corresponding diurnal and annual variations of the dynamic characteristics of the system. Especially the transport time from the heat producing unit(s) to the consumers is affected by the variations of the consumption. Approaches for both identification of time-varying models and control of systems with time-varying parameters are proposed. Some

of the methods consider explicit embedded models of the parameters while others are modelling the embedded variation implicitly through algorithms that tracks the variations.

A substantial part of the diurnal and annual variations of the heat consumption is a result of variations of climatic quantities such as ambient air temperature and incident solar radiation. To perform on-line control of the heat production it is essential to get forecasts of the heat demand since there is a delay between time of production and time of consumption, and due to the relation between the heat demand and the climate, there is a need for forecasts of climatic variables as well. The weather forecasts which can be obtained from the official meteorological services are not sufficiently detailed to be used in a limited geographical area which is supplied with district heating. Therefore, various models for short term forecasting of the local ambient air temperature are considered.

The models and methods are based upon statistical approaches combined with knowledge of the physical systems. This means that information from time series data and knowledge of the physical systems combined to specify suitable model structures. The values of the model parameters are estimated by suitable statistical methods such as the maximum likelihood method, the prediction error method and the least squares method.

An important quality of the statistical models and methods presented in this thesis is that they are or can be formulated as recursive algorithms which can be used for on-line applications. With regular time-intervals the recursive algorithms use new measurements to update a number of state variables which hold aggregate information extracted from earlier measurements. The state variables are typically estimates of model parameters and related variables which describe the current dynamic characteristics of the system. In order to track changes in the dynamics of the real system, some of the al-

gorithms are made adaptive. This means that the largest weight is attributed to recently recorded measurements while the influence of older measurements decrease as time elapses.

Although the identification of model structures and development algorithms have been carried out with a view to district heating systems and climatic systems, the resulting methods have relevance to several other dynamic systems. No specific technical configuration of the dynamic system is assumed.

1.1 Minimization of Supply Temperature in District Heating Systems

For many district heating systems, e.g. the system in Esbjerg/Varde, it is desirable to minimize the supply temperature from the heat production unit(s). Lower supply temperature implies lower costs in connection with the production and distribution of heat. For the combined heat and power plant, Vestkraft, in Esbjerg it has been estimated that decreasing the supply temperature by 1 °C implies a reduction of the production costs by about 1 % (for temperatures between 80 °C and 95 °C). The reduction of heat loss from the distribution network is about 0.5 %/°C.

Traditionally the supply temperature from the district heating plant is controlled as a function of the current ambient air temperature. The relation between the ambient air temperature and the supply temperature reflects that the supply temperature must be increased with decreasing ambient air temperatures due to increasing heat demand. A relation of this kind is illustrated by the control curve in Figure 1.1. This curve has previously been used at Vestkraft in Esbjerg as a minimum curve. The upper and lower limits of the supply

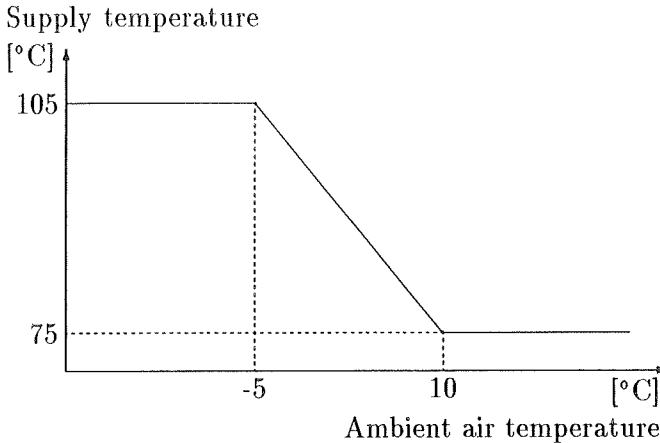


Figure 1.1: *Control curve for supply temperature which has previously been used at Vestkraft in Esbjerg.*

temperature (105 and 75 °C respectively) are primarily due to the technical design of the production system. The lower limit also ensures that the temperature of the hot tap water at the consumers is sufficiently high.

Letting the supply temperature depend on the current ambient air temperature is obviously more optimal than keeping it constant at its upper limit. However, provided that a control curve like the one in Figure 1.1 actually ensures that the heat supply is sufficient at any time, it will often lead to temperatures that are unnecessarily high. This is because the heat demand does not only depend on the current ambient air temperature. Other factors having impact on the heat demand are for instance solar radiation, wind speed, wind direction and a climate independent part which is a function of the time of the day/week/year. Furthermore the various factors including the ambient air temperature affect the heat demand in an indirect manner due to the dynamics of the distribution network and the supplied buildings which act as heat reservoirs. By taking such aspects into

consideration it is possible to obtain a more optimal control of the supply temperature, i.e. generally lower supply temperatures.

The models and methods described in the thesis take such aspects into account, and can therefore be used as elements in a more efficient minimization of the supply temperature. Some of them have been implemented in the computer program PRESS which is a tool for optimal control of supply temperature and forecasting of heat demand in district heating systems. A brief description of the program is given in Appendix B, and Figure 1.2 shows a simplified diagram of the data flow and the data processing taking place in the program. How the models and methods are related to the boxes A, B and C in the diagram is mentioned in the outline of the thesis in the following section.

1.2 Outline of the Thesis

Chapter 2 deals mainly with time-varying time-delays in dynamic input-output systems, but the variation of the ordinary dynamics is also considered. Various methods for tracking variations of the delay and the dynamics are proposed, and both simulated and real data are used to test the methods. An ARMAX model structure is assumed in connection with all the methods.

Two types of tracking methods are considered: forgetting factor based methods and methods based upon explicit embedded models of the time-varying parameters. Some of the methods are tested on simulated data generated by a model with time-invariant dynamics and a slowly varying delay. It appears that the district heating data which is used for some of the tests represents a system with very quick time-variation, and the time-delay shows a predominating diurnal variation. As a final case study, business cycle data is used to test one

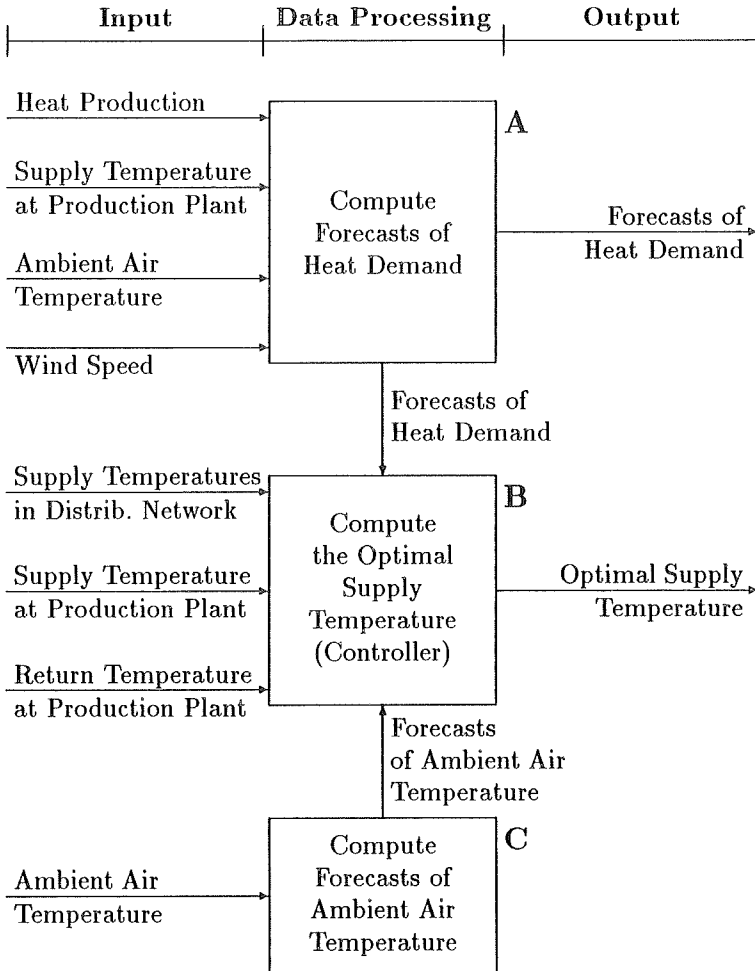


Figure 1.2: *Simplified diagram of the data flow and data processing in the computer program PRESS.*

of the methods. This case study clearly indicates that methods for tracking time-variation are also relevant to dynamic systems which are not of a technical nature.

A combination of two of the proposed methods has been implemented in PRESS. In the diagram in Figure 1.2 these methods are integrated parts of box B. Time-varying transfer functions describing the relationship between the supply temperature at the combined heat and power plant and the supply temperatures in the distribution network are estimated adaptively. For control purposes it is often required that the current estimates of the time-delays should be very accurate.

Chapter 3 describes dynamic models of the variations of the local ambient air temperature.

The first half part of the chapter deals with exponential smoothing approaches for prediction of the temperature up to 24 hours ahead. These approaches are based on empirical descriptions of the variations of the air temperature using components as level, trend and cyclic diurnal variation, and they do not consider the physical mechanisms of the underlying climatic system.

At first, some features of Winters' seasonal forecast procedure are studied, and it is verified that the resulting predictor is equivalent to an ARIMA predictor and that the time-variation of the components in Winters' model corresponds to embedded ARIMA models.

Then four alternative forecast procedures are proposed. Two of them are linear like Winters' procedure and two of them are non-linear. By testing the procedures on real air temperature data, it is shown that the non-linear procedures give most accurate forecasts. As expected a traditional ARIMA model, which is fitted to the data, is supe-

rior to the non-linear procedures as regards short term forecasting, but becomes inferior if the forecast horizon is prolonged. Finally a method is proposed that combines forecasts from two simple exponential smoothing procedures by means of an adaptive regression model. For forecast horizons longer than 5 hours this approach gives better results than all the other approaches.

One of the linear smoothing procedures is implemented in PRESS (cf. box C in Figure 1.2).

The second part of Chapter 3 relates the variations of the temperature to the net radiation and other input variables by means of stochastic state-space models in continuous time which are motivated from a physical viewpoint. By viewing the local climatic system as consisting of heat reservoirs connected by thermal resistances, the system can be described by linear differential equations with the net radiation as the main input. The model error is described by introducing stochastic terms in the resulting state-space model. Hourly climate observations from a period of almost eight years are used to obtain maximum likelihood estimates of the model parameters, and it is concluded that the most advanced models, which include a description of the energy balance at the surface of the ground, give the best results.

Chapter 4. The theme in this chapter is multi-step predictive control with a special emphasis on predictive control permitting embedded time-variation of the model parameters.

At first a control strategy that is called weighted predictive control (WPC) is proposed. In cases where the output process of a dynamic input-output system has been described by a stochastic model, which does not specify the relationship between input and output sufficiently well, use of the real generalized predictive control method is not feasible. In such situations, the WPC method is relevant, since

it lets the user specify the control weights which are found by optimization in the real generalized predictive control method. A WPC strategy for the district heating system in Esbjerg/Varde is proposed. This strategy is implemented in PRESS (as box B in Figure 1.2), and some results achieved by using the strategy in Esbjerg are shown.

The ordinary generalized predictive control (OGPC) method proposed by Clarke *et al.* (1987A) is reviewed, and several extensions of this method are suggested. First of all, contrary to the OGPC method, the extended method (XGPC) can be used to control systems described by models with embedded models of the parameter variation. Another extension concerns the cost function which is minimized by the controller. Finally the XGPC method permits the projected control signal to be subject to restrictions, which can be described by a system of linear equations.

By simulation experiments it is verified that, for a non-minimum phase system, the XGPC method is able to keep the output closer to the reference signal with less control effort than a suboptimal minimal variance controller.

Chapter 2

Time-Varying Time-Delay in Dynamic Systems

IN many dynamic input-output systems, transport of mass takes place. This very often implies a time-delay from input to output. The work in this chapter has been inspired by time-delays in district heating systems where transport of water takes place. Therefore a brief description of the use of time-varying, stochastic transfer function models of the supply temperature in district heating systems is given in Section 2.1 (cf. Madsen *et al.* (1990)).

The models and methods that have mainly been developed and applied for district heating systems are also applicable for other systems exhibiting time-varying delay and dynamics. In economic theory, business cycle mechanisms constitute such systems. As a case study, one of the adaptive methods is tested using business cycle data. This is done in Section 2.3.

This chapter is divided into three main parts:

Section 2.1: A brief description of the use of time-varying transfer function models in district heating systems.

Section 2.2: The theory of some techniques for tracking time-varying delays is reviewed. The first part deals with methods based on forgetting factor estimation, while embedded models of the variation of time-delays are described in the last part.

Section 2.3: The main purpose of this section is to describe experimental studies of the various models and methods. At first, some results of simulation experiments are given. Then data from a district heating system is used to test the methods, and finally time-delays in business cycle models are tracked by use of financial data series.

2.1 Transfer Function Models for District Heating Systems

Transfer functions from the supply temperature at the district heating plant to the supply temperature in individual supply pipes of the distribution network are of interest and usable for on-line applications as forecasting and control. If, for instance, minimum supply temperature is required for a point in the network, one solution might be to identify a transfer function model for that point and design a controller for the supply temperature at the district heating plant in order to meet the requirements.

At the combined heat and power (CHP) plant, Vestkraft in Esbjerg, Denmark, adaptive estimation of transfer function models has been implemented to control the supply temperature. The purpose of the

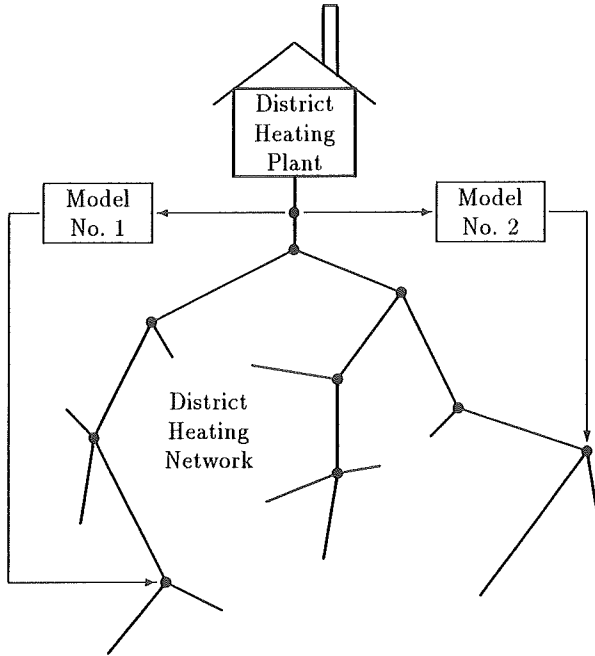


Figure 2.1: A district heating system with transfer function models for two points in the distribution network.

controller is to ensure that specified minimum temperature requirements at the service pipes are met. This is done by supervising the supply temperatures at a number of representative points in the distribution network. These points should be representative in the sense that if the temperature there is acceptable at any time then all consumers are supplied with sufficiently high temperature. Figure 2.1 illustrates a system (not the system in Esbjerg) with two such points. For those points, minimal variance controllers based on adaptively estimated transfer function models have been developed. The supply temperature at the CHP plant is controlled so that the point currently requiring the highest supply temperature from the CHP

plant is met. To make this control strategy work, it is important to know the transport times (delays) from the CHP plant to each of the points. The presence of time-variations of the time-delays is a major problem for automatic control of district heating systems. If the current delays are not known very well, the controller at the CHP plant may not react in time. Therefore, this chapter mainly deals with the description of methods for tracking time-varying delays.

The estimation methods might also be employed as a part of a large-scale on-line optimization strategy for a district heating system. Compared to models for the entire distribution network (e.g. the pipe element model and the pipe node model described by Benonysson (1991)), the methods used here are rather inexpensive as regards computations and data requirements. Therefore, they are well suited for on-line applications.

2.2 Methods for Tracking Time-Delay

This section describes various methods for tracking a time-varying time-delay in a dynamic input-output system. The section is subdivided in two parts: The first part presents algorithms based on forgetting factor estimation of the model parameters, while the second part presents methods relying on an explicit model of the embedded variation of the delay.

2.2.1 Forgetting Factor Methods

The idea of forgetting factor estimation¹ is based upon the exponential smoothing approaches by Winters (1960) and Brown (1963). After that time a great deal of literature on forgetting factor estimation have been published – e.g., Abraham and Ledolter (1983), Ljung and Söderström (1983), Young (1984) and Ljung (1987). One of the reasons why the forgetting factor techniques have been used so widely for tracking time-varying systems is that they do not rely on any specific knowledge of the type of the time-variation. On the other hand, it is well-known that fast changes of the system or a low signal-to-noise ratio may reduce the estimation performance significantly. In this presentation, however, it is assumed that the time-variation and the signal-to-noise ratio are of a kind that makes forgetting factor methods applicable.

Several authors have proposed methods for recursive estimation of an unknown and time-varying time-delay in dynamic systems. Most of these methods are based on the RLS algorithm. Bányász and Keviczky (1988) considered the discrete delay to be a continuous parameter. As the delay is included non-linearly in the ARMAX model, Bányász and Keviczky developed a RLS algorithm by using a Gauss-Newton approximation. Kurz and Goedecke (1981) presented a method which implies that a model with more parameters than the real system is estimated recursively. Elnaggar, Dumont and Elshafei (1989) presented a method which, to some extent, is similar to but not as efficient as the “parallel model estimation” method presented later in this section (and by Sjøgaard and Madsen (1989)). Chen and Zhang (1990) proposed a method for estimating the parameters, the orders and the time-delay of an ARX model.

¹In this chapter the terms “recursive estimation” and “forgetting factor estimation” are used interchangeably.

Formulation and Recursive Estimation of the ARMAX Model

Before the algorithms for tracking a time-varying delay are presented, the model structure and the recursive estimation scheme making the basis for those algorithms are reviewed.

The transfer function models considered belong to the ARMAX(n, m, r) model structure. The model can be expressed either as a difference equation,

$$\begin{aligned} y(t) + a_1 y(t-1) + \dots + a_n y(t-n) = \\ b_0 u(t-k) + b_1 u(t-k-1) + \dots + b_m u(t-k-m) \\ + e(t) + c_1 e(t-1) + \dots + c_r e(t-r), \end{aligned}$$

or by means of polynomials,

$$A(q^{-1})y(t) = B(q^{-1})q^{-k}u(t) + C(q^{-1})e(t), \quad (2.1)$$

where $\{y(t)\}$ and $\{u(t)\}$ are the output and the input signals, respectively (e.g. $u(t)$ = supply temperature at the district heating plant, $y(t)$ = supply temperature at a point in the distribution network). $\{e(t)\}$ is a Gaussian white noise process with the variance σ_e^2 (sequence of independent normally distributed stochastic variables with mean zero). The integer, k , denotes a time-delay (dead time) from input to output, and q^{-1} is the back shift operator defined by: $q^{-1}x_t = x_{t-1}$. $A(q^{-1})$, $B(q^{-1})$ and $C(q^{-1})$ are polynomials in q^{-1} :

$$\begin{aligned} A(q^{-1}) &= 1 + a_1 q^{-1} + \dots + a_n q^{-n} \\ B(q^{-1}) &= b_0 + b_1 q^{-1} + \dots + b_m q^{-m} \\ C(q^{-1}) &= 1 + c_1 q^{-1} + \dots + c_r q^{-r}. \end{aligned}$$

a_i ($i = 1, \dots, n$), b_j ($j = 0, \dots, m$) and c_l ($l = 1, \dots, r$) are the parameters of the model, and $b_0 \neq 0$.

Assume that the orders n , m and r of the polynomials and the time-delay k are known and constant, and that the parameters a_1, \dots, a_n , b_0, \dots, b_m and c_1, \dots, c_r are unknown and time-varying. Then the parameters can be estimated adaptively using a recursive forgetting factor method.

Here the recursive extended least squares (ELS) method is used (see Ljung (1987)):

$$\hat{\boldsymbol{\theta}}(t) = \hat{\boldsymbol{\theta}}(t-1) + \mathbf{L}(t)\boldsymbol{\varepsilon}(t) \quad (2.2)$$

$$\mathbf{L}(t) = \frac{\mathbf{P}(t-1)\boldsymbol{\varphi}(t)}{\lambda_\theta + \boldsymbol{\varphi}^T(t)\mathbf{P}(t-1)\boldsymbol{\varphi}(t)} = \mathbf{P}(t)\boldsymbol{\varphi}(t) \quad (2.3)$$

$$\mathbf{P}(t) = \left(\mathbf{P}(t-1) - \frac{\mathbf{P}(t-1)\boldsymbol{\varphi}(t)\boldsymbol{\varphi}^T(t)\mathbf{P}(t-1)}{\lambda_\theta + \boldsymbol{\varphi}^T(t)\mathbf{P}(t-1)\boldsymbol{\varphi}(t)} \right) / \lambda_\theta, \quad (2.4)$$

where

$$\hat{\boldsymbol{\theta}}(t) = (\hat{a}_1(t) \dots \hat{a}_n(t) \hat{b}_0(t) \dots \hat{b}_m(t) \hat{c}_1(t) \dots \hat{c}_r(t))^T,$$

is the estimated parameter vector at time t and

$$\boldsymbol{\varphi}(t) = (-y(t-1) \dots -y(t-n) \ u(t-k) \dots u(t-k-m) \ \bar{\varepsilon}(t-1) \dots \bar{\varepsilon}(t-r))^T.$$

The residual, $\bar{\varepsilon}(t)$, and the one-step-ahead prediction error, $\boldsymbol{\varepsilon}(t)$, are computed as

$$\begin{aligned} \bar{\varepsilon}(t) &= y(t) - \boldsymbol{\varphi}^T(t)\hat{\boldsymbol{\theta}}(t) \\ \boldsymbol{\varepsilon}(t) &= y(t) - \boldsymbol{\varphi}^T(t)\hat{\boldsymbol{\theta}}(t-1). \end{aligned} \quad (2.5)$$

In the ELS algorithm, λ_θ denotes the forgetting factor ($0 < \lambda_\theta \leq 1$), and $\sigma_e^2 \mathbf{P}(t)$ estimates the covariance matrix of the parameter estimate $\hat{\boldsymbol{\theta}}(t)$.

A method for recursive estimation of the variance σ_e^2 based on exponential smoothing is (Søgaard and Madsen (1989))

$$\hat{\sigma}_e^2(t) = \lambda_e \hat{\sigma}_e^2(t-1) + (1 - \lambda_e) \boldsymbol{\varepsilon}^2(t), \quad (2.6)$$

where λ_e is a forgetting factor. Since the prediction error in (2.5) is computed from the estimated and not the real parameter values, there is a tendency that the estimate of σ_e^2 is larger than the real value (the estimate is biased). Therefore $\hat{\sigma}_e^2(t)$ is an estimate of the prediction error variance rather than of the noise variance.

Poulsen (1985) and Poulsen and Holst (1988) present an alternative way of estimating the noise variance. The parameter variation is described by a linear state space model and the noise variance is assumed to be distributed as a reciprocal gamma distribution. Estimates of the parameters and the noise variance are obtained as conditional expectations (conditioned on the observations). It turns out that the optimal estimation consists of a Kalman filter supplied with a superstructure which performs the estimation of the noise variance.

In the sections below, the ELS algorithm (2.2)-(2.4) will be written in a concise way as follows:

$$(\hat{\theta}(t), P(t)) = \text{ELS}(\hat{\theta}(t-1), P(t-1), \varphi(t), y(t), \lambda_\theta). \quad (2.7)$$

Given constant system parameters, Ljung (1987) showed that the convergence properties of the ELS approach depend on the C -polynomial of the real system description. See also Holst (1977) for a discussion of the convergence properties of the ELS algorithm. In the case of $C(q^{-1}) = 1$, the ELS algorithm is reduced to the recursive least squares (RLS) algorithm which will converge to the real parameter values for $\lambda_\theta = 1$. This is due to the fact that the RLS algorithm is based upon a model which is linear in the parameters whereas the ELS algorithm is based upon a pseudo linear model. Note that only the $\hat{\theta}(t)$ vector and the $\varphi(t)$ vector need to be changed to obtain the RLS algorithm. Equations (2.2)-(2.4) remain unchanged.

The ELS algorithm requires initial values of $\hat{\theta}(0)$ and $P(0)$. If no prior information is available, it seems reasonable to choose $\hat{\theta}(0) =$

$(0 \dots 0)^T$ and $P(0) = KI$ where K is a large, positive number and I denotes an identity matrix of a suitable dimension. According to a Bayesian interpretation this corresponds to assuming a diffuse *a priori* distribution of the parameter values. In (2.6) an initial value for $\hat{\sigma}_e^2(0)$ is required. However, the influence of the initializations will decrease as time elapses. See Ljung (1987) for more information on initialization.

The choice of the forgetting factors, λ_θ and λ_e , should reflect the rapidity of the changes of the real system parameters and the noise variance, respectively. In general, a system which is subject to quick dynamic changes requires low forgetting factor values while forgetting factors close to one are preferable for systems which are almost stationary. The noise variance, however, should also be considered when forgetting factor values are chosen. A low noise variance (a high signal-to-noise ratio) permits low forgetting factor values and vice versa. For practical purposes, λ_θ is often determined as the value implying minimum prediction error variance for a given sample series of input and output.

Now leave the assumption that the time-delay is known and constant. Instead suppose that both the delay, k , and the parameters, a_1, \dots, a_n , b_0, \dots, b_m and c_1, \dots, c_r , are unknown and time-varying though it is still assumed that the orders, n , m and r , of the ARMAX model (2.1) are known and constant. Three methods for simultaneous recursive estimation of time-delay and parameters are described below. The methods were also presented by Sjøgaard and Madsen (1989) and Madsen *et al.* (1990). The first method is based upon a procedure for switching between models in a collection of models which is estimated in parallel, the second upon the general recursive Gauss-Newton method, and the third method considers the embedded continuous time-delay.

Estimating Models in Parallel

Assume that the time-delay obeys the restriction $k_{min} \leq k \leq k_{max}$ ($k_{min} \geq 0$). Then, at time t , the dynamic system is described by one of the models

$$A_p(q^{-1})y(t) = B_p(q^{-1})q^{-p}u(t) + C_p(q^{-1})e(t), \quad p = k_{min}, \dots, k_{max}, \quad (2.8)$$

where p is a delay-index. In order to estimate which of the models that describes the real system at time t , a parallel recursive estimation of the models is made using the ELS algorithm. That is, estimates of the parameter vectors, $\theta_{k_{min}}, \dots, \theta_{k_{max}}$, of the models are updated separately in each step. By means of an appropriate criterion, the model giving “the best” description of the system is currently selected.

The model formulation in (2.8) indicates that the system dynamics is assumed to change with the time-delay (cf. the subscript p of the polynomials A_p , B_p and C_p). Alternatively, the dynamics could be assumed independent of the time-delay.

The selection criterion employed here is based upon the recursively estimated noise variance. For each of the $k_{max} - k_{min} + 1$ models, the noise variances, $\sigma_{e, k_{min}}^2, \dots, \sigma_{e, k_{max}}^2$, are estimated recursively according to Equation (2.6):

$$\hat{\sigma}_{e,p}^2(t) = \lambda_e \hat{\sigma}_{e,p}^2(t-1) + (1 - \lambda_e) \varepsilon_p^2(t), \quad p = k_{min}, \dots, k_{max},$$

where $\varepsilon_p(t)$ is the one-step-ahead prediction error of model p at time t . $\hat{\sigma}_{e,p}^2(t)$ may be interpreted as the current prediction ability of model p . At each step of the algorithm the lowest $\hat{\sigma}_{e,p}^2(t)$ -value is found. The corresponding model states the current estimates of k and θ .

In terms of the notation introduced in Equation (2.7), the entire

algorithm is

$$(\hat{\theta}_p(t), P_p(t)) = \text{ELS}(\hat{\theta}_p(t-1), P_p(t-1), \varphi_p(t), y(t), \lambda_\theta),$$

$$p = k_{\min}, \dots, k_{\max} \quad (2.9)$$

$$\hat{\sigma}_{e,p}^2(t) = \lambda_e \hat{\sigma}_{e,p}^2(t-1) + (1 - \lambda_e) \varepsilon_p^2(t),$$

$$p = k_{\min}, \dots, k_{\max} \quad (2.10)$$

$$\hat{k}(t) = \arg \min_{p \in [k_{\min}, k_{\max}]} \left\{ \hat{\sigma}_{e,p}^2(t) \right\} \quad (2.11)$$

$$\hat{\theta}(t) = \hat{\theta}_{\hat{k}(t)}(t). \quad (2.12)$$

For some systems the algorithm can be improved by the restriction $|\hat{k}(t) - \hat{k}(t-1)| \leq \Delta k_{\max}$ where $\Delta k_{\max} \geq 0$. Thus (2.11) can be replaced by

$$\hat{k}(t) = \arg \min_{p \in [k_{\min}, k_{\max}] \cap [|\hat{k}(t-1) - \Delta k_{\max}, \hat{k}(t-1) + \Delta k_{\max}|]} \left\{ \hat{\sigma}_{e,p}^2(t) \right\}. \quad (2.13)$$

The restriction is especially convenient for systems influenced by heavy noise. It ensures that the change of $\hat{k}(t)$ from one step to another remains within a reasonable range compared to the changes of the real delay, k . If this restriction was not present, large jumps of $\hat{k}(t)$ that could not occur in connection with the real delay might be observed.

If the number of models estimated in parallel is large, the computational effort at each step may prevent the method from being employed in real-time applications. But, in case no abrupt or total changes of the time-delay occur, it is reasonable to reduce the number of potential models that are being estimated, and hence reduce the computational work. One way to do this is to make k_{\min} and

k_{max} time-varying depending on $\hat{k}(t-1)$ - e.g.,

$$\begin{aligned} k_{min}(t) &= \max\{\hat{k}(t-1) - \kappa, 0\} \\ k_{max}(t) &= \hat{k}(t-1) + \kappa, \end{aligned} \quad (2.14)$$

where $\kappa \geq 1$ is an integer. In case \hat{k} changes from time $t-1$ to time t this solution implies that one or more of the models estimated up to time t should be replaced by others at time $t+1$. To initiate the estimation of the new models, initial values of their $\hat{\theta}$'s, P 's and $\hat{\sigma}_e$'s have to be specified - e.g. $\hat{\theta} = \mathbf{0}$ and $P = KI$ (K being a large positive value) as suggested earlier in this chapter, and $\hat{\sigma}_e$ as a large positive value. As a result of the selection equation in (2.11), a suitable large, positive value of $\hat{\sigma}_e$ ensures that the corresponding model recently initialized will not be selected until the estimates of the parameters have been stabilized (P comes close to its steady-state value).

The method described above is related to the switching autoregressions described by Holst, Lindgren and Holst (1992). They assume an autoregressive process with stochastic parameters which vary in time governed by a Markov regime process. For estimation of the autoregressive parameters, the variance of the innovations and the transition probabilities in the Markov chain they develop a recursive EM-algorithm based on the maximum likelihood method. However, the idea of switching autoregressions is not directly applicable to ARMAX models with time-varying time-delay. In the first place, an AR model and not an ARMAX model is assumed. In the second place, the regimes belong to a discrete distribution while a multivariate continuous distribution for the parameters in the B -polynomial is needed in order to describe the variation of the embedded continuous time-delay of the system.

Bányász and Keviczky's Method

The General Gauss-Newton Method. Bányász and Keviczky (1988) have proposed an algorithm for simultaneous recursive estimation of the parameters and a time-varying delay of an ARMAX model. They are applying a general Gauss-Newton technique to perform a recursive minimization of the least squares loss function $V(\boldsymbol{\theta})$:

$$\min_{\boldsymbol{\theta}} V(\boldsymbol{\theta}), \quad (2.15)$$

where

$$V(\boldsymbol{\theta}) = E[Q(\boldsymbol{\theta})] \quad \text{and} \quad Q(\boldsymbol{\theta}) = \frac{1}{2}\varepsilon^2(t; \boldsymbol{\theta}).$$

$E[\cdot]$ denote the expectation operator, and $\varepsilon(t; \boldsymbol{\theta})$ is the one-step-ahead prediction error with the real parameter vector, $\boldsymbol{\theta}$, given. This vector contains the usual ARMAX parameters together with the time-delay.

Assume that the first and second order derivatives of $\varepsilon(t; \boldsymbol{\theta})$ exist for all $\boldsymbol{\theta}$. Then the general Gauss-Newton algorithm becomes

$$\hat{\boldsymbol{\theta}}(t) = \hat{\boldsymbol{\theta}}(t-1) - \mathbf{P}(t)\mathbf{g}(t), \quad (2.16)$$

where $\mathbf{g}(t)$ is the gradient of $Q(\boldsymbol{\theta})$,

$$\mathbf{g}(t) = \left[\frac{dQ(\boldsymbol{\theta})}{d\boldsymbol{\theta}} \right]_{\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}(t-1)}^T = \frac{d\varepsilon(t; \boldsymbol{\theta})}{d\boldsymbol{\theta}} \varepsilon(t; \boldsymbol{\theta}) \Big|_{\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}(t-1)}, \quad (2.17)$$

and

$$\mathbf{P}(t) = \left[\sum_{j=1}^t \mathbf{H}(j; \boldsymbol{\theta}) \right]_{\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}(t-1)}^{-1}. \quad (2.18)$$

The Hessian $\mathbf{H}(t; \boldsymbol{\theta})$ is defined as the second order derivative of $Q(\boldsymbol{\theta})$,

$$\mathbf{H}(t; \boldsymbol{\theta}) = \frac{d^2 Q(\boldsymbol{\theta})}{d\boldsymbol{\theta}^2} = \left[\frac{d\varepsilon(t; \boldsymbol{\theta})}{d\boldsymbol{\theta}} \right]^T \frac{d\varepsilon(t; \boldsymbol{\theta})}{d\boldsymbol{\theta}} + \frac{d^2 \varepsilon(t; \boldsymbol{\theta})}{d\boldsymbol{\theta}^2} \varepsilon(t; \boldsymbol{\theta}). \quad (2.19)$$

Since the evaluation of $P(t)$ in (2.18) is based upon all observations up to time t , the algorithm is not fully recursive. However, using the Hessian approximation in (2.20) and the matrix inversion lemma, a recursive version of (2.18) can easily be obtained (Ljung and Söderström (1983)).

If $\varepsilon(t; \theta)$ is linear in θ (as, e.g., in the case of an ARX model), then the last term of the expression in (2.19) will vanish. If $\varepsilon(t; \theta)$ is not linear in θ , the expected value of the last term will approach zero close to the real minimum of (2.15). For that reason Ljung and Söderström (1983) suggest that

$$H(t; \theta) = \left[\frac{d\varepsilon(t; \theta)}{d\theta} \right]^T \frac{d\varepsilon(t; \theta)}{d\theta} \quad (2.20)$$

is used in order to ensure that a positive semidefinite approximation of the Hessian is obtained. This approximation corresponds to a linearization of $\varepsilon(t; \theta)$ around $\hat{\theta}(t-1)$ before taking the second derivative.

Tracking a Pure Time-Delay. Now, according to the Bányász-Keviczky algorithm it is possible to make an estimation of the time-delay, k , in a pure time-delay system:

$$y(t) = q^{-k}u(t) + e(t).$$

This is a special case of the ARMAX model in (2.1), ($A(q^{-1}) = B(q^{-1}) = C(q^{-1}) = 1$), and the only parameter is k ($= \theta$). The one-step-ahead predictor and the associated prediction error of this model are

$$\hat{y}(t; k) = q^{-k}u(t) \quad \text{and} \quad \varepsilon(t; k) = y(t) - \hat{y}(t; k). \quad (2.21)$$

Consequently

$$\begin{aligned} \frac{d\varepsilon(t; k)}{dk} &= -\frac{d\hat{y}(t; k)}{dk} = -(-\ln q)q^{-k}u(t) \\ &\approx (1 - q^{-1})\hat{y}(t; k) = \nabla\hat{y}(t; k), \end{aligned} \quad (2.22)$$

where the approximation $\ln x \approx 1 - x^{-1}$ for $x \approx 1$ has been used. Using this result, it is furthermore found that

$$\frac{d^2\varepsilon(t; k)}{dk^2} \approx \frac{d}{dk} [(1 - q^{-1})\hat{y}(t; k)] \approx -\nabla^2\hat{y}(t; k). \quad (2.23)$$

As the parameter space is one-dimensional, (2.18) becomes scalar and can easily be put into recursive form:

$$\begin{aligned} \frac{1}{P(t)} &= \sum_{j=1}^t \lambda^{t-j} H(j; k) = \lambda \sum_{j=1}^t \lambda^{t-j-1} H(j; k) \\ &= \lambda \left(\frac{1}{\lambda} H(t; k) + \sum_{j=1}^{t-1} \lambda^{t-j-1} H(j; k) \right) \\ &= H(t; k) + \frac{\lambda}{P(t-1)}, \end{aligned}$$

where $k = \hat{k}_{t-1}$. Notice that a forgetting factor, λ , has been introduced. Now replace $H(t; k)$ by the expression in (2.19) and use Equations (2.22) and (2.23). Then an updating equation for $P(t)$ is obtained:

$$P(t) = \frac{P(t-1)}{\lambda + P(t-1)[(\nabla\hat{y}(t; k))^2 - \nabla^2\hat{y}(t; k)\varepsilon(t; k)]} \Big|_{k=\hat{k}_{t-1}}. \quad (2.24)$$

Using (2.16), (2.17) and (2.22), we obtain an updating equation for \hat{k}_t :

$$\hat{k}_t = \hat{k}_{t-1} - P(t)\nabla\hat{y}(t; \hat{k}_{t-1})\varepsilon(t; \hat{k}_{t-1}). \quad (2.25)$$

Equations (2.24) and (2.25) complete the Bányász-Keviczky's algorithm.

Remarks:

1. In (2.22) the complex backward shift operator q^{-1} is treated as if it was a real valued variable. But Bányász and Keviczky do not state any reasons for the legality of this. Furthermore, they

do not explain the meaning of the approximation introduced in (2.22). These uncertainties make it doubtful whether the resulting algorithm is reasonable.

2. The Hessian matrix,

$$H(t; k) = (\nabla \hat{y}(t; k))^2 - \nabla^2 \hat{y}(t; k) \varepsilon(t; k), \quad (2.26)$$

which is a part of the denominator in (2.24), may in some cases become negative. Simulation studies have shown that this may cause serious convergence problems in the resulting algorithm. Therefore two modified versions of (2.24) are considered.

As mentioned earlier, the problem can be solved by omitting the last term of the Hessian matrix. Hence (2.24) will be changed to

$$P(t) = \frac{P(t-1)}{\lambda + P(t-1)(\nabla \hat{y}(t; k))^2} \Big|_{k=\hat{k}_{t-1}}.$$

Another alternative, which has not been considered by Bányász and Keviczky, is

$$P(t) = \begin{cases} \frac{P(t-1)}{\lambda + P(t-1)[(\nabla \hat{y}(t; k))^2 - \nabla^2 \hat{y}(t; k) \varepsilon(t; k)]} & , \quad H(t; k) \geq 0 \\ P^* & , \quad H(t; k) < 0 \end{cases}, \quad (2.27)$$

where $k = \hat{k}_{t-1}$, and P^* is a suitable large positive number. This is the equation being used together with (2.25) later in this chapter when simulation experiments are carried out.

In the above algorithm \hat{k}_t is a real number. However, in (2.21) an integral estimate of k is needed to calculate the prediction error, $\varepsilon(t; \hat{k}_{t-1})$, which is included in the updating equations for $P(t)$ and \hat{k}_t . Bányász and Keviczky (1988) suggest that

$$\hat{k}_t^* = \text{round}(\hat{k}_t) \quad (2.28)$$

is used for prediction purposes, where $\text{round}(\cdot)$ denotes the function which rounds off to the nearest integer.

A Continuous Time-Delay Approach

The Weighted Pure Time-Delay Model. Reconsider the pure time-delay model,

$$y(t) = u(t - k) + e(t). \quad (2.29)$$

Søgaard and Madsen (1989) proposed a new method for recursive estimation of the time-delay, k , in this model. This method introduces the *weighted* pure time-delay model,

$$y(t) = w_{-1}(\alpha)u(t - d + 1) + w_0(\alpha)u(t - d) + w_1(\alpha)u(t - d - 1) + e(t), \quad (2.30)$$

where

$$w_{-1}(\alpha) = \frac{1}{2}\alpha(\alpha - 1), \quad w_0(\alpha) = 1 - \alpha^2, \quad w_1(\alpha) = \frac{1}{2}\alpha(\alpha + 1)$$

are weighting functions in $\alpha \in [-1, 1]$. Introducing the notation $M(d, \alpha)$ for the model (2.30), the following is true:

1. $w_{-1}(\alpha) + w_0(\alpha) + w_1(\alpha) = 1$ for all α .
2. There is symmetry around $\alpha = 0$: $w_1(\alpha) = w_{-1}(-\alpha)$ and $w_0(\alpha) = w_0(-\alpha)$ (see Figure 2.2).
3. The derivatives of the weight functions are continuous.
4. $M(d, -1)$, $M(d, 0)$ and $M(d, 1)$ represent 3 different, pure time-delay models with the time-delays $d-1$, d and $d+1$, respectively.
5. $M(p + 1, -1) = M(p, 0) = M(p - 1, 1)$ - i.e. the same model can be obtained for three different pairs of α and d . For these pairs the sum $d + \alpha$ remains constant - in this case $d + \alpha = p$.
6. The quantity $d + \alpha$ can be considered an approximate measure for an embedded continuous delay of model $M(d, \alpha)$. Thus the continuous delay is partitioned into an integral part, d , and a real part, α .

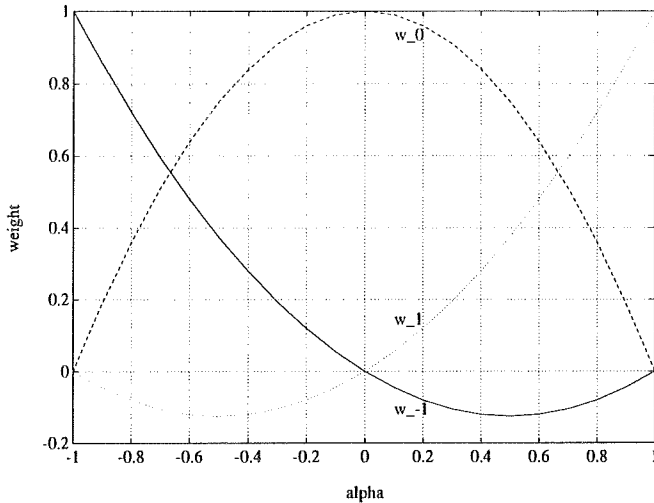


Figure 2.2: *Weight functions in the weighted pure time delay-model.*

Alternative parametrizations of the weight functions might of course be suggested. However, it is difficult to find alternative weight functions which are simple and possess the qualities stated above.

The introduction of the weight functions in (2.30) enables us to move continuously between three different pure time-delay models by varying α from -1 to 1 (d fixed). An adaptive estimation of the continuous delay parameter, α , is the basic idea behind the algorithms presented below.

Figure 2.2 shows the weight functions. For $\alpha \in [-1, 0]$ it is seen that $w_{-1}(\alpha)$ and $w_0(\alpha)$ represent the largest weights, whereas $w_1(\alpha)$ is close to zero. This means that the input is included in the weighted pure time-delay model mainly as a weighted sum of $u(t - d + 1)$ and $u(t - d)$ reflecting a time-delay between d and $d - 1$ ($\approx d + \alpha$). Due to symmetry, the time-delay is approximately $d + \alpha$ for $\alpha \in [0, 1]$ as

well.

Recursive Estimation of α . As the first step, an algorithm for a recursive estimation of α is derived. The integral part, d , of the delay is considered to be known constant during this derivation.

The one-step-ahead prediction error of model (2.30) becomes

$$\begin{aligned} \varepsilon(t; \alpha, d) = \\ y(t) - [w_{-1}(\alpha)u(t-d+1) + w_0(\alpha)u(t-d) + w_1(\alpha)u(t-d-1)] . \end{aligned} \quad (2.31)$$

As this expression is non-linear in α , the RLS algorithm cannot be used directly for recursive estimation of α . Therefore, the Gauss-Newton approximation (see Ljung (1987)) is adopted. This method implies that the derivative of $\varepsilon(t)$ with respect to α should be computed. That is,

$$\begin{aligned} \varphi(t; \alpha, d) = -\frac{\partial \varepsilon(t; \alpha, d)}{\partial \alpha} \\ = (\alpha - \frac{1}{2})u(t-d+1) - 2\alpha u(t-d) + (\alpha + \frac{1}{2})u(t-d-1) . \end{aligned} \quad (2.32)$$

By using (2.31), (2.32) and $\hat{\theta}(t) = \hat{\alpha}(t)$ together with (2.2)-(2.4), an approximate RLS algorithm for recursive estimation of α is obtained:

$$\hat{\alpha}(t) = \hat{\alpha}(t-1) + L(t)\varepsilon(t; \hat{\alpha}(t-1), d) \quad (2.33)$$

$$L(t) = P(t)\varphi(t; \hat{\alpha}(t-1), d) \quad (2.34)$$

$$P(t) = \frac{P(t-1)}{\lambda_\alpha + \varphi^2(t; \hat{\alpha}(t-1), d)P(t-1)} . \quad (2.35)$$

This algorithm based on the model (2.30) is well suited for a recursive estimation of the pure time-delay model (2.29) with the time-delay k being exactly one of the three consecutive values: $k_0 - 1$, k_0 and $k_0 + 1$. In (2.30) those values of k correspond to a fixed $d = k_0$ and

$\alpha \in \{-1, 0, 1\}$. Since α is a continuous parameter, the rule,

$$\hat{k}(t) = \begin{cases} k_0 - 1, & \text{if } \hat{\alpha}(t) < -\frac{1}{2} \\ k_0 + 1, & \text{if } \hat{\alpha}(t) > \frac{1}{2} \\ k_0 & \text{otherwise,} \end{cases}$$

can be applied in order to obtain an estimate of k at time t .

Recursive Estimation of d and α . If the time-delay is not limited to three consecutive discrete values, $k_0 - 1$, k_0 and $k_0 + 1$, as assumed above, the algorithm must be extended in order to estimate d and α simultaneously:

$$(\hat{\alpha}^*(t), P(t)) = \text{RLS}_\alpha(\alpha, P(t-1), \varphi(t; \alpha, d), \varepsilon(t; \alpha, d), \lambda_\alpha) \Big|_{\alpha=\hat{\alpha}(t-1), d=\hat{k}(t-1)} \quad (2.36)$$

$$\Delta \hat{k}(t) = \text{round}(\hat{\alpha}^*(t)) \quad (2.37)$$

$$(\hat{\alpha}(t), \hat{k}(t)) = (\hat{\alpha}^*(t) - \Delta \hat{k}(t), \hat{k}(t-1) + \Delta \hat{k}(t)), \quad (2.38)$$

where $\text{RLS}_\alpha(\dots)$ denotes the algorithm (2.33)-(2.35) and $\text{round}(\cdot)$ is a function rounding off its argument to the nearest integer:

$$\text{round}(x) = \begin{cases} i & \text{if } i - \frac{1}{2} \leq x < i + \frac{1}{2}, \quad \text{where } i = 1, 2, 3, \dots \\ j & \text{if } j - \frac{1}{2} < x \leq j + \frac{1}{2}, \quad \text{where } j = -1, -2, -3, \dots \\ 0 & \text{if } -\frac{1}{2} < x < \frac{1}{2}. \end{cases}$$

After having computed $\hat{\alpha}^*(t)$ in (2.36), the sum $\hat{k}(t-1) + \hat{\alpha}^*(t)$ can be regarded as an estimate of the continuous time-delay at time t . Equations (2.37) and (2.38) ensure that $\hat{\alpha}(t)$ remains in the interval $[-\frac{1}{2}, \frac{1}{2}]$ without the estimate of the continuous delay is changed - i.e. the estimate, $\hat{k}(t)$, of the discrete delay complies with the equation $\hat{k}(t) + \hat{\alpha}(t) = \hat{k}(t-1) + \hat{\alpha}^*(t)$.

In order to make the algorithm more robust to data influenced by noise, the restrictions on $\hat{k}(t)$ considered earlier can be introduced:

$k_{min} \leq \hat{k}(t) \leq k_{max}$ and $0 \leq |\hat{k}(t) - \hat{k}(t-1)| \leq \Delta k_{max}$. In this case (2.37) is still used if $\hat{k}(t-1) + \text{round}(\hat{\alpha}^*(t))$ is within the interval

$$[k_{min}, k_{max}] \cap [\hat{k}(t-1) - \Delta k_{max}, \hat{k}(t-1) + \Delta k_{max}]$$

Otherwise the nearest end point of this interval is used.

Recursive Estimation of the ARMAX model. The algorithm can be further extended to include the general ARMAX model (2.1). One way to do this is to make use of the fact that the ARMAX model can be formulated as a pure time-delay model with filtered input and output:

$$\tilde{y}(t) = \tilde{u}(t-k) + e(t), \quad (2.39)$$

where

$$\tilde{y}(t) = A(q^{-1})y(t) + (1 - C(q^{-1}))e(t) \quad \text{and} \quad \tilde{u}(t) = B(q^{-1})u(t). \quad (2.40)$$

Using (2.39), an algorithm for recursive estimation of the time-delay, k , and the parameters $a_1, \dots, a_n, b_0, \dots, b_m$ and c_1, \dots, c_r is obtained:

$$\varepsilon_\theta(t) = y(t) - \varphi_\theta^T(t) \hat{\theta}(t-1) \Big|_{k=\hat{k}(t-1)} \quad (2.41)$$

$$\begin{aligned} (\hat{\theta}(t), \mathbf{P}_\theta(t)) &= \text{ELS}(\hat{\theta}(t-1), \mathbf{P}_\theta(t-1), \\ &\quad \varphi_\theta(t), y(t), \lambda_\theta) \Big|_{k=\hat{k}(t-1)} \end{aligned} \quad (2.42)$$

$$\hat{y}(t) = \hat{A}_i(q^{-1})y(t) + (1 - \hat{C}_i(q^{-1}))\varepsilon_\theta(t) \quad (2.43)$$

$$\hat{u}(t) = \hat{B}_i(q^{-1})u(t) \quad (2.44)$$

$$\begin{aligned} (\hat{\alpha}(t), \hat{k}(t), P_\alpha(t)) &= \text{RLS}_{\alpha, k}(\alpha, d, P_\alpha(t-1), \varphi_\alpha(t; \alpha, d), \\ &\quad \varepsilon_\alpha(t; \alpha, d), \lambda_\alpha) \Big|_{\alpha=\hat{\alpha}(t-1), d=\hat{k}(t-1)} \end{aligned} \quad (2.45)$$

Equation (2.45) is a short notation for the algorithm (2.36)-(2.38) with $\varepsilon_\alpha(t; \hat{\alpha}(t-1), \hat{k}(t-1))$ and $\varphi_\alpha(t; \hat{\alpha}(t-1), \hat{k}(t-1))$ computed from (2.31) and (2.32), where $u(t)$ and $y(t)$ are replaced by $\hat{u}(t)$ and $\hat{y}(t)$, respectively.

The algorithm (2.41)-(2.45) combines the ELS algorithm for estimation of the continuous parameters of an ARMAX model and the RLS algorithm for estimation of the discrete time-delay of a pure time-delay model. To avoid mixing up the φ 's, P 's and ε 's of the two algorithms they have been given the subscripts θ and α .

Equations (2.41) and (2.42) are the updating steps of the ELS algorithm (see (2.2)-(2.4) and (2.5)). Note that $\hat{\theta}(t)$ is updated with $k = \hat{k}(t-1)$. Equations (2.43) and (2.44) constitute the filters introduced in (2.40) with the coefficients in A , B and C being replaced by the estimates computed in (2.42).

Notice that the algorithm (2.41)-(2.45) splits the estimation of the parameters and the time-delay into two separate computations. Also notice that the parameters are not assumed to be directly dependent on the time-delay as they are in the "parallel model" approach (see algorithm (2.9)-(2.12)).

2.2.2 Modelling the Variations of the Time-Delay

In the above sections, the time-delay was basically considered a constant parameter of the ARMAX model. Although the real delay was known to be time-varying, no model of this variation was specified. The variation was taken into consideration at the estimation phase where a forgetting factor technique was used. However, for systems showing rapidly changing dynamics, forgetting factor techniques might give bad results. The reason is that such methods rely on a local approximation of the real time-varying system by means of a time-invariant model (e.g. a transfer function model). Due to estimates based on past observations of the input and output processes, the parameter estimates are already "out of date" when calculated. For rapidly changing systems it would be preferable to choose a small

forgetting factor value to forget old dynamics quickly, but by doing that the estimator is made rather sensitive to noise. Therefore the forgetting factor must be chosen as a compromise between two conflicting interests. Such a compromise might lead to a poor estimator.

Explicit models of the variation of the parameters and the time-delay of transfer functions are considered below. The first category of models includes a deterministic description of the time-varying system, and in a later section the embedded time-variation is described by means of a stochastic model.

Both the deterministic and the stochastic approaches are based on a time-varying version of the ARMAX model in (2.1):

$$A_t(q^{-1})y(t) = B_t(q^{-1})q^{-k_t}u(t) + C_t(q^{-1})e(t), \quad (2.46)$$

where the time-delay and the coefficients of the polynomials A_t , B_t and C_t are time-varying:

$$\begin{aligned} A_t(q^{-1}) &= 1 + a_{1,t}q^{-1} + \cdots + a_{n,t}q^{-n} \\ B_t(q^{-1}) &= b_{0,t} + b_{1,t}q^{-1} + \cdots + b_{m,t}q^{-m} \\ C_t(q^{-1}) &= 1 + c_{1,t}q^{-1} + \cdots + c_{r,t}q^{-r}. \end{aligned} \quad (2.47)$$

The orders, n , m and r , of the polynomials are assumed constant and known.

The following sections describe how to specify models of the embedded variation of delay and coefficients.

Deterministic Variations

In this section the discrete delay parameter, k_t , is assumed to be a constant, k . The embedded variation of the real continuous time-

delay is modelled through the variation of the coefficients of the B_t polynomial.

Let $d_i(\boldsymbol{\theta})$ denote any coefficient of the polynomials in (2.47). The vector

$$\boldsymbol{\theta} = (\theta_1 \theta_2 \cdots \theta_p)^T$$

makes up the set of parameters that is included in the deterministic function describing the time-variation of $d_i(\boldsymbol{\theta})$. Here only functions being linear in the parameters are discussed – i.e.

$$d_i(\boldsymbol{\theta}) = \theta_1 f_1(t) + \theta_2 f_2(t) + \cdots + \theta_p f_p(t), \quad (2.48)$$

where $f_i(t)$ ($i = 1, \dots, p$) is assumed to be a known function of time. It can be an explicit function of time or it can be a function of a time-varying variable – e.g. $u(t)$ or another sampled input signal.

The linear parametrization of the time-variation in (2.48) implies that the resulting time-varying ARMAX model becomes linear in the B_t -parameters, hence permitting ELS estimation. An example illustrates the point.

Example 2.1 Consider the following time-varying ARX model,

$$y(t) = -a_t y(t-1) + b_{1,t} u(t-1) + b_{2,t} u(t-2) + e(t), \quad (2.49)$$

with

$$\begin{aligned} a_t &= \alpha f(t) \\ b_{1,t} &= \beta_1 g_1(t) + \beta_2 g_2(t) \\ b_{2,t} &= \gamma h(t), \end{aligned} \quad (2.50)$$

where $f(t)$, $g_1(t)$, $g_2(t)$ and $h(t)$ are known functions of time. Insertion of (2.50) into (2.49) results in

$$\begin{aligned} y(t) &= -\alpha f(t) y(t-1) + \beta_1 g_1(t) u(t-1) + \beta_2 g_2(t) u(t-1) \\ &\quad + \gamma h(t) u(t-2) + e(t), \end{aligned}$$

or

$$y(t) = \varphi^T(t)\theta + e(t), \quad (2.51)$$

where

$$\begin{aligned} \varphi(t) = & \\ & (-f(t)y(t-1) \quad g_1(t)u(t-1) \quad g_2(t)u(t-1) \quad h(t)u(t-2))^T \\ \theta(t) = & (\alpha \quad \beta_1 \quad \beta_2 \quad \gamma)^T. \end{aligned}$$

Since the model can be put into the form of (2.51) it belongs to the class of general linear models. Consequently the recursive least squares algorithm in (2.2)-(2.4) is directly applicable.

□

Time-Varying Transfer Functions for a District Heating System. Sjøgaard (1988) identified transfer function models for points in the distribution network of the district heating system in Esbjerg. One of the these points is called Hjerting Strandvej. The input series, $\{u(t)\}$, consists of hourly averages of the supply temperature at the combined heat and power plant, Vestkraft. The output series, $\{y(t)\}$, consists of hourly averages of the supply temperature at Hjerting Strandvej. The following models give a reasonable description of the input-output system.

$$(1 - 0.838q^{-1})y(t) = (0.314 - 0.172q^{-2})u(t-5) + e(t) \quad (2.52)$$

$$\begin{aligned} (1 - 0.835q^{-1})y(t) & \\ & = (0.320 - 0.177q^{-2})u(t-5) \\ & + (0.136 - 0.263q^{-1} + 0.136q^{-2})u_s(t-5) + e(t) \quad (2.53) \end{aligned}$$

The “pseudo” input signal $\{u_s(t)\}$ is defined as

$$u_s(t) = u(t) \sin\left(\frac{\pi}{12}t\right),$$

and $\{e(t)\}$ denotes a white noise process. Due to the sine terms in model (2.53), it is necessary to relate the time index, t , to the hour of the day. This is done as follows: $t = 1 \sim 01:00$.

The model (2.52) was identified as an ordinary time-invariant ARMAX model, while the model (2.53) was identified as a time-varying ARMAX model (from the same data set). Søgaard (1988) showed that the model (2.53) which has terms involving $u_s(t)$ is equivalent to a model with time-varying coefficients in the B -polynomial. For this model it is easily shown that

$$\begin{aligned} b_{0,t} &= 0.320 + 0.136 \sin\left(\frac{\pi}{12}(t - 5)\right) \\ b_{1,t} &= -0.263 \sin\left(\frac{\pi}{12}(t - 6)\right) \\ b_{2,t} &= -0.177 + 0.136 \sin\left(\frac{\pi}{12}(t - 7)\right). \end{aligned} \quad (2.54)$$

Time-variation of the type $\theta_1 + \theta_2 \cos\left(\frac{\pi}{12}t\right) + \theta_3 \sin\left(\frac{\pi}{12}t\right)$ was tried for all the parameters in the model but statistical significance tests showed that the time-variation was limited to the B_t -polynomial and that only sine terms should be included. Notice that none of the parameters in the C_t -polynomial became significant.

A graph of the diurnal variation of the coefficients is shown in Figure 2.3. The coefficients of B_t are modelled by means of the first harmonic of the Fourier expansion with a basis period of 24 hours. In this way, a major part of the diurnal variation of the delay from input to output is described (a better approximation of the diurnal variation could of course have been achieved by including higher order harmonics of the Fourier expansion). The diurnal variation of the delay is due to variations of the flow of water in the distribution network, which in turn is induced by the varying heat consumption (low at night and high in the day-time).

The three time-varying coefficients correspond to two derived coefficients and a time-varying delay. This is explained below.

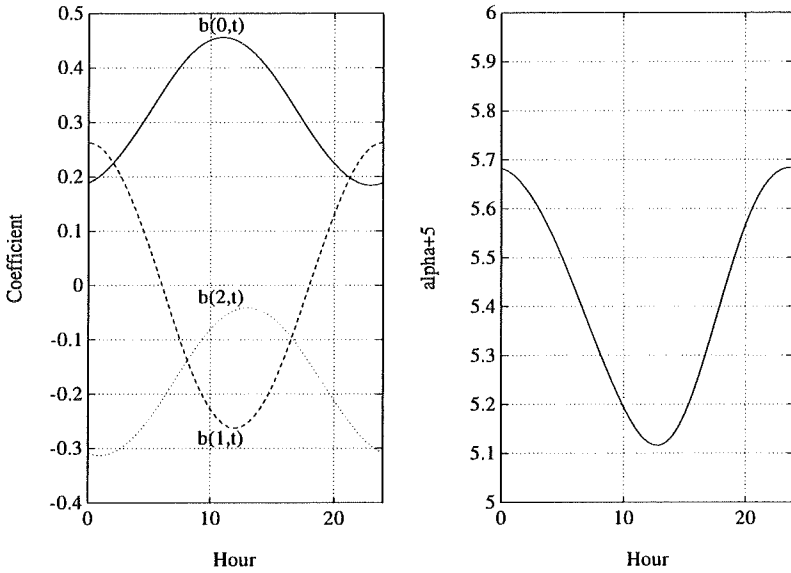


Figure 2.3: *Left: Diurnal variation of the coefficients of the B_i polynomial. Right: Diurnal variation of the time-delay.*

The input series $\{u(t)\}_{t=1,2,3,\dots}$ is obtained by sampling the corresponding continuous input signal. An approximate reconstruction of the continuous signal can be created by linear interpolation between the sampled input values:

$$u(t - \alpha) \approx (1 - \alpha)u(t) + \alpha u(t - 1), \quad 0 \leq \alpha \leq 1. \quad (2.55)$$

The sampled input values are included in the model in (2.53) as

$$w(t) = b_{0,t}u(t - 5) + b_{1,t}u(t - 6) + b_{2,t}u(t - 7), \quad (2.56)$$

but we could obtain the same result by using an input series shifted α_t of an hour backwards in time (by means of (2.55)):

$$\begin{aligned} w(t) &= \gamma_{1,t}u(t - \alpha_t - 5) + \gamma_{2,t}u(t - \alpha_t - 6) \\ &= \gamma_{1,t}(1 - \alpha_t)u(t - 5) + (\gamma_{1,t}\alpha_t + \gamma_{2,t}(1 - \alpha_t))u(t - 6) \\ &\quad + \gamma_{2,t}\alpha_t u(t - 7), \end{aligned} \quad (2.57)$$

where $\gamma_{1,t}$ and $\gamma_{2,t}$ are new coefficients and $\alpha_t + 5$ is the approximate continuous time-delay ($0 \leq \alpha_t \leq 1$). These two coefficients and α_t can be found by equating the right hand side of (2.56) with the last expression in (2.57). This leads to three equations with three unknowns:

$$\begin{aligned} b_{0,t} &= \gamma_{1,t}(1 - \alpha_t) \\ b_{1,t} &= \gamma_{1,t}\alpha_t + \gamma_{2,t}(1 - \alpha_t) \\ b_{2,t} &= \gamma_{2,t}\alpha_t. \end{aligned}$$

This system of equations results in a quadratic equation in α_t :

$$(b_{0,t} + b_{1,t} + b_{2,t})\alpha_t^2 - (b_{1,t} + 2b_{2,t})\alpha_t + b_{2,t} = 0,$$

From (2.54) it is evident that $b_{0,t} > 0$ and $b_{2,t} < 0$ for all values of t . Consequently the discriminant, $D_t = b_{1,t}^2 - 4b_{0,t}b_{2,t}$, of the quadratic equation is positive for all values of t , implying that there

are two solutions. Only one of them, however, satisfies the condition $0 \leq \alpha_t \leq 1$.

In Figure 2.3 the diurnal variation of the delay, $\alpha_t + 5$, is shown. Note that the maximum delay occurs at night, while the minimum is found in the day-time – just as expected. Also note that the stationary gain from input to output becomes time-varying and can be computed as $B_t(1)/(1 - A(1))$. From the parameter estimates it is found that the variation of the stationary gain is rather insignificant: it varies from 0.868 (maximum) in the the day-time to 0.865 (minimum) at night. This means that the district heating water is cooled about 13 % on its way from Vestkraft to Hjerting Strandvej. This percentage is confirmed by temperature measurements from the district heating system in Esbjerg (see Søggaard (1988)).

The estimated standard deviations of the one-step-ahead prediction error (the standard error) for the models (2.52) and (2.53) are 0.497 and 0.484, respectively. Thus the standard error of model (2.53) is $(0.497 - 0.484)/0.497 = 2.6\%$ less than the standard error of model (2.52). This reduction is due to the introduction of time-varying coefficients.

The transfer function structure (2.53) has been used as a basis for the models which have been implemented at the combined heat and power plant Vestkraft in Esbjerg. In this implementation the parameters and the discrete time-delays of the models are estimated adaptively. Since the model structure is rather general and the parameters and delay are estimated adaptively it would be reasonable to apply this method for other district heating systems.

Stochastic State-Space Model

In this section, an explicit model of the parameters and the time-delay is given. The model is then formulated as a non-linear state-space model with the parameters and the time-delay being the state variables. Finally, an extended Kalman filter is employed in order to estimate the state variables. Non-linear state-space models and the related extended Kalman filters are, for instance, described by Harvey (1989).

A State-Space Formulation. The following model structure is considered:

$$y_t = a_t y_{t-1} + b_t((1 - \alpha_t)u_{t-k_t} + \alpha_t u_{t-k_t-1}) + e_t \quad (2.58)$$

$$\alpha_t = \delta_t - k_t \quad (2.59)$$

$$k_t = \text{int}(\delta_t) \quad (2.60)$$

$$\delta_t = \phi \delta_{t-1} + \mu_t + \epsilon_t \quad (2.61)$$

$$a_t = a_{t-1} + \eta_t \quad (2.62)$$

$$b_t = b_{t-1} + \xi_t \quad (2.63)$$

$$\mu_t = \mu_{t-1} + \kappa_t . \quad (2.64)$$

This model structure appears to be reasonable for the district heating system in Ishøj which is studied later in this chapter.

The model in (2.58) is a time-varying ARMAX model with $n = 1$, $m = 2$ and $r = 0$ (an ARX model). This special structure is based upon experiments with data from a district heating system. The unusual non-linear parametrization of the B_t -polynomial,

$$B_t(q^{-1}) = b_t(1 - \alpha_t) + b_t \alpha_t q^{-1} ,$$

permits an approximate description of the embedded continuous time-delay. Below it is explained why.

Let $\{\tilde{u}(\tau)\}$ denote the continuous input signal which has been sampled to get $\{u_i\}$. Thus, the following approximation can be made:

$$\tilde{u}([t - (k_t + \alpha_t)]T) \approx (1 - \alpha_t)u_{t-k_t} + \alpha_t u_{t-k_t-1} \quad (2.65)$$

where $0 \leq \alpha_t \leq 1$, and T is the sampling interval. The right hand side of (2.65) expresses a linear interpolation between u_{t-k_t} and u_{t-k_t-1} . Using this continuous time approximation, the model in (2.58) can be rewritten as

$$y_t = a_t y_{t-1} + b_t \tilde{u}([t - (k_t + \alpha_t)]T) + e_t .$$

For this reason, (2.58) should be thought of as an ARMAX model with $m = 1$ and a continuous time-delay $\delta_t = k_t + \alpha_t$ rather than an ARMAX model with $m = 2$. The relationship $\delta_t = k_t + \alpha_t$ and the restriction $0 \leq \alpha_t \leq 1$ are expressed in Equations (2.59) and (2.60). In (2.60) $\text{int}(\cdot)$ denotes the integral part of the argument - i.e. $\text{int}(x)$ is the integer i for which $0 \leq x - i < 1$.

In (2.61), the model of the variations of the continuous time-delay, δ_t , is given. It is an ARMAX(1,1,0) model driven by μ_t and white noise $\{\epsilon_t\}$. μ_t has been introduced to make the mean of δ_t non-zero.

The ARMAX parameters, a_t and b_t , as well as μ_t are assumed to be the result of random walk processes (Equations (2.62)-(2.64)). These processes are driven by white noise: $\{\eta_t\}$, $\{\xi_t\}$ and $\{\kappa_t\}$. Note that the random walk model (2.62) may result in values of the AR-parameter, a_t , which are outside the interval $[-1, 1]$. Consequently the time-varying ARX-model in (2.58) may become unstable.

An alternative model of a_t which reduces the probability of instability is

$$a_t = \beta a_{t-1} + \gamma + \eta_t ,$$

where $-1 < \beta < 1$ and $-1 < \gamma < 1$. Contrary to the random walk model this model generates a stochastic process with finite variance,

and the probability of a_t values outside the stable interval can be made arbitrarily small by choosing β , γ and the variance of η_t suitable close to zero.

The correlation structure of the various noise components which have been introduced is assumed to be as follows

$$\text{Var}[(e_t, \epsilon_t, \eta_t, \xi_t, \kappa_t)^T] = \begin{pmatrix} \sigma_e^2 & 0 & 0 & 0 & 0 \\ 0 & \sigma_\epsilon^2 & \sigma_{\epsilon\eta} & \sigma_{\epsilon\xi} & \sigma_{\epsilon\kappa} \\ 0 & \sigma_{\eta\epsilon} & \sigma_\eta^2 & \sigma_{\eta\xi} & \sigma_{\eta\kappa} \\ 0 & \sigma_{\xi\epsilon} & \sigma_{\xi\eta} & \sigma_\xi^2 & \sigma_{\xi\kappa} \\ 0 & \sigma_{\kappa\epsilon} & \sigma_{\kappa\eta} & \sigma_{\kappa\xi} & \sigma_\kappa^2 \end{pmatrix}. \quad (2.66)$$

Through the introduction of the state vector $\mathbf{x}_t = (x_{1,t}, x_{2,t}, x_{3,t}, x_{4,t})^T$, where $x_{1,t} = \delta_t$, $x_{2,t} = a_t$, $x_{3,t} = b_t$ and $x_{4,t} = \mu_t$, a state-space formulation of the model (2.58)-(2.64) is obtained:

$$\mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \boldsymbol{\nu}_t \quad (2.67)$$

$$y_t = \left. \begin{aligned} &x_{2,t}y_{t-1} + x_{3,t}((1 - \alpha_t)u_{t-k_t} + \alpha_t u_{t-k_t-1}) + e_t \\ &\text{with } k_t = \text{int}(x_{1,t}) \text{ and } \alpha_t = x_{1,t} - k_t \end{aligned} \right\}, \quad (2.68)$$

where

$$\mathbf{A} = \begin{pmatrix} \phi & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}, \quad \boldsymbol{\nu}_t = \begin{pmatrix} \epsilon_t + \kappa_t \\ \eta_t \\ \xi_t \\ \kappa_t \end{pmatrix} \quad (2.69)$$

$$\begin{aligned} \text{Var}[\boldsymbol{\nu}_t] &= \Sigma_1 \\ &= \begin{pmatrix} \sigma_\epsilon^2 + 2\sigma_{\epsilon\kappa} + \sigma_\kappa^2 & \sigma_{\epsilon\eta} + \sigma_{\kappa\eta} & \sigma_{\epsilon\xi} + \sigma_{\kappa\xi} & \sigma_{\epsilon\kappa} + \sigma_\kappa^2 \\ \sigma_{\eta\epsilon} + \sigma_{\eta\kappa} & \sigma_\eta^2 & \sigma_{\eta\xi} & \sigma_{\eta\kappa} \\ \sigma_{\xi\epsilon} + \sigma_{\xi\kappa} & \sigma_{\xi\eta} & \sigma_\xi^2 & \sigma_{\xi\kappa} \\ \sigma_{\kappa\epsilon} + \sigma_\kappa^2 & \sigma_{\kappa\eta} & \sigma_{\kappa\xi} & \sigma_\kappa^2 \end{pmatrix} \end{aligned} \quad (2.70)$$

$$\text{Var}[e_t] = \Sigma_2 = \sigma_e^2. \quad (2.71)$$

Equation (2.67) is called the system equation and (2.68) the observation equation. The system equation is linear in \mathbf{x}_t while the observation equation is not – i.e. it *cannot* be written as

$$y_t = \mathbf{C}_t \mathbf{x}_t + e_t, \quad (2.72)$$

where \mathbf{C}_t is a time-varying vector.

A Kalman Filter Approach. Due to the non-linearity of (2.68), the ordinary Kalman filter (see Appendix A) is not directly applicable for tracking the state variables of the system described by (2.67)-(2.68). To solve this problem, a linearization of (2.68) around the current state estimates is made. This implies that the \mathbf{C} -vector of the Kalman filter equations (A.9) and (A.10) is replaced by

$$\begin{aligned} \hat{\mathbf{C}}_{t|t-1} &= \left. \frac{\partial \hat{y}_{t|t-1}}{\partial \mathbf{x}_t} \right|_{\mathbf{x}_t = \hat{\mathbf{x}}_{t|t-1}} \\ &= (x_{3,t}(u_{t-k_t-1} - u_{t-k_t}), y_{t-1}, \\ &\quad (1 - \alpha_t)u_{t-k_t} + \alpha_t u_{t-k_t-1}, 0) \Big|_{\mathbf{x}_t = \hat{\mathbf{x}}_{t|t-1}}, \end{aligned} \quad (2.73)$$

where

$$\hat{y}_{t|t-1} = x_{2,t}y_{t-1} + x_{3,t}((1 - \alpha_t)u_{t-k_t} + \alpha_t u_{t-k_t-1}),$$

and

$$k_t = \text{int}(x_{1,t}) \quad \text{and} \quad \alpha_t = x_{1,t} - k_t.$$

Furthermore, the one-step-ahead prediction, $\mathbf{C}\hat{\mathbf{x}}_{t|t-1}$, in Equation (A.8) is replaced by

$$\hat{y}_{t|t-1} = x_{2,t}y_{t-1} + x_{3,t}((1 - \alpha_t)u_{t-k_t} + \alpha_t u_{t-k_t-1}) \Big|_{\mathbf{x}_t = \hat{\mathbf{x}}_{t|t-1}}. \quad (2.74)$$

In both Equations (2.73) and (2.74) the following substitutions should be done: $k_t = \text{int}(x_{1,t})$ and $\alpha_t = x_{1,t} - k_t$.

The modifications of the Kalman filter introduced in (2.73) and (2.74) lead to an extended Kalman filter for the tracking of \mathbf{x}_t ,

$$\hat{\mathbf{x}}_{t|t-1} = A\hat{\mathbf{x}}_{t-1|t-1} \quad (2.75)$$

$$\mathbf{P}_{t|t-1} = A\mathbf{P}_{t-1|t-1}A^T + \Sigma_1 \quad (2.76)$$

$$\hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{K}_t(y_t - \hat{y}_{t|t-1}) \quad (2.77)$$

$$\mathbf{P}_{t|t} = \mathbf{P}_{t|t-1} - \mathbf{K}_t\hat{\mathbf{C}}_{t|t-1}\mathbf{P}_{t|t-1} \quad (2.78)$$

$$\mathbf{K}_t = \mathbf{P}_{t|t-1}\hat{\mathbf{C}}_{t|t-1}^T(\hat{\mathbf{C}}_{t|t-1}\mathbf{P}_{t|t-1}\hat{\mathbf{C}}_{t|t-1}^T + \Sigma_2)^{-1}. \quad (2.79)$$

An aspect which has not been taken into consideration in this modified Kalman filter is the additional variance introduced by the use of a stochastic \mathbf{C} -vector instead of a deterministic one. However, this lack of variance can in part be compensated by adjusting the values of Σ_1 and Σ_2 . In Section 2.3, Σ_1 , Σ_2 and ϕ are estimated by optimization of an objective function. Thus the estimates of Σ_1 and Σ_2 are expected to be biased.

2.3 Experimental Results

Below are given some results from experiments with the various delay tracking methods from the previous sections. The presentation is organized as follows:

Section 2.3.1: Results from the tracking of time-delay and dynamics of simulated input-output systems are reviewed. These results are also found in Sjøgaard and Madsen (1989).

Section 2.3.2: Some of the methods have been used for modelling a district heating system. The results in this section were first presented in Madsen *et al.* (1990).

Section 2.3.3: One of the methods has been used to track the time-delays of business cycle forecasting models. The results given in this section were first presented in Edlund and Sjøgaard (1993).

2.3.1 Simulation Results

In this section the three forgetting factor methods from Section 2.2.1 are tested with simulated data. The first part of the investigations concerns simulated data from a time-delay model,

$$y(t) = u(t - k_t) + e(t), \quad \text{Var}[e(t)] = \sigma_e^2 = 0.5^2. \quad (2.80)$$

The second part concerns an ARX model,

$$\begin{aligned} y(t) &= ay(t-1) + b_0u(t - k_t) + b_1u(t - k_t - 1) + e(t), \\ a &= 0.75, \quad b_0 = 0.10, \quad b_1 = 0.14, \quad \text{Var}[e(t)] = \sigma_e^2 = 0.5^2. \end{aligned} \quad (2.81)$$

In both cases, the noise process, $\{e(t)\}$, is Gaussian white noise with the mean zero. For the pure time-delay model, the time-varying delay, k_t , is estimated recursively, and for the ARX model both the delay and the parameters (a , b_0 and b_1) are subject to recursive estimation.

The input signal, $\{u(t)\}$, for the simulation examples is an AR(1) process:

$$u(t) = 0.5u(t-1) + v(t), \quad \text{Var}[v(t)] = \sigma_v^2 = 10^2, \quad (2.82)$$

where $\{v(t)\}$ is Gaussian white noise with mean zero.

The time-varying delay is generated by a random walk type of process:

$$k_t = \begin{cases} k_{t-1} + 1 & \text{with probability } p \\ k_{t-1} & \text{with probability } 1 - 2p \\ k_{t-1} - 1 & \text{with probability } p \end{cases}, \quad (2.83)$$

where $0 \leq p \leq \frac{1}{2}$ and $k_0 = 8$.

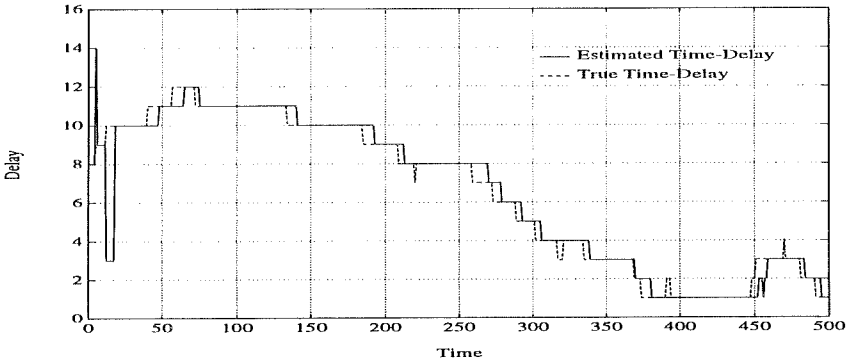


Figure 2.4: *Tracking a pure time-delay by estimation of models in parallel. Simulation: Eqn. (2.80), (2.82), (2.83). Tracking: Eqn. (2.9)-(2.12) with $k_{min} = 1$, $k_{max} = 20$, $\lambda_{\theta} = \lambda_e = 0.9$.*

Tracking a Pure Time-Delay

Recursive estimation of the time-delay in (2.80) gives the results shown in Figures 2.4, 2.5, 2.6 and 2.7. For the transition probability in (2.83), $p = 0.03$ is used.

As it appears from Figure 2.4, the pure time-delay is tracked quite well by recursive estimation of 20 models in parallel. Since the method is based upon a forgetting factor approach, it is unable to catch some of the high frequency variation of the real delay. Hence the changes of the real delay are not discovered immediately when they occur, and narrow peaks are not discovered at all. Generally, the simulation studies show that the method is rather robust and reliable as far as slow variations are concerned.

The only difference between the method used in Figure 2.4 (method A) and the method used in Figure 2.5 (method B) is that the latter

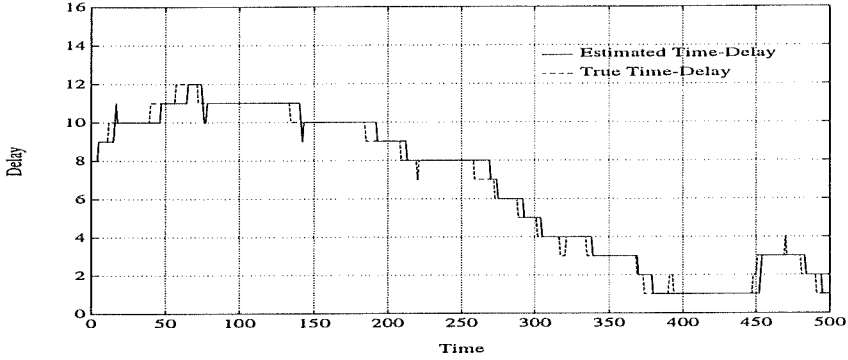


Figure 2.5: *Tracking a pure time-delay by estimation of models in parallel (modified). Simulation: Eqn. (2.80), (2.82), (2.83). Tracking: Eqn. (2.9)-(2.12) with $\lambda_\theta = \lambda_e = 0.9$ and (2.14) with $\kappa = 1$ and $\hat{k}(0) = 8$.*

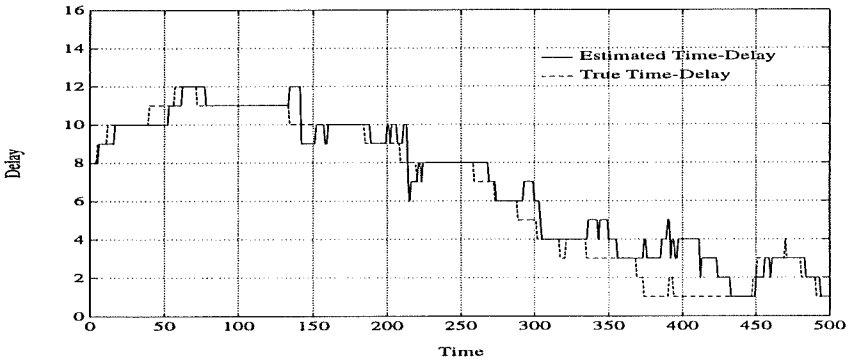


Figure 2.6: *Tracking a pure time-delay by Bányász and Keviczky's method. Simulation: Eqn. (2.80), (2.82), (2.83). Tracking: Eqn. (2.25), (2.26), (2.27), (2.28) with $\lambda = 0.9$ and $\hat{k}_0 = 8$.*

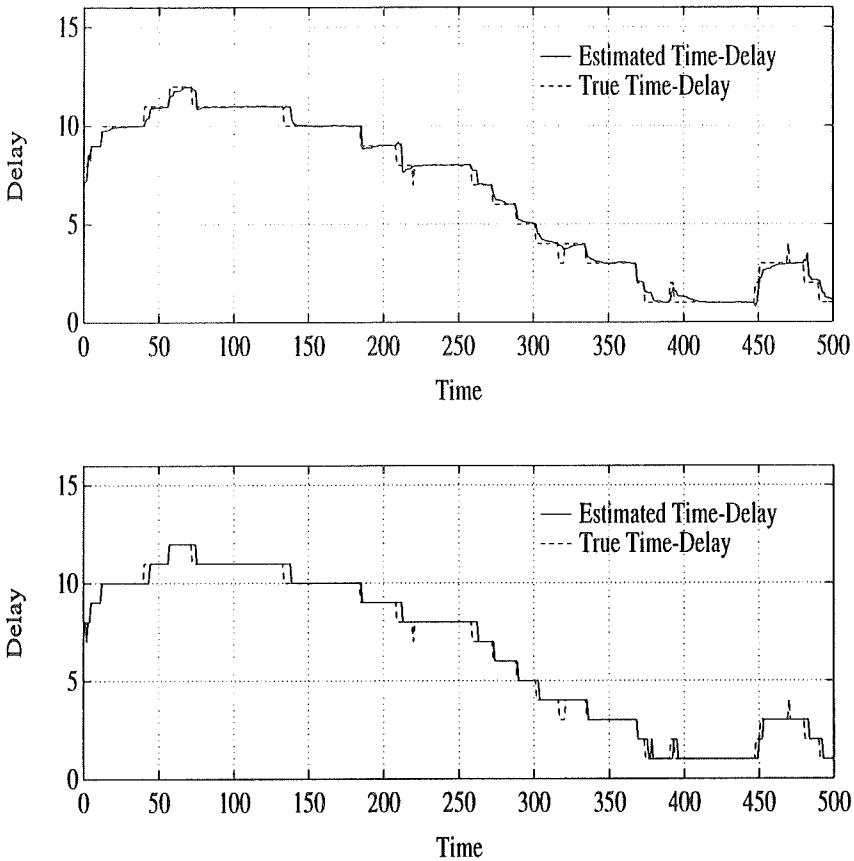


Figure 2.7: *Tracking a pure time-delay by estimation of a continuous time-delay. Simulation: Eqn. (2.80), (2.82), (2.83). Tracking: Eqn. (2.36)-(2.38) with $\lambda_\alpha = 0.8$, $\hat{k}(0) = 9$ and $\hat{\alpha}(0) = 0$. Top: $\hat{k}(t) + \hat{\alpha}(t)$ and the real delay vs. t . Bottom: $\hat{k}(t)$ and the real delay vs. t .*

only considers a subset of the models considered by the former. In this case, a subset of three models are estimated in parallel: One represents the current estimate of the delay, and the two others have delays which are one unit less and one unit greater than the current estimate, respectively. Because the real delay varies very slowly, the two methods give very similar results. In general, however, method A is much more robust than method B (but requires more computing). Very frequent changes of the real delay may cause a divergence of the estimate produced by method B. This also means that the method requires that the initial delay estimate is close to the real delay. Finally, the method is rather sensitive to the initialization values of the new models being introduced when the delay estimate changes.

It has not been possible to get the original Bányász-Keviczky algorithm in (2.24) and (2.25) to work properly. Therefore it has been modified by replacing (2.24) by (2.27). In this way the problems concerning negative definite Hessian matrices have been avoided. Figure 2.6 reveals that the algorithm is able to track an increasing delay. On the contrary, when the real delay decreases, the estimate very often increases a few samples later. This phenomenon seems to be one of the most serious drawbacks of the method.

Figure 2.7 shows that the continuous time-delay approach adapts very quickly to changes in the real delay. To illustrate how the algorithm works, the trajectory of the continuous estimate, $\hat{k}(t) + \hat{\alpha}(t)$, has been plotted. It can be seen that a change in the real delay results in an exponential transient of the continuous estimate. Through the simulation experiments it has been found that the method is fairly robust.

Tracking the Time-Delay of an ARX Model

Figures 2.8, 2.9 and 2.10 show the results obtained by recursive estimation of the time-delay and the parameters of the ARX model in (2.81). The variations of the delay have been simulated with a transition probability of $p = 0.02$ in (2.83).

First of all it should be noticed that no results from the Bányász-Keviczky algorithm are given. Bányász and Keviczky (1988) present an example of how to use their algorithm for estimation of the parameters and the time-delay of an ARMAX model. For the model in (2.81), however, an attempt to apply a similar algorithm fails totally because of divergence of the estimates.

As regards the time-delay, most of the comments which have earlier been made on Figures 2.4, 2.5 and 2.7 also apply to Figures 2.8, 2.9 and 2.10. Notice, however, that the method used in Figure 2.9 does not adapt to changes in the delay quite as fast as the method used in Figure 2.8. Notice also that the estimation of a continuous time-delay in Figure 2.10 results in a very quick adaption.

Concerning the process parameters, a , b_0 and b_1 , the estimates provided by the algorithm for estimation of models in parallel exhibit large fluctuations, while the continuous delay approach gives very smooth trajectories close to the real values.

The results and discussion above indicate that the estimation of a continuous delay makes up a method which is superior to the rest of the methods considered. A possible explanation for this is that the algorithm operates in two steps: The first step deals with the normal ARMAX parameters, and the second step deals with the time-delay. Furthermore, the forgetting factors used in those two steps may be different, thus permitting the frequency of the variation of the parameters and the delay to be different. The advantage of dividing

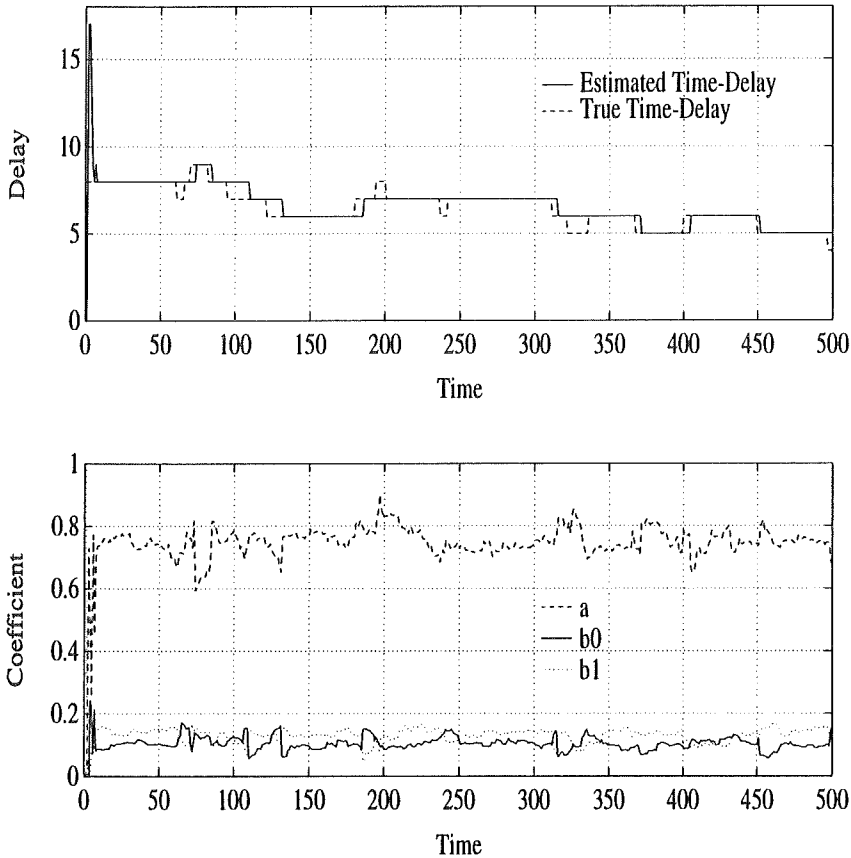


Figure 2.8: Recursive estimation of the time-delay and the parameters of an ARX model by estimation of models in parallel. **Simulation:** Eqn. (2.81), (2.82), (2.83). **Tracking:** Eqn. (2.9)-(2.12) with $k_{min} = 1$, $k_{max} = 20$, $\lambda_\theta = \lambda_e = 0.9$. **Top:** $\hat{k}(t)$ and the real delay vs. t . **Bottom:** \hat{a}_t , $\hat{b}_{0,t}$, $\hat{b}_{1,t}$ vs. t .

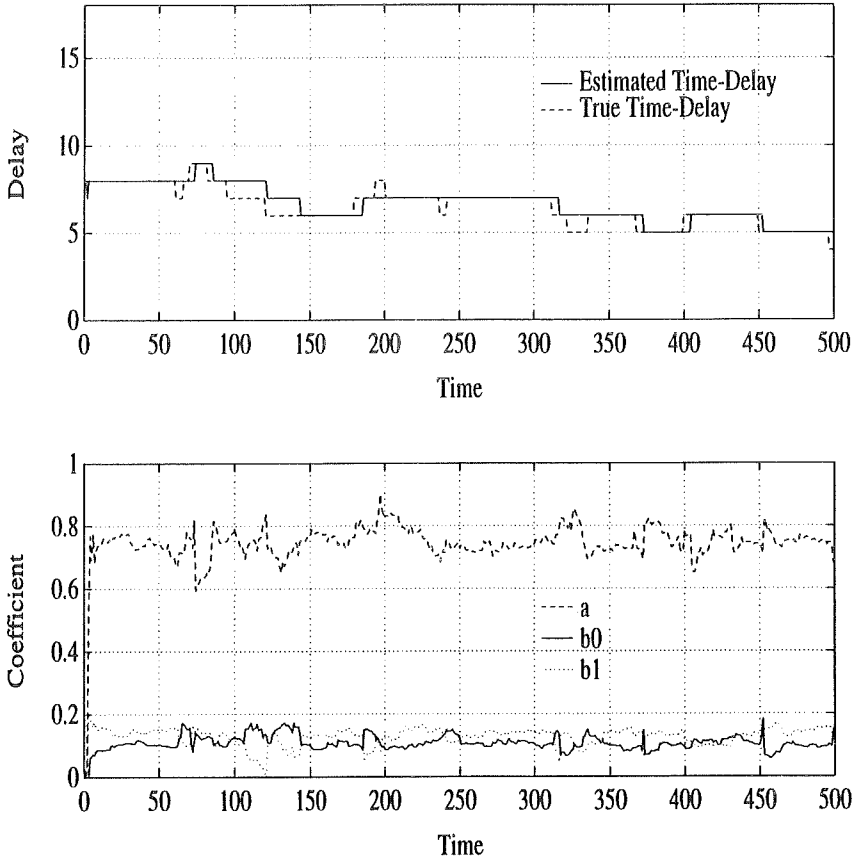


Figure 2.9: Recursive estimation of the time-delay and the parameters of an ARX model by estimation of models in parallel (modified). **Simulation:** Eqn. (2.81), (2.82), (2.83). **Tracking:** Eqn. (2.9)-(2.12) with $\lambda_\theta = \lambda_e = 0.9$ and (2.14) with $\kappa = 1$ and $\hat{k}(0) = 9$. **Top:** $\hat{k}(t)$ and the real delay vs. t . **Bottom:** \hat{a}_t , $\hat{b}_{0,t}$, $\hat{b}_{1,t}$ vs. t .

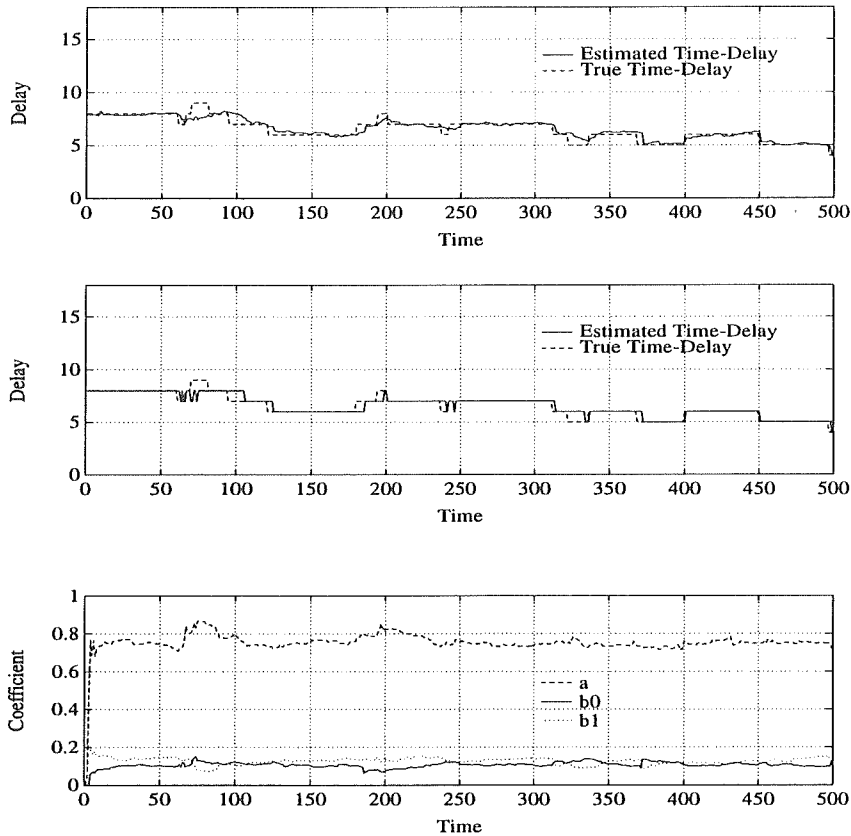


Figure 2.10: Recursive estimation of the time-delay and the parameters of an ARX model by estimation of a continuous time-delay. **Simulation:** Eqn. (2.81), (2.82), (2.83). **Tracking:** Eqn. (2.41)-(2.45) with $\lambda_\theta = 0.95$, $\lambda_\alpha = 0.85$, $\hat{k}(0) = 8$ and $\hat{\alpha}(0) = -0.1$. **Top:** $\hat{k}(t) + \hat{\alpha}(t)$ and the real delay vs. t . **Middle:** $\hat{k}(t)$ and the real delay vs. t . **Bottom:** \hat{a}_t , $\hat{b}_{0,t}$, $\hat{b}_{1,t}$ vs. t .

the estimation into two steps has become clear from the simulation experiments: Since the real process parameters are constant while the delay is time-varying, the chosen λ_θ should be higher than the chosen λ_α to obtain the best overall prediction ability.

2.3.2 Tracking a Time-Varying Time-Delay in a District Heating System

As mentioned in Section 2.1, there is an extensive need for on-line methods to track the transport times between different points in the distribution network of district heating systems. In order to perform on-line control, prediction or optimization of the system, it is extremely useful to know about the time distances from the heat production plant to various places in the network. One way to obtain this knowledge is to gather information from flow-meters placed throughout the network. Most frequently, however, such flow measurements are not available. Moreover, flow-meters are known to give rather uncertain measurements. Therefore, if there is only a small number of places to which the time distances should be tracked, it may be less expensive to use measurements of the supply temperature and statistical methods to establish an on-line estimation of the distances.

In this section, supply temperature data from the district heating system in Ishøj near Copenhagen in Denmark is used. Figure 2.11 shows the primary distribution network of this system. At the bottom of the figure the district heating plant is seen. The rest of the terminal points of the network are heat exchanger stations numbered 1 to 17. The figures in brackets (e.g. (100m)) are pipe lengths in meters, and the figures preceded by a 'ø' denote the inside diameter of the pipes in millimetres. The heat exchanger stations are the link between the primary and the secondary distribution networks.

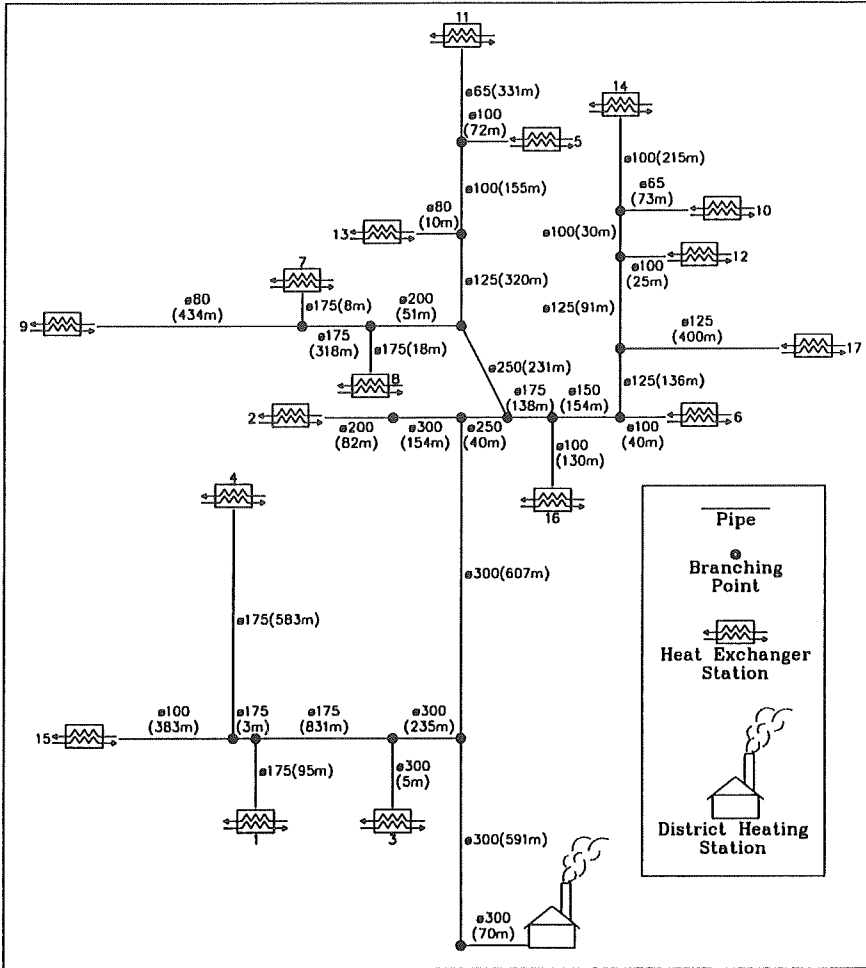


Figure 2.11: The distribution network of the district heating system in Ishøj, Denmark.

The district heating system in Ishøj was built in 1982, and today it supplies 5400 dwellings, five public schools and the city centre. The total length of the primary network is 7059 meters and all connected buildings are situated within an area of 1.2 km². At the district heating plant the heat is produced by three 17 MW coal-fired boilers, and by one smaller gas-fired boiler (see Benonysson (1991)).

Below special attention will be paid to heat exchanger station number 7. The total pipe length from the district heating plant to this heat exchanger station is 1916 m. Due to the transportation time of water over this distance, a changes in the supply temperature at the district heating plant does not immediately affect the supply temperature at heat exchanger station 7. After some delay, however, the change will occur. The time-delay of the “supply temperature signal” depends on the flow in the 8 different pipes through which the water must pass on its way to the heat exchanger station. This flow is among other things dependent on the heat load, the supply temperature and the temperature in return pipes at each of the 17 heat exchanger stations (not only at heat exchanger station 7!). Thus a “correct” physical based deterministic model of the variations of the time-delay should involve the entire distribution system.

However, the scope here is not to model the entire district heating system. (Furthermore such a model may be so comprehensive as regards numerical computation that it would be useless for on-line application.) The scope is to estimate the time-variation of the time-delay from the district heating plant to heat exchanger station 7 assuming that only measurements of the supply temperature at these two places are available. For that purpose transfer function models relating the supply temperature, $u(t)$, at the district heating plant to the supply temperature, $y(t)$, at heat exchanger station 7 are identified and estimated by means of the methods presented in Sections 2.2.1 and 2.2.2. The influence from the rest of the district heating system is considered to be noise, possibly correlated.

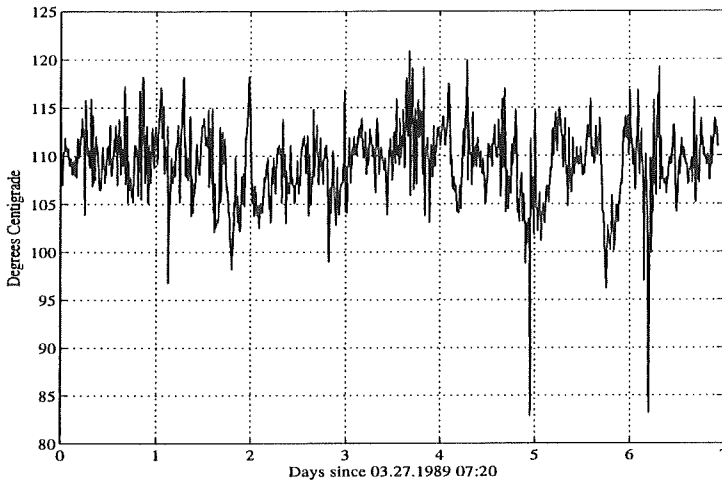


Figure 2.12: *The supply temperature at the district heating plant.*

Figures 2.12 and 2.13 show the supply temperature at the district heating plant and at heat exchanger station 7 in the period from 03.27.1989 07:20 to 04.03.1989 13:20 (2000 observations in each time series; sampling interval $T = 5$ min.). In the following these supply temperatures are denoted as input, $u(t)$, and output, $y(t)$, respectively.

Unfortunately, it has not been possible to get measurements of flows in the network from the same period. Therefore, the real time-delay from the district heating plant to heat exchanger station 7 cannot be computed. However, a table in Benonysson (1991) which shows the average mass flows in a very heavily loaded period (30 December 1985 - 8 January 1986) indicates that the minimum time-delay is about 28 minutes in this period. Furthermore, there is a figure in the same report which shows that the time-delay to heat exchanger 9 typically varies between 90 and 170 minutes (and that the minimum

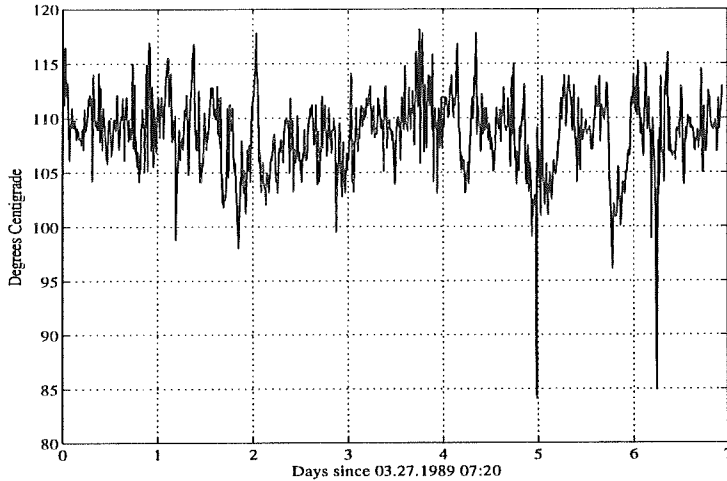


Figure 2.13: *The supply pipe temperature at heat exchanger station number 7.*

time-delay is reached in the morning between 7 and 7:30). Therefore, the time-delay to heat exchanger 7 which is about half the time-delay to heat exchanger 9 can be expected to vary between 45 and 85 minutes. However, this estimate of the maximum and minimum of the time-delay to heat exchanger 7 is rather uncertain and can only give an idea of how the variation has been in the period represented by Figures 2.12 and 2.13.

Estimating Models in Parallel

The recursive method based on parallel estimation of different models presented in Section 2.2.1 is tested with the above supply temperature data. The algorithm consisting of Equations (2.9), (2.10), (2.12) and (2.13) is employed. In order to use this method it is necessary to

specify n , m and r of the ARMAX(n, m, r) model, the forgetting factors, λ_θ and λ_ϵ , and the delay limits, k_{min} , k_{max} and Δk_{max} . As there is no information present on these parameters they are chosen so that the estimated variance of the one-step-ahead prediction error over the last 1800 observations is minimized,

$$\hat{\sigma}_\epsilon^2 = \min_{n, m, r, \lambda_\theta, \lambda_\epsilon, k_{min}, k_{max}, \Delta k_{max}} \left\{ \frac{1}{1800} \sum_{t=201}^{2000} \varepsilon_{\hat{k}(t-1)}^2(t) \right\}. \quad (2.84)$$

During the first 200 observations the recursive parameter estimates are stabilized; and after that point the influence of the initial values, $\hat{\theta}_p(0)$ and $P_p(0)$, $p = k_{min} \dots k_{max}$, are negligible. As an example, if $\lambda_\theta = 0.96$ then $\hat{\theta}_p(0)$ and $P_p(0)$ influence less than 0.03% on the criterion in (2.84) for $t > 200$.

Note that the objective function in (2.84) has a well defined minimum with respect to model order, (n, m, r) . This would not be the case if the model parameters were estimated as time-invariant parameters by an off-line estimation method. In this case the variance of the prediction error could be made arbitrarily small by increasing the model order. For recursive estimation methods, however, the variance of the prediction error is not minimized with respect to the traditional ARMAX parameters but with respect to the forgetting factors. In this method the ARMAX parameters should be considered as state variables in an extended model consisting of the original ARMAX model together with the recursive estimation algorithm. In this embedded model the actual parameters are the forgetting factors. This means that increasing the model order, (n, m, r) , does not imply an increase of the number of parameters in the extended model.

The minimization problem (2.84) is not easy to solve. One difficulty is due to the objective function being a function of discrete variables $(n, m, r, k_{min}, k_{max}$ and $\Delta k_{max})$. Another difficulty is due to the objective function being a discontinuous function of the continuous

variables λ_θ and λ_e (for a certain value of, say, λ_θ , an infinitesimal change of λ_θ may give a new trajectory of $\hat{k}(t)_{t=201, \dots, 2000}$ implying that the sequence of models being used for prediction changes. This in turn implies that the sum of squared prediction errors in (2.84) is subject to a non-infinitesimal change).

As a consequence of the problems mentioned above, it is difficult to find the global minimum. However, several local minima can be found; Table 2.1 shows the best one found in the present investigation. The corresponding trajectories of the estimated delay, $\hat{k}(t)$,

Table 2.1: *A local minimum of the minimization problem (2.84).*

n	m	r	λ_θ	λ_e	k_{min}	k_{max}	Δk_{max}	$\hat{\sigma}_e^2$
1	1	0	0.865	0.700	9 45 min.	21 105 min.	3 15 min.	1.023

parameters, $\hat{a}_1(t)$ and $\hat{b}_0(t)$, stationary gain, $\hat{G}(t) = \hat{b}_0(t)/(1 - \hat{a}_1(t))$, and residuals (one-step-ahead prediction errors) are plotted in Figures 2.14 to 2.18.

The estimated delay is seasonal with a period of 24 hours. The delay seems to be relatively low early in the morning (notice that 0, ..., 6 and 7 on the time axes of the graphs are at 07:20 in the morning). This can be explained by the fact that the heat consumption reaches its maximum in the morning and thereby implies a large flow and a short delay from the district heating plant to the consumer (the heat exchanger station). The last two days of observations were a weekend where the consumption is more irregular. In addition to the seasonality a long term change of the delay is seen.

The estimated time-delay is in reasonable accordance with the results from Benonysson (1991) which were mentioned previously – namely

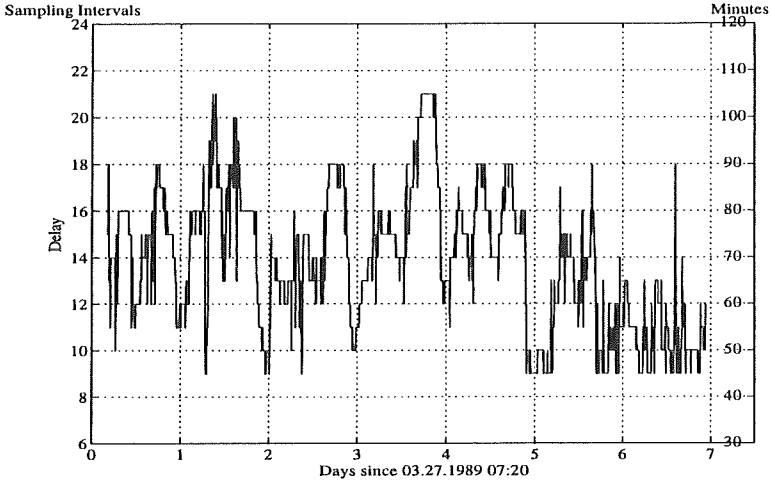


Figure 2.14: Trajectory of the estimated delay, $\hat{k}(t)$, corresponding to the solution of (2.84) given in Table 2.1.

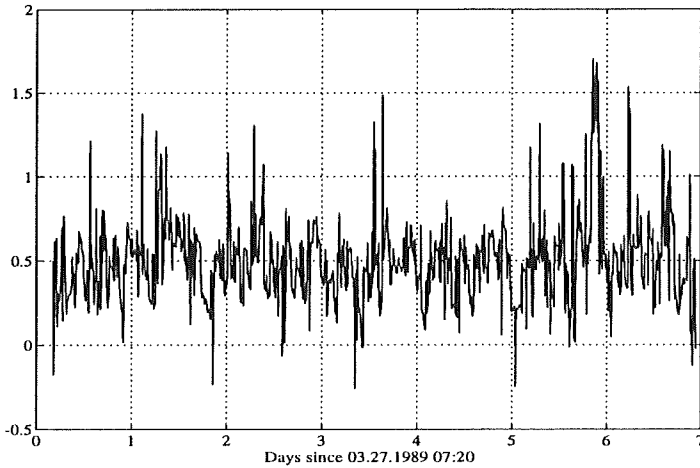


Figure 2.15: Trajectory of the estimated parameter $\hat{a}_1(t)$ corresponding to the solution of (2.84) given in Table 2.1.

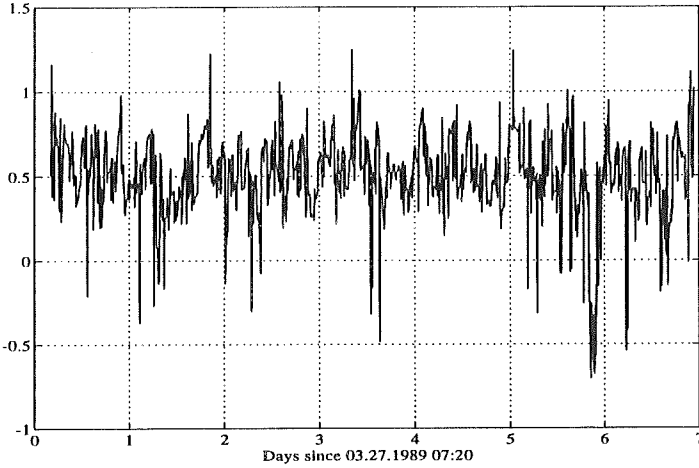


Figure 2.16: Trajectory of the estimated parameter $\hat{b}_0(t)$ corresponding to the solution of (2.84) given in Table 2.1.

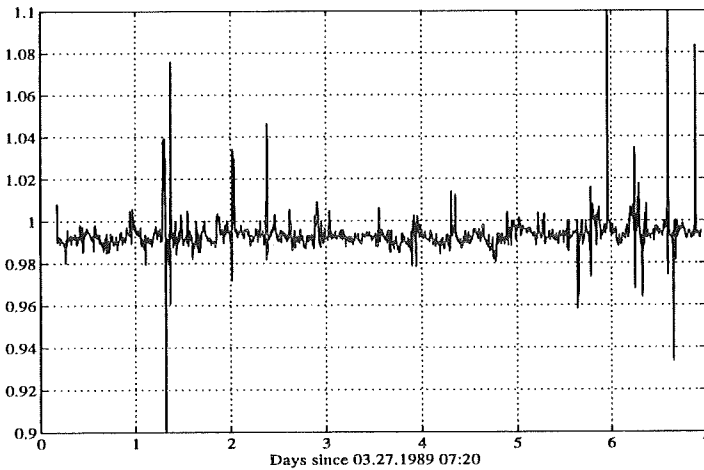


Figure 2.17: Trajectory of the estimated stationary gain, $\hat{G}(t)$, corresponding to the solution of (2.84) given in Table 2.1.

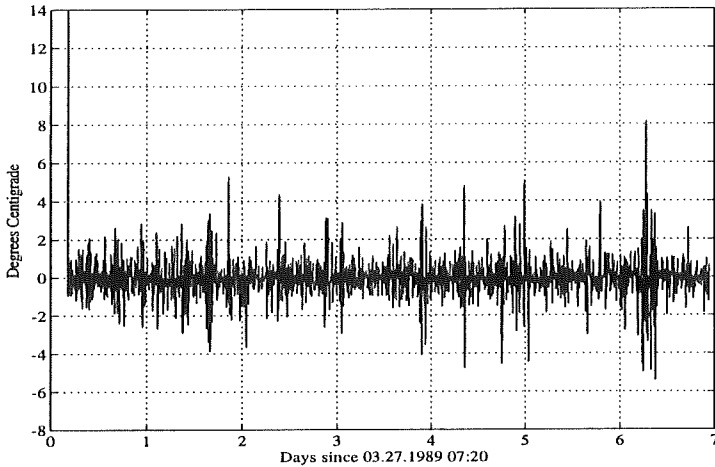


Figure 2.18: *Residuals corresponding to the solution of (2.84) given in Table 2.1.*

that the minimum of the time-delay is about 45 minutes and occurs in the morning between 7 and 7:30.

The optimal forgetting factors, λ_θ and λ_e , are quite small. Hence, the parameter estimates and the variance estimate are based only on $1/(1 - 0.865) = 7.41$ and $1/(1 - 0.700) = 3.33$ observations, respectively. In other words the dynamics of the system can be considered constant over a range of only about 7 sampling intervals (~ 35 min.). This indicates very small persistence of the system. The strong time-variation is confirmed by the estimated delay which varies over a considerable range within a few sampling intervals. Totally, the estimated time-delay varies from 9 to 21 sampling intervals – that is over a range of one hour. The time-variation also becomes apparent from the trajectories of the parameter estimates (see Figure 2.15 and 2.16).

When estimating the autocorrelation function of the one-step-ahead prediction errors, significant autocorrelations at the 5 % level are found at lag 2, 3, 4 and 5 (-0.161, -0.104, -0.074 and -0.058). This may be due to the algorithm being a little behind the real system dynamics because of the quick time-variation.

The stationary gain (Figure 2.17) seems to be fairly constant and very close to 1 (mean ≈ 0.993). Thus the water is approximately cooled 0.7% on its way from the district heating plant to heat exchanger station 7. The large peaks occurring every now and then are caused by large residuals and changes from one model (one value of the time-delay) to another.

A certain part of the variation of the estimated delay is presumably due to noise. In order to obtain a certain suppression of the noise, the input- and the output-observations are averaged in groups of three. That is, $\tilde{x}(\tilde{t}) = \frac{1}{3} \sum_{t=3\tilde{t}-2}^{3\tilde{t}} x(t)$, $\tilde{t} = 1, \dots, 667$, where $x(t)$ denotes $u(t)$ or $y(t)$, and $\tilde{x}(\tilde{t})$ is the corresponding series of average observations ($\tilde{u}(\tilde{t})$ or $\tilde{y}(\tilde{t})$). Thus $3 \times 667 = 2001$ original observations of both $u(t)$ and $y(t)$ are used; and the sampling interval is 15 minutes.

The algorithm consisting of Equations (2.9), (2.10), (2.12) and (2.13), i.e. the same as previously, is now applied with $u(t)$, $y(t)$ and t replaced by $\tilde{u}(\tilde{t})$, $\tilde{y}(\tilde{t})$ and \tilde{t} , respectively. A criterion corresponding to (2.84) is minimized:

$$\hat{\sigma}_e^2 = \min_{n, m, r, \lambda_\theta, \lambda_e, k_{\min}, k_{\max}, \Delta k_{\max}} \left\{ \frac{1}{600} \sum_{\tilde{t}=68}^{667} \varepsilon_{\hat{k}(\tilde{t}-1)}^2(\tilde{t}) \right\}. \quad (2.85)$$

A number of local minima of this criterion have been found, and the one resulting in the lowest value of the criterion is listed in Table 2.2. The trajectories of the estimated delay, $\hat{k}(\tilde{t})$, parameters, $\hat{b}_0(\tilde{t})$ and $\hat{b}_1(\tilde{t})$, and stationary gain ($\hat{G}(\tilde{t}) = \hat{b}_0(\tilde{t}) + \hat{b}_1(\tilde{t})$) are plotted in Figures 2.19 to 2.22.

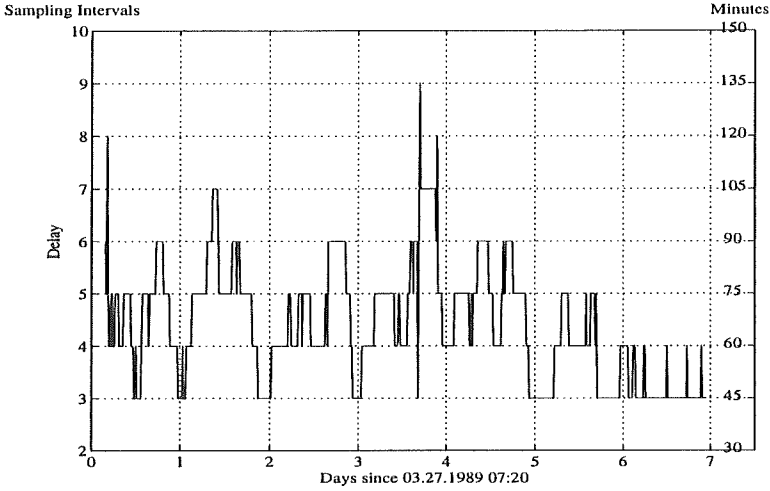


Figure 2.19: Trajectory of the estimated delay, $\hat{k}(\tilde{t})$, corresponding to the solution of (2.85) given in Table 2.2.

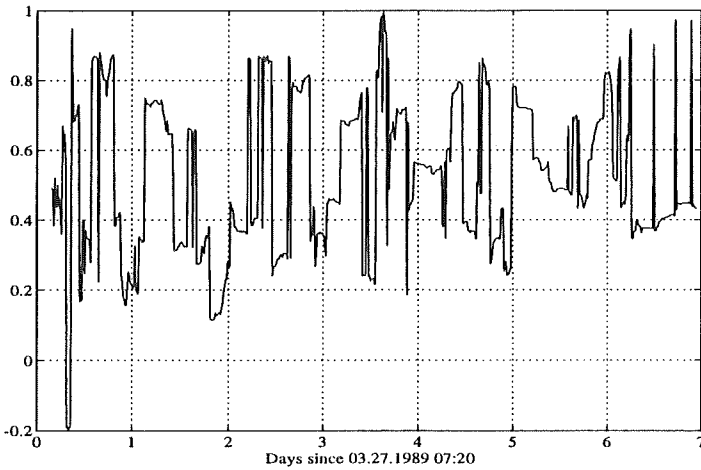


Figure 2.20: Trajectory of the estimated parameter $\hat{b}_0(\tilde{t})$ corresponding to the solution of (2.85) given in Table 2.2.

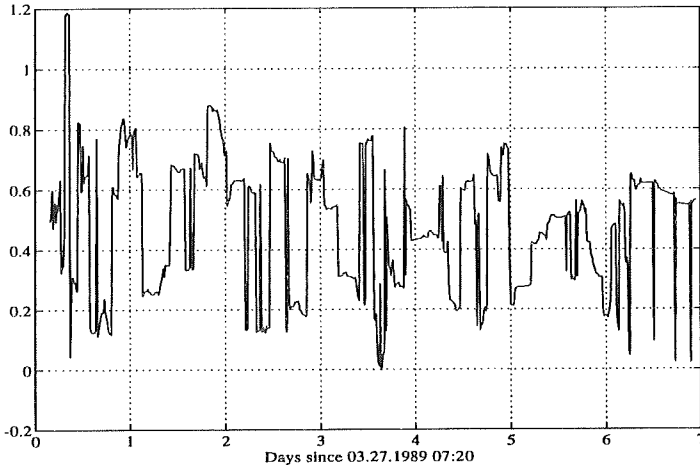


Figure 2.21: Trajectory of the estimated parameter $\hat{b}_1(\tilde{t})$ corresponding to the solution of (2.85) given in Table 2.2.

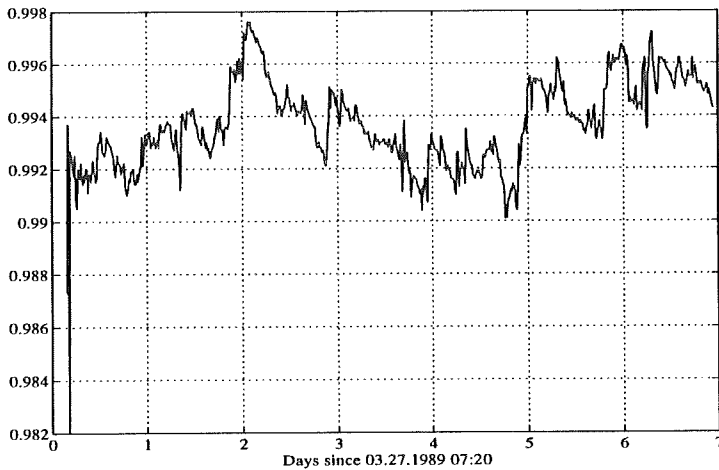


Figure 2.22: Trajectory of the estimated stationary gain, $\hat{G}(\tilde{t})$, corresponding to the solution of (2.85) given in Table 2.2.

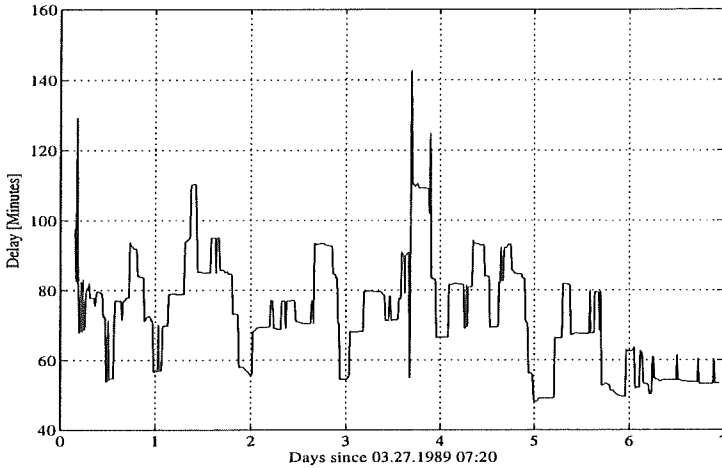


Figure 2.23: Trajectory of the estimated continuous delay, $\hat{\delta}(\tilde{t})$, corresponding to the solution of (2.85) given in Table 2.2.

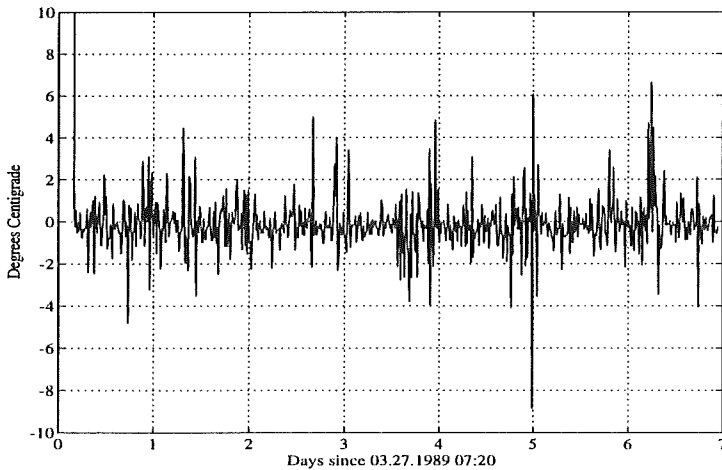


Figure 2.24: Residuals corresponding to the solution of (2.85) given in Table 2.2.

Table 2.2: A local minimum of the minimization problem (2.85).

n	m	r	λ_θ	λ_e	k_{min}	k_{max}	Δk_{max}	$\hat{\sigma}_e^2$
0	2	0	0.96	0.66	3 45 min.	9 135 min.	3 45 min.	1.583

Since the model

$$\tilde{y}(\tilde{t}) = b_0 \tilde{u}(\tilde{t} - k) + b_1 \tilde{u}(\tilde{t} - k - 1) + \tilde{e}(\tilde{t})$$

has two b -coefficients, an approximate continuous delay can be computed as

$$\hat{\delta}(\tilde{t}) = \hat{k}(\tilde{t}) + \hat{b}_1(\tilde{t}) / (\hat{b}_0(\tilde{t}) + \hat{b}_1(\tilde{t})). \quad (2.86)$$

Actually this approximation corresponds to the approximation introduced in Equation (2.65). The trajectory of $\hat{\delta}(\tilde{t})$ is plotted in Figure 2.23.

As anticipated, the estimated trajectory of the delay in Figure 2.19 is more smooth than the delay in Figure 2.14. The estimated continuous delay (Figure 2.23) seems to be even smoother. However, the seasonality and the long term changes of the delay remain the same as in Figure 2.14. For the stationary gain the smoothing effect of using average observations is very clear (compare Figure 2.17 with Figure 2.22). Actually the fact that the estimated stationary gain is almost constant while the parameters b_0 and b_1 varies very quickly suggests that the model should be re-parametrized as follows

$$\tilde{y}(\tilde{t}) = G((1 - \alpha)\tilde{u}(\tilde{t} - k) + \alpha\tilde{u}(\tilde{t} - k - 1)) + \tilde{e}(\tilde{t}).$$

Here the stationary gain G is considered constant and α is a time-varying parameter which corresponds to $b_1/(b_0 + b_1)$ and describes a part of the continuous time-delay (cf. Equation (2.86)).

The autocorrelation function of the one-step-ahead prediction error has a significant autocorrelation at lag 1 at the 5% level ($r(1) = -0.177$). This autocorrelation is probably due to problems arising from the very quick time-variation.

Continuous Time-Delay Approach

Application of the algorithm (2.41)-(2.45) to the district heating data requires a specification of n , m and r of the ARMAX(n, m, r) model, the forgetting factors, λ_θ and λ_α , and the limits, k_{min} , k_{max} and Δk_{max} , of the delay estimate. These parameters are chosen so that the estimated variance of the one-step-ahead prediction error over the last 1800 observations is minimized,

$$\hat{\sigma}_e^2 = \min_{n, m, r, \lambda_\theta, \lambda_\alpha, k_{min}, k_{max}, \Delta k_{max}} \left\{ \frac{1}{1800} \sum_{t=201}^{2000} \varepsilon_\theta^2(t) \right\}. \quad (2.87)$$

The remarks on the criterion in (2.84) also applies to (2.87).

A local minimum of (2.87) is listed in Table 2.3. The estimated trajectory of the delay corresponding to this solution is shown in Figure 2.25. The estimated delay (Figure 2.25) does not show the character-

Table 2.3: *A local minimum of the minimization problem (2.87).*

n	m	r	λ_θ	λ_α	k_{min}	k_{max}	Δk_{max}	$\hat{\sigma}_e^2$
1	2	2	0.97	0.85	9 45 min.	21 105 min.	2 10 min.	1.468

istic seasonal pattern seen in Figures 2.14 and 2.19. Moreover, in the latter half of the period the estimated delay stays at its lower limit (9 sampling intervals \sim 45 min.). The conclusion drawn from the

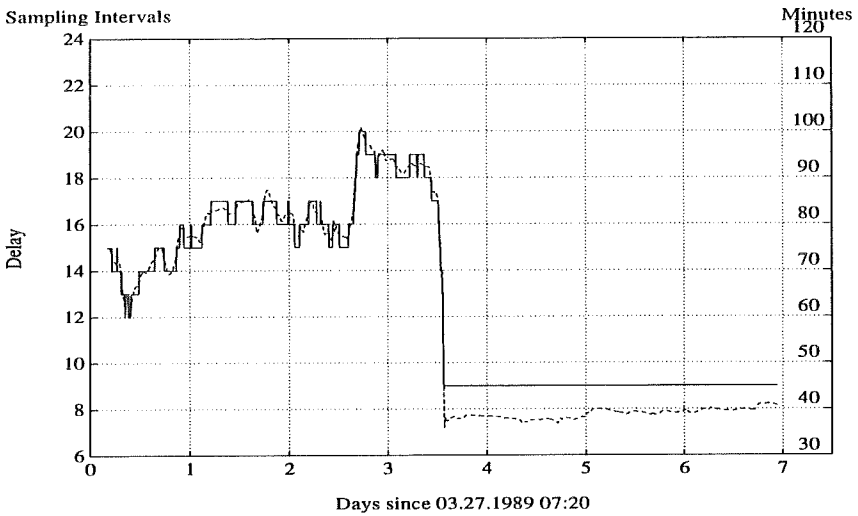


Figure 2.25: Trajectories of the estimated delay, $\hat{k}(t)$ (solid curve), and the estimated continuous delay, $\hat{\delta}(t) = \hat{k}(t) + \hat{\alpha}(t)$ (dashed curve), corresponding to the solution of (2.87) given in Table 2.3.

results is that the algorithm based upon continuous parametrization of the delay is not able to track the delay-variations of the present district heating data. The insufficiency of the algorithm is probably due to the time-variations of the delay being too quick and great.

The algorithm was also tested in Section 2.3.1 with great success. However, the investigations there were carried out with simulated data with slow variation. Furthermore, the delay changed only one sampling interval at a time at the most.

Seen as a whole, the above discussion indicates that the method is quite sensitive to the extent of non-stationarity of the system. Therefore, in on-line applications the algorithm should only be applied to slowly varying systems and perhaps in combination with a simple and robust technique that keeps the algorithm on track. For instance, an adaptively estimated cross-correlation function, $\hat{\rho}[u(t), y(t+l)]$, between input and output can provide a rough estimate of the delay as the lag, l , maximizing $\hat{\rho}[u(t), y(t+l)]$. Notice that such an estimate cannot stand alone since it may be affected by the phase shift introduced by dynamic relationship between the input and the output.

A Stochastic State-Space Model

In this section the state-space model presented in Section 2.2.2 is employed for modelling the district heating system (see Equations (2.67) and (2.68)). The extended Kalman filter (2.75)-(2.79) is applied in order to track the state variables of this model. The covariance matrix Σ_1 , the variance Σ_2 and the parameter ϕ are estimated by minimization of the estimated variance of the one-step-ahead prediction error over the last 1800 observations,

$$\hat{\sigma}_\epsilon^2 = \min_{\Sigma_1, \Sigma_2, \phi} \left\{ \frac{1}{1800} \sum_{t=201}^{2000} (y_t - \hat{y}_{t|t-1})^2 \right\}, \quad (2.88)$$

where Σ_1 is symmetric. This objective function has several local minima and furthermore it is discontinuous (due to (2.60)). Therefore it is not simple to find the global minimum. The lowest of the local minima found in the present investigation results in:

$$\Sigma_1 = \begin{pmatrix} 0.63 & -0.0023 & 0.010 & 0.56 \\ -0.0023 & 0.0077 & -0.0062 & -0.0028 \\ 0.010 & -0.0062 & 0.021 & 0.074 \\ 0.56 & -0.0028 & 0.074 & 0.62 \end{pmatrix} \quad (2.89)$$

$$\Sigma_2 = 1.65 \quad (2.90)$$

$$\phi = 0.84 . \quad (2.91)$$

This optimum corresponds to the following prediction error variance

$$\hat{\sigma}_e^2 = 1.758 . \quad (2.92)$$

Estimates of the non-zero elements in (2.66) can be obtained by comparing (2.89) with (2.70).

The number of different variance parameters estimated in the model is 11 (10 elements in Σ_1 plus Σ_2). However, from the innovations form of a state-space model with n state variables and s output variables it can be verified that the total number of identifiable variance parameters is $ns + \frac{1}{2}s(s+1)$ (see e.g. Ljung (1987)) – or for the present model $4 \times 1 + \frac{1}{2} \times 1 \times (1 + 1) = 5$. This means that the estimates of Σ_1 and Σ_2) given here represents only one out of an infinite number of possible solutions which result in the same sequence of one-step-ahead predictions of the output (i.e. the same value of the criterion in (2.88)) and the same estimate of ϕ .

Figures 2.26 to 2.30 show the reconstructed state variables ($\hat{x}_{t|t}$) corresponding to this solution of (2.88). Figure 2.31 shows the residual sequence.

By comparing Figure 2.26 with Figure 2.14 the trajectories of the estimated delay are found to be very similar. The Kalman filter

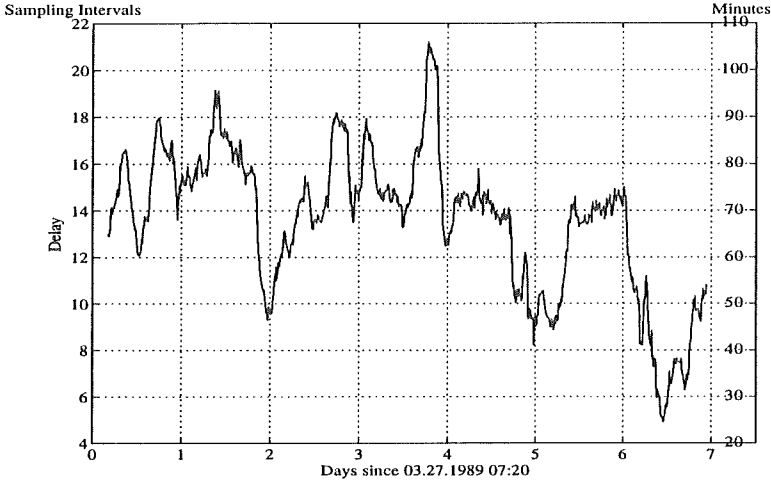


Figure 2.26: *The reconstructed time-delay, $\hat{\delta}_{i|t}$, corresponding to the solution of (2.88) given in (2.89) to (2.92).*

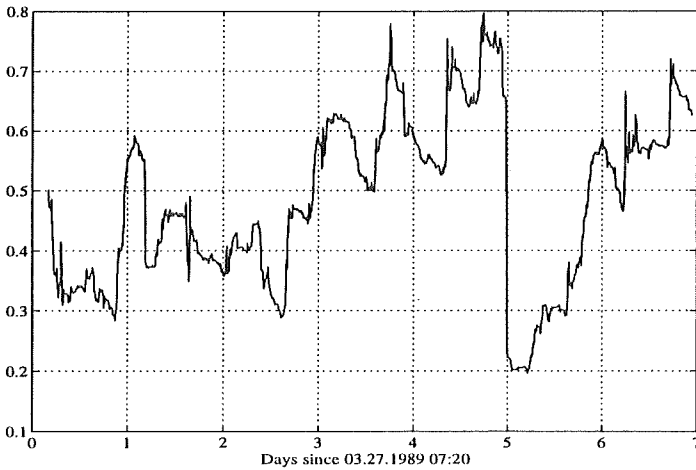


Figure 2.27: *The reconstructed a_t -parameter, $\hat{a}_{i|t}$, corresponding to the solution of (2.88) given in (2.89) to (2.92).*

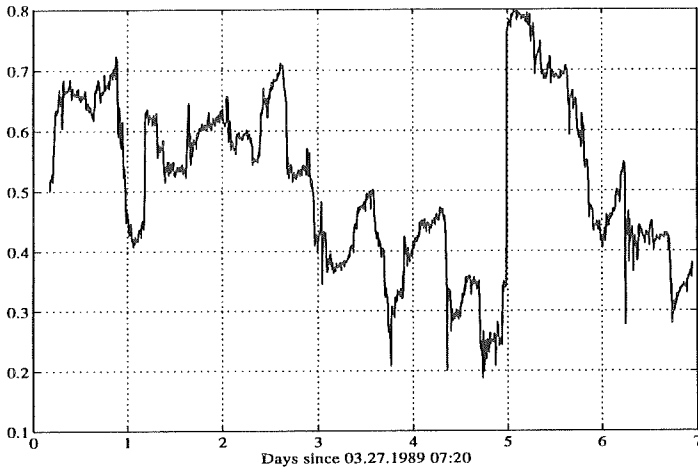


Figure 2.28: The reconstructed b_t -parameter, $\hat{b}_{t|t}$, corresponding to the solution of (2.88) given in (2.89) to (2.92).

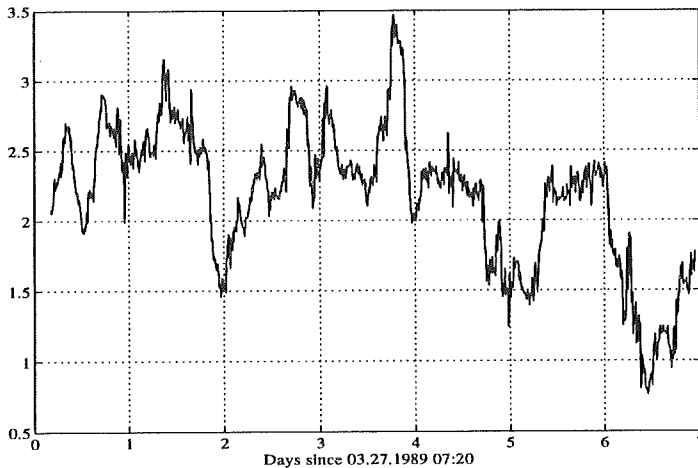


Figure 2.29: The reconstructed μ_t -state, $\hat{\mu}_{t|t}$, corresponding to the solution of (2.88) given in (2.89) to (2.92).

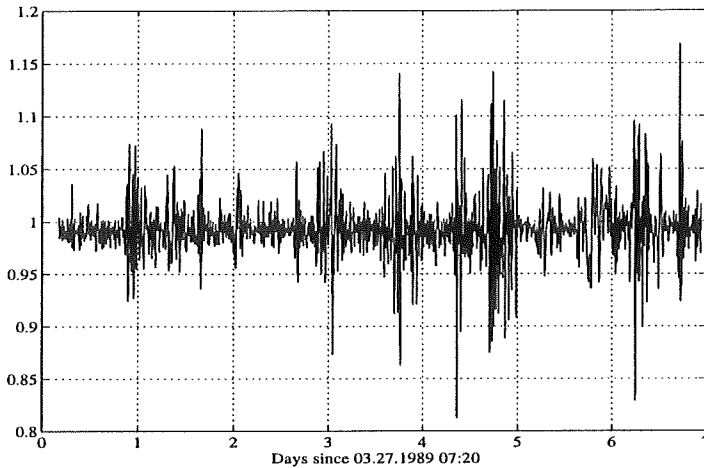


Figure 2.30: The reconstructed stationary gain, $\hat{G}_{t|t} = \hat{b}_{t|t}/(1 - \hat{a}_{t|t})$, corresponding to the solution of (2.88) given in (2.89) to (2.92).

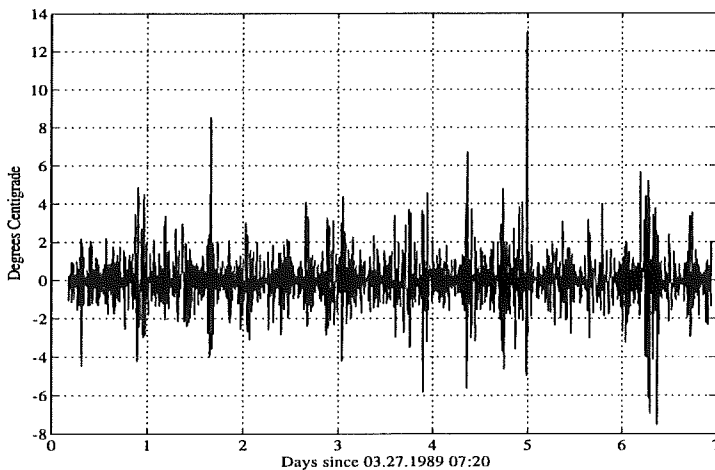


Figure 2.31: Residuals corresponding to the solution of (2.88) given in (2.89) to (2.92).

seems to miss some of the variation of the delay but the basic pattern is obvious. Even a weak indication of the 24 hour seasonality is seen. During half of the last day the estimated time-delay is less than the minimum delay estimated by the “parallel model” method (nine sampling intervals or 45 minutes). Actually it is likely that the true delay gets less than 45 minutes during the last day but the “parallel model” method was limited by a minimum value of 45 minutes in order to ensure that the overall prediction ability was maximized.

Due to large residuals about $t = 5.0$ the trajectories of the a_t and b_t parameters exhibit very drastic changes. The reason why the residuals become large is that peaks in the input and output signals occur at this time (see Figures 2.12 and 2.13). In this situation even a minor difference between the true and the estimated time-delay would imply large residuals.

The stationary gain (see Figure 2.30) shows very large fluctuations. The fluctuations may be due to the uncertainties of the estimated parameters, $\hat{a}_{t|t}$ and $\hat{b}_{t|t}$, but might as well indicate inadequacies in the model. The mean which is about 0.99 seems to be reasonable since a similar value was previously found by the “parallel model” method. But due to the fluctuations the estimated stationary gain exceeds one several times. This indicates a problem since the district heating system implies a reduction of the mean level of the temperature signal (cooling) and this in turn give rise to a stationary gain less than one. One way to solve the problem is to re-parametrize the model in (2.58) as follows,

$$y_t = a_t y_{t-1} + (1 - g_t^2)(1 - a_t)((1 - \alpha_t)u_{t-k_t} + \alpha_t u_{t-k_t-1}) + e_t .$$

The b_t parameter has been replaced by $G_t(1 - a_t)$ (using the relationship $G_t = b_t/(1 - a_t)$) and then the stationary gain, G_t , has been parametrized through g_t as $G_t = 1 - g_t^2$ in order to ensure that it remains less than one.

The structure of the model (2.58)-(2.64) was chosen in a very *ad hoc* manner. The order of the ARMAX model (2.58) seems reasonable in the light of the results found in Table 2.1, where the system is described by an ARMAX(1,1,0) model. The continuous delay approximation in (2.58) is introduced to achieve an estimate of the underlying continuous input process, $u_c(\tau)$ (also see Equation (2.65)). A more complex function of the u_t samples might give a better estimate and an improved overall description of the system. The continuous delay is modelled by (2.61) and (2.64) in order to describe both the stochastic nature and the sluggishness of the delay variations. However, other relevant models of the delay could easily be proposed.

The above discussion suggests that further investigations of state-space models ought to involve a structural analysis of the model in order to obtain a better description of the system. The potential improvement is confirmed by the estimated autocorrelation function of the one-step-ahead prediction errors since large correlations at lag 2, 3 and 4 are found. Furthermore it should be noted that the estimated variance of the prediction error ($1.758 \text{ }^\circ\text{C}^2$) is much larger than the variance obtained when models are estimated in parallel ($1.023 \text{ }^\circ\text{C}^2$). This also indicates that it might be possible to improve the prediction performance by a more suitable model structure.

However, although the structure of the model has not been optimized, the results obtained in this section illustrates the possibility of embedded stochastic models of the delay.

Benonysson (1991) tested two deterministic approaches for modelling the district heating system in Ishøj.

1. Pipe element model: Each pipe in the primary network is divided into a number of elements and for each of these elements dynamic heat loss equations are established and solved numerically. The solution of these equations is characterized by rela-

tively long calculation time compared to real time.

2. Node model: The temperature history at the inlet and outlet ends (nodes) of each pipe in the primary network is considered. The outlet temperature from a pipe is calculated from past inlet temperatures taking the rate of mass flow, the heat capacity of the steel pipe and heat loss from the pipe into account. The calculation time of this method is much shorter than for the element method.

Benonysson (1991) used real data (supply temperature time series from the district heating plant and mass flow time series from each pipe in the primary network) for simulation of the temperature variations at various heat exchanger stations in the system. For station 7 the standard deviation of the residual series was about 2.0 °C using the element model and 1.25 °C using the node model (sampling interval = 5 min. in both cases). The corresponding standard deviation obtained by estimating models in parallel was 1.01 °C and for the state space model the standard deviation was 1.33 °C. In other words, statistically identified time-varying transfer function models perform just as well or even better than the deterministic element and node models. This result is surprising since the latter models describe more details of the system than the former. However, the relatively bad performance of the deterministic models is probably due to accumulation of errors coming from the sub-models of the pipes or pipe elements.

2.3.3 Time-Delays in Models for Business Cycle Forecasting

In business cycle forecasting, dynamic input-output models are often used for prediction of the business cycle reference series. In

these models the predictor is a function of lagged values of input variables (leading indicators). In most cases, economic theory does not provide the exact lag structure, and identification of the dynamic structure is therefore usually based upon empirical approaches. Box and Jenkins (1976) proposed a pre-whitening method for identification of fixed lags in transfer function models. This also includes the ARMAX model used previously in this chapter.

The Box-Jenkins transfer function model assumes a time-invariant dynamic relationship between input and output. In practice, however, this assumption is frequently violated since both the time-delay (lead) and the dynamics can exhibit variation over time. This is the case for the Swedish business cycle data studied in this section. (see also Edlund and Sjøgaard (1993)).

The dependent variable of the data set is the Swedish index of industrial production which is used as an indicator of the business cycle. By identification and estimation of transfer function models, dynamic relationships between this variable and leading indicators will be established. First the time-invariant transfer function models will be stated and then the time-variation of the models will be estimated by the algorithm from Section 2.2.1 which performs recursive estimation of a collection of models in parallel. The results below were originally presented in Edlund and Sjøgaard (1993).

The Business Cycle Data Set

Two leading financial indicators which are part of the set of time series used by the OECD to forecast the Swedish business cycle are used. The data set covers the period January 1970 - September 1988 (225 monthly observations) and is collected from the Main Economic Indicators database at Statistics Sweden. The following two indi-

cators are used as explanatory variables for the index of industrial production (IIP), $Y(t)$:

$U_1(t)$ = Money supply, M1, with minor adjustments and deflated by the consumer price index.

$U_2(t)$ = Yield of long term government bonds.

The time series are shown in Figures 2.32-2.34. The IIP series has been adjusted for a major strike in Sweden in May 1980: the value for May 1980 has been replaced by an average of the May values from 1979 and 1981. In this way the IIP value for May 1980 has been treated as a missing observation. An alternative and statistically more correct way of replacing this missing observation would for instance be to apply the EM algorithm (Shumway (1988), Shumway and Stoffer (1982)). The observation from May 1980 could, as another alternative, be treated as an additive outlier. In this case the influence of the outlier could be suppressed by using robust estimation (Martin and Yohai (1985), Sejling (1993)).

In the models, the following transformed variables are used:

$$u_1(t) = \nabla_{12} \nabla \ln U_1(t) \quad (2.93)$$

$$u_2(t) = \nabla_{12} \nabla \ln U_2(t) \quad (2.94)$$

$$y(t) = \nabla_{12} \nabla \ln Y(t) . \quad (2.95)$$

These transformations remove seasonality and trend from the series, and the variance is made more homogeneous in time due to the logarithm function.

According to the OECD (1987) the money supply, $U_1(t)$, can be seen as an “expansionary accommodating stance of monetary policy”. Therefore, it is reasonable to assume that an increase in $U_1(t)$

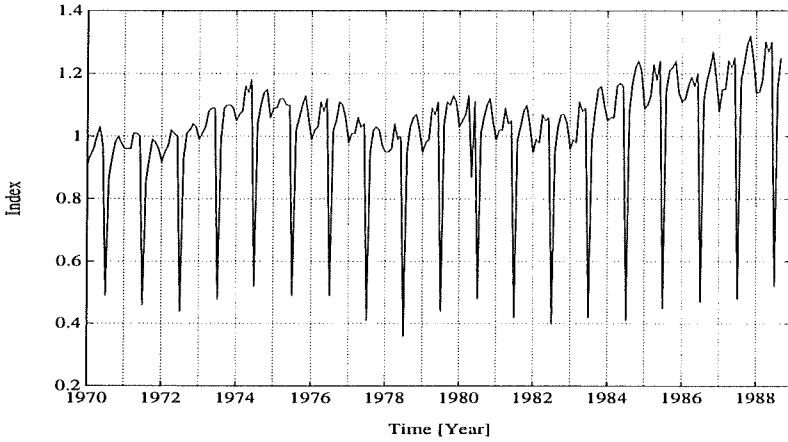


Figure 2.32: *Swedish index of industrial production (IIP). Original series January 1970 - September 1988.*

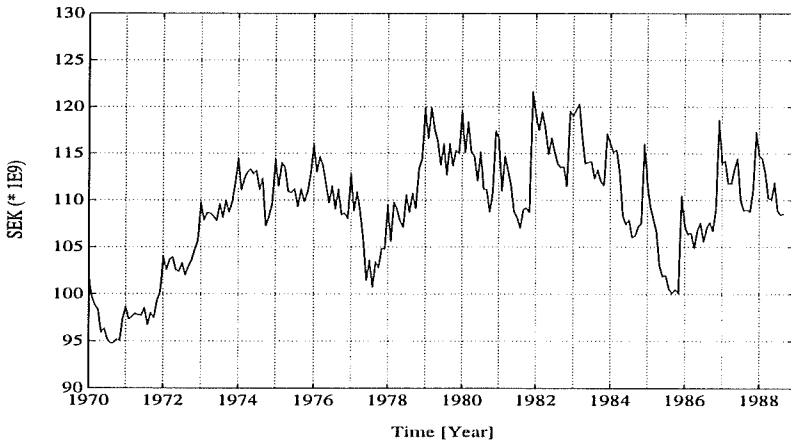


Figure 2.33: *Money supply, M1. Deflated series January 1970 - September 1988.*

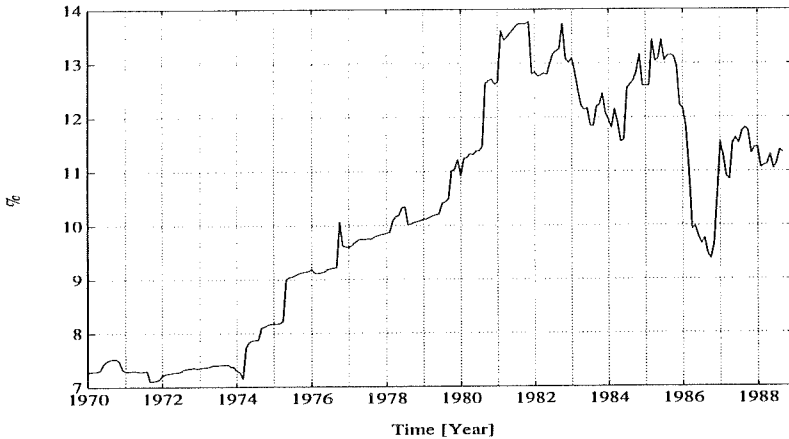


Figure 2.34: Yield of long term government bonds. Original series January 1970 - September 1988.

is followed by an increase in the IIP, $Y(t)$. The yield variable, $U_2(t)$, is an indicator of “the stimulus to consume or invest” and an increase in $U_2(t)$ should lead to a decrease in $Y(t)$.

Time-Invariant Models

The multiple-input single-output version of the Box and Jenkins (1970) transfer function model is considered,

$$y(t) = c + \sum_{i=1}^m H_i(q^{-1})u_i(t - k_i) + \frac{C(q^{-1})}{D(q^{-1})}e(t), \quad (2.96)$$

where $H_i(q^{-1}) = \frac{B_i(q^{-1})}{A_i(q^{-1})} = h_{i0} + h_{i1}q^{-1} + h_{i2}q^{-2} + \dots$ ($i = 1, \dots, m$). Here h_{ij} is the impulse response weight for input i at lag j . This is a generalization of the ARMAX model in (2.1) since a constant c and

the possibility of having multiple input and different poles for the input and the noise ($A_i(q^{-1}) \neq D(q^{-1})$) have been introduced.

Two types of transfer function models of the transformed IIP variable, $y(t)$, are considered: One type only contains the transformed M1 variable, $u_1(t)$, as an input, and another type contains both $u_1(t)$ and the transformed yield of long term government bonds, $u_2(t)$, as input variables. A two-step ridge regression method is used to identify noise models and estimate the impulse response weights (see e.g. Edlund (1984)). For the single-input case, the following approximation of the transfer function model in (2.96) is considered:

$$y(t) = c + (h_0 + h_1q^{-1} + \dots + h_Mq^{-M})u(t) + n(t), \quad (2.97)$$

where $n(t) = \frac{C(q^{-1})}{D(q^{-1})}e(t)$. In this approximation it has been assumed that the impulse response weights, h_j , are ≈ 0 for lags $j > M$, where M is a suitable large positive integer. Since the time-delay (k_i) in (2.96) is unknown all the impulse response weights are estimated although some of them may turn out to be zero.

The two-step ridge regression method goes through the following steps:

1. Estimate c and the impulse response weights, h_1, h_2, \dots, h_M , in (2.97) by ridge regression, and compute the residuals, $\{\hat{n}(t)\}$, of this model. If the residual series is white noise then stop. Otherwise identify and estimate a noise model,

$$\hat{n}(t) = \frac{\hat{C}(q^{-1})}{\hat{D}(q^{-1})}\hat{e}(t).$$

Use this model as a filter for the input and output series:

$$\begin{aligned} \hat{C}(q^{-1})y'(t) &= \hat{D}(q^{-1})y(t) \\ \hat{C}(q^{-1})u'(t) &= \hat{D}(q^{-1})u(t). \end{aligned}$$

2. Replace $u(t)$ and $y(t)$ in (2.97) by $u'(t)$ and $y'(t)$, respectively, and go to 1.

The two-step ridge regression method can easily be extended to the multiple-input case.

The reason why ridge regression is used for estimation of the impulse response weights is that there is multicollinearity between input values with different lags in (2.97) (the input series is autocorrelated). For one-input models the method can be regarded as an alternative to the pre-whitening cross-correlation approach proposed by Box and Jenkins (1976).

The model in (2.97) belongs to the output error model structure where $n(t)$ is the output error, and in step one of the two-step ridge regression method the impulse response weights are estimated using the output error method. This method and especially its convergence properties are discussed by Söderström and Stoica (1989).

Single-Input: Money Supply. For the model with only $u_1(t)$ as input, the following noise model is obtained by the two-step ridge regression method (using AUTOBOX PLUS, version 2.0 from Automatic Forecasting Systems Inc.):

$$(1+0.306q^{-12})\hat{n}(t) = (1-0.651q^{-1})(1-0.714q^{-24})\hat{e}(t), \quad \hat{\sigma}_e = 0.0285. \quad (2.98)$$

The estimated impulse response weights suggest two alternative lag structures – one with lead time (time-delay k) 6 and one with lead time 13 (see Edlund and Sjøgaard (1993)). Therefore the two follow-

ing models are estimated (using AUTOBOX):

Model 1:

$$y(t) = 0.326u_1(t - 6) + n(t) ,$$

where

$$(1 + 0.299q^{-12})\nabla_{12}n(t) = (1 - 0.660q^{-1})(1 - 0.564q^{-24})e(t) ,$$

and

$$\hat{\sigma}_e = 0.02976 .$$

Model 2:

$$y(t) = 0.277u_1(t - 13) + n(t) ,$$

where

$$(1 + 0.335q^{-12})\nabla_{12}n(t) = (1 - 0.641q^{-1})(1 - 0.632q^{-24})e(t) ,$$

and

$$\hat{\sigma}_e = 0.02929 .$$

All parameter estimates are significant (at the 5% level), and for both models the sample autocorrelation function of the residuals shows no significant correlations.

Double-Input: Money Supply and Yield of Long Term Government Bonds. For the model which has both $u_1(t)$ and $u_2(t)$ as input, the following noise model is obtained (using AUTOBOX):

$$(1 + 0.324q^{-12})\hat{n}(t) = (1 - 0.669q^{-1})(1 - 0.707q^{-24})\hat{e}(t) , \quad \hat{\sigma}_e = 0.0275 , \quad (2.99)$$

The estimated impulse response weights suggest two alternative lag structures in this case too, namely $(k_1, k_2) = (6, 16)$ or $(6, 21)$ (see Edlund and Sogaard (1993)). The following two models are esti-

mated (using AUTOBOX):

Model 1:

$$y(t) = 0.281u_1(t - 6) - 0.107u_2(t - 16) + n(t) ,$$

where

$$(1 + 0.322q^{-12})\nabla_{12}n(t) = (1 - 0.683q^{-1})(1 - 0.693q^{-24})e(t) ,$$

and

$$\hat{\sigma}_e = 0.02896 .$$

Model 2:

$$y(t) = 0.382u_1(t - 6) - 0.146u_2(t - 21) + n(t) ,$$

where

$$(1 + 0.380q^{-12})\nabla_{12}n(t) = (1 - 0.671q^{-1})(1 - 0.628q^{-24})e(t) ,$$

and

$$\hat{\sigma}_e = 0.02828 .$$

All parameter estimates are significant (at the 5% level), and for both models the sample autocorrelation function of the residuals shows no significant correlations.

Time-Varying Models

Single-Input: Money Supply. The original M1 and IIP series are not directly used for estimation of the time-variation. Instead the series filtered by the noise model (2.98) are used. By use of filtered variables, a model without seasonality is approximately obtained. Furthermore, the noise included in the model is approximately white. According to the results of the two-step ridge regression method, a suitable model structure would be

$$y'(t) = bu'_1(t - k) + e(t) ,$$

Table 2.4: Algorithm parameters and residual standard deviation for the model $IIP=f(M1)$.

k_{min}	k_{max}	Δk_{max}	λ_θ	λ_e	$\hat{\sigma}_e$
5	9	3	0.93	0.97	0.03150

where $\{y'(t)\}$ and $\{u'_1(t)\}$ are the filtered series and $\{e(t)\}$ is approximately white noise.

The continuous parameter b and the delay parameter k are estimated recursively using a slightly modified version of the algorithm consisting of Equations (2.9), (2.10), (2.12) and (2.13). The modification is introduced to ensure that, if possible, $\hat{b}(t)$ is positive because the sign of b is expected to be positive due to the reasons mentioned previously (increasing money supply leads to increasing IIP). The modification concerns Equation (2.13): at time t ($t = 1, \dots, 225$) the model with minimum prediction error variance among the models with a positive estimate of b is selected. If no such model is present, the model with the largest estimate of b is chosen.

The algorithm parameters, k_{min} , k_{max} , Δk_{max} , λ_θ and λ_e , are found by minimization of the estimated variance of the one-step-ahead prediction error during the period January 1970 - September 1988 (To let the estimates of b and k stabilize before January 1970, data from 1960 and onwards are used):

$$\hat{\sigma}_e^2 = \min_{k_{min}, k_{max}, \Delta k_{max}, \lambda_\theta, \lambda_e} \left\{ \frac{1}{225} \sum_{t=1}^{225} \varepsilon_{k(t-1)}^2(t) \right\} .$$

A local minimum is shown in Table 2.4. The series of residuals (one-step-ahead prediction errors) shows a significant autocorrelation at lag 24. This indicates inadequacies of the model. The problem

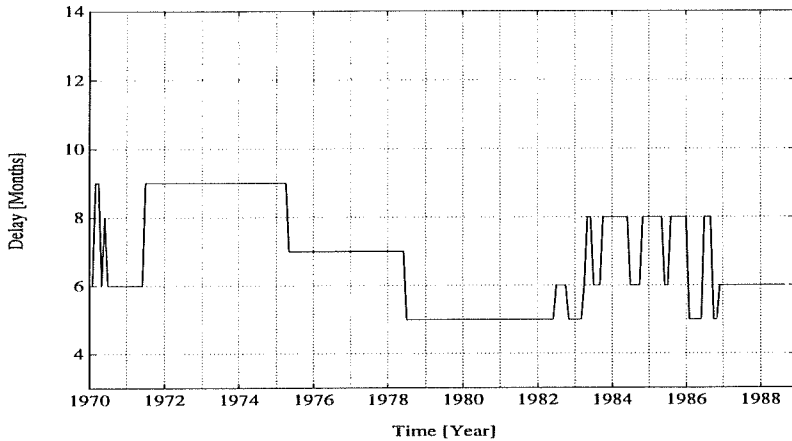


Figure 2.35: *Estimated delay between deflated money supply ($M1$) and the Swedish index of industrial production (IIP) using the the model $IIP=f(M1)$.*

may be due to the noise model used to filter the u and y variables. This noise model is a fixed parameter model developed to the time-invariant case. Therefore the results might be improved by including the noise model in the recursive estimation.

The trajectories of $\hat{k}(t)$ and $\hat{b}(t)$ are shown in Figures 2.35 and 2.36. Note that $\hat{b}(t)$ is very close to zero in the beginning of 1970 and in a period from 1982 to 1987. This indicates that there is a bad correlation between $U_1(t)$ and $Y(t)$ during those two periods. Consequently it is very difficult to estimate the lead time (see Figure 2.36).

Double-Input: Money Supply and Yield of Long Term Government Bonds. For this double-input case, an approach similar to the one used for the single-input case is applied. Hence the IIP , $M1$

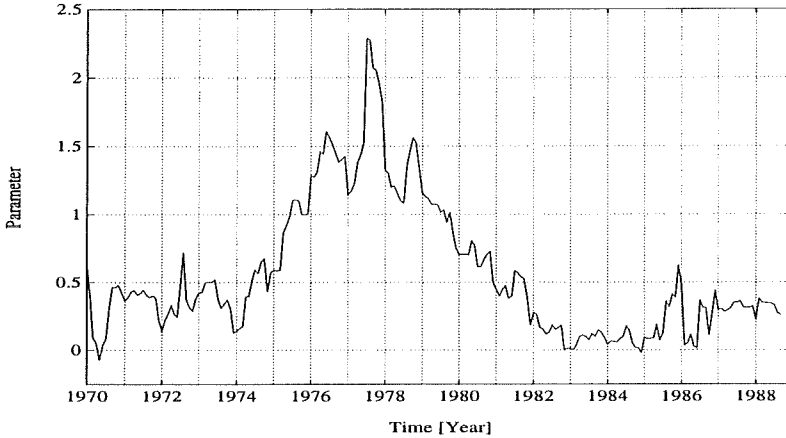


Figure 2.36: *Estimated trajectory of the parameter b .
Model: $IIP=f(M1)$.*

and Yield series are filtered by the noise model in (2.99). Denoting the filtered variables $y'(t)$, $u'_1(t)$ and $u'_2(t)$, the model can be written as

$$y'(t) = b_1 u'_1(t - k_1) + b_2 u'_2(t - k_2) + e(t),$$

i.e. same model structure as in the time-invariant case.

Since the model contains two input variables, the parallel-models algorithm used for the single-input model is extended correspondingly. This extension implies an estimation of $(k_{1,max} - k_{1,min} + 1) \times (k_{2,max} - k_{2,min} + 1)$ models in parallel. In order to achieve the expected signs of $\hat{b}_1(t)$ and $\hat{b}_2(t)$, the selection of a model at time t ($t = 1, \dots, 225$) is made according to the following rule: the model with minimum prediction error variance among those having a positive estimate of b_1 and a negative estimate of b_2 is selected. If no model with this property is found, the model with minimum prediction error variance in the entire set of models is selected (b_1 and b_2 unconstrained).

Table 2.5: Algorithm parameters and residual standard deviation for the model $IIP=f(M1, Yield)$.

$k_{1,min}$	$k_{1,max}$	$\Delta k_{1,max}$	$k_{2,min}$	$k_{2,max}$	$\Delta k_{2,max}$
5	9	2	16	16	1
λ_θ	λ_e	$\hat{\sigma}_e$			
0.96	0.85	0.03050			

The algorithm parameters, $k_{1,min}$, $k_{1,max}$, $\Delta k_{1,max}$, $k_{2,min}$, $k_{2,max}$, $\Delta k_{2,max}$, λ_θ and λ_e , are found by minimization of the estimated variance of the one-step-ahead prediction error during the period January 1970 - September 1988. The local minimum shown in Table 2.5 is found.

The autocorrelation function of the residuals has significant values at lag 24 and lag 48. The explanation for this is probably the same as for the single-input model, namely that the noise model is not adaptive.

Plots of $\hat{k}_1(t)$, $\hat{k}_2(t)$, $\hat{b}_1(t)$ and $\hat{b}_2(t)$ are shown in Figures 2.37-2.40. Comparing Figure 2.37 with Figure 2.35 shows that the estimated M1-IIP delay is very much alike in the two cases. However, the estimation of the delay during the years 1981 - 1987 is difficult because of low M1-IIP gain (see Figure 2.39). The Yield-IIP delay shows no variation at all.

Comparison of Time-Invariant and Time-Varying Models

The results of the recursive estimation method agree fairly well with those of the two-step ridge regression method. A summary of the

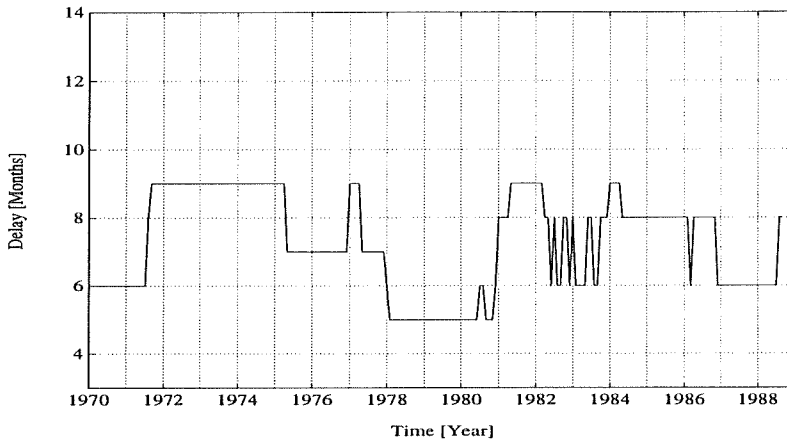


Figure 2.37: *Estimated delay between deflated money supply (M1) and Swedish index of industrial production (IIP) when using the the model $IIP=f(M1, Yield)$.*

results is shown in Table 2.6. In this table the delays shown for the recursive estimators are averages during the period January 1970 - September 1988 (the min-max ranges are shown in parenthesis). The stationary gains for the recursive estimators are averages as well (for the single-input model the stationary gain is b , and for the double-input model the stationary gains are b_1 and b_2).

The estimated M1-IIP delay from the recursive estimators varies from 5 to 9 months while the two-step ridge regression method result in 6 or 13 months (in the single-input case). This indicates that the adaptive model should be compared to the time-invariant model having delay 6. The reason why the long delay is not represented in the time-varying model is possibly that the variation of the true delay is seasonal, and since the number of observations within a season is only 12 the recursive estimation procedure does not adapt quick

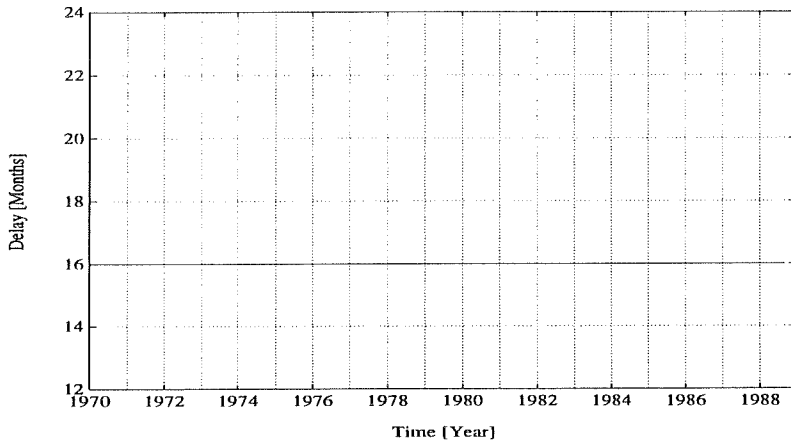


Figure 2.38: *Estimated delay between yield of long term government bonds (Yield) and Swedish index of industrial production (IIP) when using the the model $IIP=f(M1, Yield)$.*

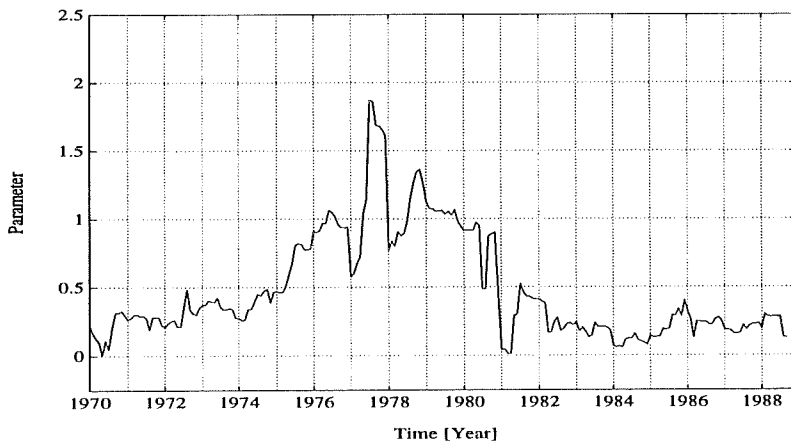


Figure 2.39: *Estimated trajectory of the parameter b_1 . Model: $IIP=f(M1, Yield)$.*

Table 2.6: *Summary of the results of the two-step ridge regression estimator and the recursive estimators.*

	Two-Step Ridge Regression	Recursive Estimation
Money Supply		
Time-delay	6 or 13	6.79 (5-9)
Gain	0.326 or 0.277	0.567
Standard Error	0.0298 or 0.0293	0.0315
Money Supply and Yield of Bonds		
<i>Money supply:</i>		
Time-delay	6	7.25 (5-9)
Gain	0.281 or 0.382	0.477
<i>Yield of Bonds:</i>		
Time-delay	16 or 21	16
Gain	-0.107 or -0.146	-0.175
Standard Error	0.0290 or 0.0283	0.0305

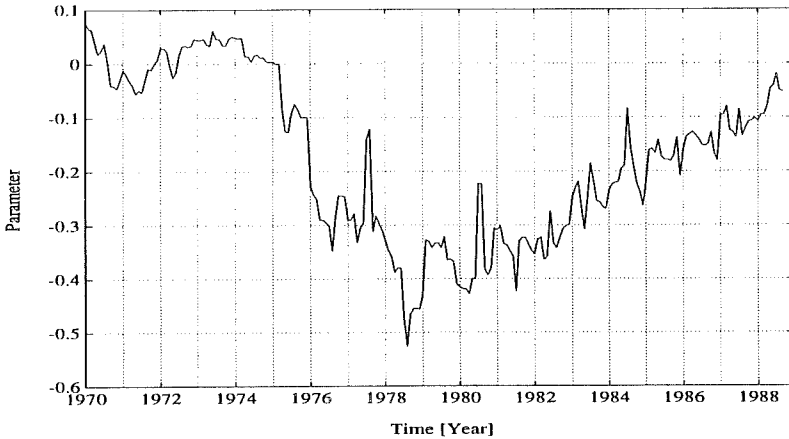


Figure 2.40: *Estimated trajectory of the parameter b_2 .
Model: $IIP=f(M1, Yield)$.*

enough to track the whole range of delay values. A clear indication of the presence of seasonality in the data has been obtained through experiments with quick adapting models.

The constancy of the estimated Yield-IIP delay gives reason to believe that there is no long term variations. However, experiments with quick adapting models show that there is a seasonal variation between 15 and 22 months.

The comparison of stationary gains in Table 2.6 is questionable. It is not certain that the estimated gain in a time-invariant model is comparable with a simple average of a time-varying gain. This may explain the significant differences in the table.

In general, the recursive estimation gives larger standard errors than the two-step ridge regression method. Two explanations for this are:

1. A noise model fitted for the time-invariant case was used for the time-varying case as well. A more appropriate approach would probably be to include the noise model in the recursive estimation.
2. The current parameter estimates of a recursive forgetting factor method will necessarily be a certain number of sampling intervals behind the real parameters. This is because the estimates are based on only past observations. The lag between the estimates and the real values is proportional to $1/(1 - \lambda)$, where λ is the forgetting factor. For a slowly varying system, this lag would be rather unimportant, but for the business cycle data, which seems to represent a rapidly changing system, the parameter estimates cannot be expected to be fully up-to-date, and the accuracy of the forecasts is consequently affected.

2.4 Conclusion

Various methods for tracking the time-delay and the other parameters of a dynamic input-output system have been proposed and tested on simulated data, on data from a district heating system and on business cycle data. An ARMAX model structure was assumed.

At first three methods based upon recursive forgetting factor estimation were presented:

Method A: Parallel forgetting factor estimation of a number of different models. Each model corresponds to one of the possible discrete values of the time-delay of the system. At each sampling instant, the model which represent the current state of the system in the best way is selected.

Method B: An approach proposed by Bányász and Keviczky (1988).

The discrete time-delay is considered as a continuous parameter, and by means of the general Gauss-Newton algorithm a recursive least squares algorithm is derived.

Method C: The embedded continuous time-delay is partitioned into a discrete part and a continuous part. The continuous part is estimated using a recursive least squares algorithm, and the estimate is used as an indicator of changes of the discrete part.

As expected, these methods showed the best results for slowly varying time-delays and dynamics. For the simulated data the real time-variation of the discrete delay was restricted: At each sampling instant, the delay had changed to be one time unit longer or shorter at the most, and the probability of transition to another value was only 0.04. Furthermore, the remaining dynamics were time-invariant. For this slowly varying case, methods A and C were able to follow the variations very well. Method B failed totally, and it was not used in the subsequent case studies.

When applied to the district heating data, only method A gave good results. Method C showed tendencies towards divergence. The results achieved by using method A indicated that the variations of the time-delay were wide and very quick. This is probably the reason for the failure of method C. In method C the continuous time-delay is approximated by the sum of an integer variable and a continuous variable. The approximation is based upon the assumption that the continuous variable is limited to the interval $[-1,1]$. However, if the real time-delay changes too much and too quickly, it is likely that the estimate of the continuous variable falls outside this interval, and as a result the estimate of the discrete variable is corrupted.

In connection with the business cycle data method A was used for both a single-input model and a double-input model. In both cases time-variation of the time-delays as well as the dynamics was re-

vealed.

It is beyond any doubt that method A is more robust than the other forgetting factor methods. Method B seems to work reasonable only for a slowly varying pure time-delay system. Method C is rather sensible to the degree of time-variation. One way to improve the robustness is to use another approximation of the continuous time-delay in order to allow for wider and more abrupt variations of the continuous variable.

Notice that forgetting factor estimation is based on local approximations of the real system by means of a time-invariant model. Furthermore, the current parameter estimates of this model are based on past observations only. The estimator thus has a low-pass filtering effect, i.e. the estimator lags behind the real variations, and a large proportion of the high-frequency variations is filtered out (smoothing). If the forgetting factor is close to one, the smoothing effect is very strong and the influence of noise is suppressed. If the forgetting factor is small, the opposite effect is observed. Therefore, the choice of the forgetting factor should reflect both the degree of parameter variation and the signal-to-noise ratio of the data: the optimal value is decreasing with increasing parameter variation or increasing signal-to-noise ratio of the data. This is illustrated in Figure 2.41.

For the district heating data, a strong cyclic variation of the time-delay was seen, and thus rather low values of the forgetting factors for the parameter estimation turned out to be optimal (typically around 0.95, and in a single case below 0.9).

Two types of explicit models of the embedded parameter variation were considered:

Method D: Deterministic models.

Method E: A stochastic model.

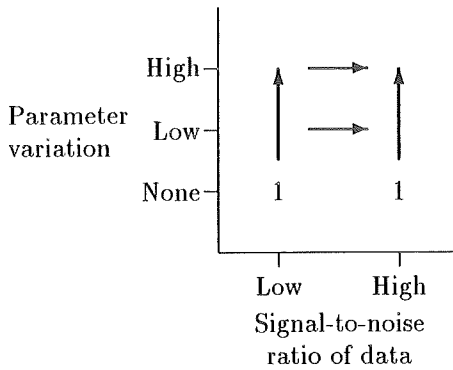


Figure 2.41: *Qualitative illustration of the optimal choice of the forgetting factor. The optimal value is decreasing in the arrow directions.*

The class of deterministic models of embedded parameters considered is linear in the parameters. For a time-varying ARX model this means that the total model is a linear regression model, and that the ordinary recursive least squares method can be used to estimate the parameters. Subsequently the variation of the time-delay can be computed from the variation of the parameters. As an example, if the real variation of the time-delay is known to be cyclic, the coefficients of the numerator polynomial of the transfer function could be modelled by the first harmonics of a Fourier series. It has proved to be very convenient to model diurnal variations of a district heating system in this way since the models are directly applicable for design of an embedded model based predictive controller.

The stochastic description of the embedded parameter variation was formulated as a non-linear state-space model. The state variables of this model were the continuous time-delay and the other time-varying

parameters of an ARX model. By means of an extended Kalman filter the state variables were estimated recursively. Data from a district heating system was used to test this approach, and the results were very similar to those found by using method A. The time-variation of the parameters were described by random walk models and an AR(1) model. As the time-delay showed a clear diurnal variation, the model can probably be improved by the use of a seasonal model of the variations of the time-delay. Furthermore, time-delays in a district heating system are known to depend on the supply temperature. Therefore it might be beneficial to use the supply temperature as an explanatory variable in the embedded model of the time-delay.

Below some remarks and conclusions relating to all or some of the approaches are given.

- Forgetting factor based estimation is easy to use since no specification of the embedded models is required. However, the estimator acts as low-pass filter and a delay. Therefore, in general, a combination of explicit models of the high-frequency components of the variation and forgetting factor estimation of the low-frequency components seems to be a very good solution (e.g. the deterministic description of the parameter variation in method D combined with adaptive estimation).
- For all forgetting factor methods (A, B and C) and the stochastic approach in method E, there is a possibility of having different degrees of adaptability for the time-delay and the other model parameters. Through the experimental results this has proved to be a very important quality. For the district heating data, for instance, the strong seasonal variation of the time-delay required quick adaption, whereas the slower variation of the remaining parameters permitted less adaptability.
- For most of the methods (C, D and E) approximations of the

embedded continuous time-delay have been considered. Actually these approximations were equivalent to approximations of the continuous input signal by interpolation between the sampled values. Linear interpolation was used in methods D and E, while method C was based upon quadratic interpolation. The need for an approximation of the continuous time-delay is due to the fact that the estimates of the parameters (including the time-delay) are found by minimization of an objective function by means of an iterative method, which requires that the derivative of the objective function with respect to the parameters exists. If the continuous input signal is a low-frequency signal (with a negligible effect above the Nyquist frequency) it is sufficient to reconstruct it by an interpolation method. Otherwise the input signal has to be sampled between the ordinary sampling instants to get sufficient information about the variation. In connection with the district heating data, for instance, hourly averages of the original signals were computed in order to obtain low-pass filtered signals and thus make the interpolation technique more applicable. A problem of this low-pass filtering, however, is that changes of the time-delay may be detected up to one hour too late.

- The various algorithms for tracking the time-delay have certain associated algorithm parameters (forgetting factors, minimum and maximum delay, covariance matrices etc.). In order to find reasonable values of these parameters, an objective function measuring the resulting prediction performance was optimized. Since this objective function became discontinuous because of the discretely valued delay, the global optimum was not necessarily found. But several local optima were found, and the most optimal of them was used. Generally it can be concluded that it is important not to give up too early searching for new local optima – seemingly a better one can “always” be found.

When an optimal solution has been found it is important to assess whether this particular solution is reasonable from other points of view. Knowledge of the physical system to be modelled is valuable when assessing the obtained delay and parameter values.

- All methods have been formulated as recursive algorithms. Consequently they are well suited for on-line applications. At the combined heat and power plant, Vestkraft, in Esbjerg, a combination of method A and D has been implemented in association with an adaptive controller for control of the supply temperature. Until now the results of this project have been very promising.

Chapter 3

Dynamic Models for Air Temperature

THE climate system on the earth is determined by an enormous number of interacting subsystems. The energy source driving the climatic system is the sun. The radiation energy from the sun is exposed to absorbing, reflecting, transferring and storing mechanisms on earth and in the atmosphere. Thus the climatic system constitutes a very complicated dynamic system. Fig. 3.1 shows a schematic illustration of the climatic system.

The climate at a particular place on earth is characterized by climatic variables like air temperature, air pressure, wind speed etc. Due to the complexity of the climatic system it is not a simple task to model the variations of these variables. Although it might be possible to establish a physical model for each individual subsystem, it is not certain that a global model of the climate obtained by coupling the sub-models leads to a suitable description of the climatic vari-

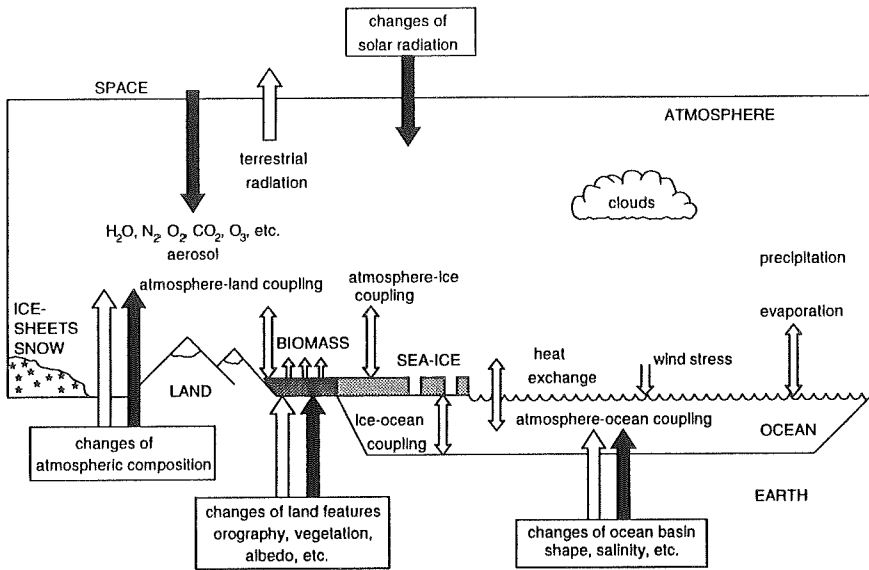


Figure 3.1: Schematic illustration of the climatic system. Reproduction of figure from Houghton (1977).

ables. Stochastic and chaos-like phenomena in the subsystems will inevitably add up in the final model, making it unusable.

The theme of this chapter is stochastic models of the variations of the air temperature. The models are intended for description of local (in a spatial sense) variations of the temperature. The exponential smoothing procedures described in Section 3.1 are empirical based approaches while the continuous time models presented in Section 3.2 are based on physical considerations.

The exponential smoothing procedures are designed for multi-step prediction of air temperature. They are based on empirical descriptions of the variations of the air temperature using components as level, trend and cyclic diurnal variation, and they do not consider the physical mechanisms of the underlying climatic system. The procedures are closely related to ARIMA modelling, and it is shown that the components correspond to embedded ARIMA model structures.

The continuous time models relate the air temperature to other climatic variables, especially the net radiation. These models are primarily intended for a physical description of a limited part of the climatic system. They involve dynamic heat balance equations describing the exchange of heat between various heat reservoirs in the climatic system. The embedded physical parameters of the system are modelled by thermal resistances and heat capacities which are directly identified using the maximum likelihood method.

Both the exponential smoothing approaches and the continuous time models are of particular interest in connection with energy supply (e.g. power and heat supply) where the energy demand is dependent on the ambient air temperature. Forecasts of the ambient air temperature are, for example, necessary for optimal control of the supply temperature in a district heating system (cf. the discussion in Chapter 1).

3.1 Exponential Smoothing

Exponential smoothing belongs to the classical tools in the time-series analysis and plays an important role in literature on time-series analysis and forecasting. Exponential smoothing has been discussed by, among others, Holt (1957), Winters (1960), Brown (1963) and Abraham and Ledolter (1983). Today several computer software packages for forecasting are based upon exponential smoothing.

Exponential smoothing techniques are characterized by being easy to implement and by being computational “inexpensive” to use. For long term forecasting of time series with seasonal components where the underlying physical mechanisms are unknown or not available, such approaches seem to provide better forecasts than the Box-Jenkins transfer function models, including ARIMA models (cf. Makridakis *et al.* (1984)). Therefore exponential smoothing is very relevant to forecasting of air temperature which shows both diurnal and annual variations. Since the air temperature observations studied in this section are sampled at intervals of one hour, it is adequate to concentrate on the diurnal variations. The annual variations are considered as a slowly drifting level which is coped with by the adaptability of the algorithms.

3.1.1 Winters’ Seasonal Forecast Procedure

Winters’ seasonal forecast procedure (Abraham and Ledolter (1983)) forms the basis of the procedures presented in Section 3.1.2. This procedure assumes that the observations belong to a discrete time series, $\{y(t)\}$, generated by the model

$$\begin{aligned}y(t) &= \mu(t) + p(t) + e(t) \\ \mu(t) &= \mu(t-1) + \beta ,\end{aligned}\tag{3.1}$$

where $\mu(t)$ is a linear trend component and $p(i) = p(i + s) = p(i + 2s) = \dots$ (for $i = 1, \dots, s$) are seasonal components which are restricted to add to zero, $\sum_{i=1}^s p(i) = 0$. The components $p(1), \dots, p(s)$ form the seasonal profile of the time series. $\mu(t)$ and β are the level and the slope respectively. s is the number of samples during one season and $\{e(t + j)\}_{j=0, \pm 1, \pm 2, \dots}$ is a sequence of independent and identically distributed random variables with the mean zero and the variance σ_e^2 (later in this chapter a normal distribution is assumed).

The data series generated by the model in (3.1) is the sum of a straight line, a periodic signal and a noise series. It is assumed that the slope and the seasonal profile remain constant in time. In practice, however, it is most likely that these components show some variation in time. Winters' forecast procedure copes with this by using local estimates of the components for the forecasts. The method used for obtaining local estimates is exponential smoothing of the data which, in turn, means that the existence of embedded models of the component variations has been implicitly assumed. In fact it will be shown that Winters' procedure constitutes an optimal predictor if the level, the slope and the seasonal profile are governed by random walk type models.

Ng and Young (1990) introduce a component model which is more general than (3.1). This model includes a trend component, a stochastic perturbation component (ARMA term), a signal component for explanatory variables, a seasonal component and a white noise component. It is shown how this model can be formulated as a state-space model, and a fully recursive algorithm for estimation of the model parameters as well as the states is proposed.

Updating Components and Forecasts

Winters' procedure updates three smoothed statistics $\hat{\mu}(i)$, $\hat{\beta}(i)$ and $\hat{\rho}(i)$ recursively in order to provide estimates of the seasonal and trend components in (3.1). The recursive equations are

$$\hat{\mu}(t) = \alpha_1(y(t) - \hat{\rho}(t-s)) + (1 - \alpha_1)(\hat{\mu}(t-1) + \hat{\beta}(t-1)) \quad (3.2)$$

$$\hat{\beta}(t) = \alpha_2(\hat{\mu}(t) - \hat{\mu}(t-1)) + (1 - \alpha_2)\hat{\beta}(t-1) \quad (3.3)$$

$$\hat{\rho}(t) = \alpha_3(y(t) - \hat{\mu}(t)) + (1 - \alpha_3)\hat{\rho}(t-s). \quad (3.4)$$

The Equations (3.2)-(3.4) all match the pattern $\hat{x}(t) = \alpha_j z(t) + (1 - \alpha_j)\hat{v}(t-q)$. α_j is a smoothing constant, and $z(t)$ represents new information provided with the most recent observation, $y(t)$. Note that α_j corresponds to $1 - \lambda_j$, where λ_j is the forgetting factor, cf. Section 2.2.1. $\hat{v}(t-q)$ which includes the previous estimate of x , is an aggregate measure of the information from past observations, $y(t-1), y(t-2), \dots$. In short, the recursive equations (3.2)-(3.4) simply update the estimates of the smoothed statistics by forming a weighted average of new and past information.

Note that in (3.2) $y(t)$ is adjusted for the seasonal component, $\hat{\rho}(t-s)$, before it is used in the equation, and $y(t)$ in (3.4) is adjusted correspondingly for the level, $\hat{\mu}(t)$. These adjustments reflect the decomposition of the time series into a level (trend) and a seasonal variation.

While $\hat{\mu}(t)$ and $\hat{\beta}(t)$ are updated at each sampling instant, the seasonal coefficients are updated only once during one season - i.e. $\hat{\rho}(t)$ is updated from the previous estimate s steps ago. Note that the updating equations do not ensure that $\sum_{i=0}^{s-1} \hat{\rho}(t+i) = 0$ ($t = 1, 2, 3, \dots$). The procedure could, however, be modified to fulfil this simply by replacing (3.4) by $\hat{\rho}(t) = -\sum_{i=1}^{s-1} \hat{\rho}(t-i)$ if $t = 1, s+1, 2s+1, 3s+1, \dots$

Once the estimates of the level, the slope and the seasonal com-

ponents have been calculated, a forecast of the future observation $y(t+k)$ from time origin t is obtained as

$$\hat{y}(t+k|t) = \hat{\mu}(t) + \hat{\beta}(t)k + \hat{p}(t+k-m), \quad (3.5)$$

where $m = (q+1)s$ for $k = qs+1, \dots, (q+1)s$ ($q = 0, 1, 2, \dots$). The notation $\hat{y}(t+k|t)$ expresses that the k -step-ahead forecast of $y(t+k)$ based on data up to and including time t is considered.

Initial Values for the Updating Equations

The updating equations (3.2)-(3.4) require that initial values of $\hat{\mu}(0)$, $\hat{\beta}(0)$ and $\hat{p}(j-s)$ ($j = 1, 2, \dots, s$) should be specified. Assume that the number of observations is N ($\geq s$). Then a possible choice is

$$\hat{\mu}(0) = \frac{1}{s} \sum_{i=1}^s \pi(i) \quad (3.6)$$

$$\hat{\beta}(0) = 0 \quad (3.7)$$

$$\hat{p}(i-s) = \pi(i) - \hat{\mu}(0), \quad i = 1, \dots, s, \quad (3.8)$$

where $\pi(i)$ is the average of observations from the seasonal period i ($i = 1, \dots, s$):

$$\pi(i) = \frac{\sum_{j=1}^K I(i,j)y(j)}{\sum_{j=1}^K I(i,j)}, \quad i = 1, \dots, s. \quad (3.9)$$

$I(i,j)$ is a seasonal indicator function (1 if j is in seasonal period i , and 0 otherwise), and K ($s \leq K \leq N$) is the number of observations used for the initialization. The optimal value of K depends on N and the smoothing constants (see Abraham and Ledolter (1983)). The choice $\hat{\beta}(0) = 0$ in (3.7) is only reasonable if the trend in the first part of the observations is negligible.

Abraham and Ledolter (1983) suggest an alternative method where the initial values are derived as ordinary least squares estimates of the parameters in the regression model

$$y(t) = \mu + \beta t + \sum_{i=1}^{s-1} \phi(i)(I(i,t) - I(s,t)) + e(t). \quad (3.10)$$

The initial values are obtained from the least squares estimates $\hat{\mu}$, $\hat{\beta}$ and $\hat{\phi}(i)$ ($i = 1, \dots, s$) as $\hat{\mu}(0) = \hat{\mu}$, $\hat{\beta}(0) = \hat{\beta}$, $\hat{p}(j-s) = \hat{\phi}(j)$ ($j = 1, \dots, s-1$) and $\hat{p}(0) = -\sum_{j=1}^{s-1} \hat{\phi}(j)$. At least the first $s+1$ observations are needed to compute these estimates. Abraham and Ledolter (1983) suggest that in case the trend and seasonal pattern are stable, all the available observations should be used.

Choice of the Smoothing Constants

The updating equations (3.2)-(3.4) depend on three smoothing constants, α_1 , α_2 , α_3 , which are often in the interval $[0, 1]$ (Abraham and Ledolter (1983)). If the components (level, slope and seasonal profile) are changing rapidly, then large values of the smoothing constants result in the most accurate forecasts while slowly varying components call for smoothing constants close to zero. If, in particular, $\alpha_i = 0$ ($i = 1, \dots, 3$), the initial values of the components are retained throughout the iterations of (3.2)-(3.4) and (3.5).

The k -step-ahead forecasts, $\hat{y}(t|t-k)$, and the real observations, $y(t)$, can be compared by simulation in order to assess the accuracy of the predictor in (3.5). A measure for this assessment can be obtained by adding the squared one-step-ahead forecast errors, $\varepsilon(t|t-1) = y(t) - \hat{y}(t|t-1)$, over the available set of observations: $\sum_{t=1}^N [\varepsilon(t|t-1)]^2$. Abraham and Ledolter (1983) suggest that the smoothing constants are chosen so that this sum of squared errors (SSE) is minimized.

This is known as the (one-step-ahead) prediction error method (see e.g. Ljung (1987)). For the one-step-ahead predictor this is naturally an appropriate method, but this is not necessarily the case for the k -step-ahead predictor in general. In fact the smoothing constants for each prediction horizon, k , should be tuned separately. Therefore, the following set of multi-step criteria for choosing the α 's are proposed:

$$\min_{\alpha_1(k), \alpha_2(k), \alpha_3(k)} V_k^N(\alpha_1(k), \alpha_2(k), \alpha_3(k)), \quad (3.11)$$

where

$$V_k^N(\alpha_1(k), \alpha_2(k), \alpha_3(k)) = \sum_{t=k}^N [\varepsilon(t|t-k)]^2,$$

and $k = 1, 2, 3, \dots$. Parameter estimation based upon multi-step criteria is well-known. Kabaila (1981) shows that for a rather wide class of models, which includes the ARIMA model structure, the multi-step SSE-criterion leads to an asymptotic normally distributed estimator with the mean equal to the real parameter values and a covariance which decays as N^{-1} (provided that the model structure agrees with the data generating mechanism). Later in this chapter it will be shown that Winters' forecast procedure corresponds to an ARIMA model parametrized by the smoothing constants. Consequently, if the data is generated by such an ARIMA process, then for a suitably large number of observations it can be expected that the multi-step criterion provides consistent parameter estimates. On the other hand, the multi-step estimator cannot be more efficient than the one-step estimator (Kabaila (1981)).

In practice it is unlikely that the data is generated by an ARIMA system. Therefore it can be questioned whether there exists any real values of the estimated parameters. Most frequently we are faced with a parameter "tuning" problem rather than an estimation problem. However, if the models are intended for multi-step forecasting, the multi-step prediction error criteria are well suited for "estimating" the

model parameters. This is the reason for using the criterion in (3.11) when estimating the smoothing constants in Winters' procedure.

By employing the criterion in (3.11), a set of smoothing constants is obtained for each individual prediction horizon. This means that individual sets of smoothed statistics (level, trend and seasonal profile) for the various prediction horizons must be updated at each sampling instant. To reduce the computations of the updating steps and in order to avoid sudden changes in the trace of the predictions, it might be preferable to have only one set of smoothing constants and one set of smoothed statistics, which can be used for all prediction horizons. A criterion for tuning such a set of smoothing constants can be obtained by combining SSE measures for several prediction horizons (see e.g. Stoica and Nehorai (1989)):

$$\min_{\alpha_1, \alpha_2, \alpha_3} F(V_1^N(\alpha_1, \alpha_2, \alpha_3), V_2^N(\alpha_1, \alpha_2, \alpha_3), \dots, V_K^N(\alpha_1, \alpha_2, \alpha_3)), \quad (3.12)$$

where

$$V_k^N(\alpha_1, \alpha_2, \alpha_3) = \sum_{t=k}^N [\varepsilon(t|t-k)]^2,$$

and F is a monotonically increasing function (in each argument). K (≥ 1) is the maximum prediction horizon. If, for instance, F forms a weighted sum (linear combination) of its arguments with positive weights, w_k ($k = 1, \dots, K$), the following criterion arise

$$\min_{\alpha_1, \alpha_2, \alpha_3} \sum_{k=1}^K w_k \sum_{t=k}^N [\varepsilon(t|t-k)]^2.$$

The compound criterion (3.12) results inevitably in compromising values of the smoothing constants. Thus the resulting multi-step predictors have reduced their prediction ability (compared with the predictors obtained with the criteria in (3.11)).

The minimization in (3.11) is unconstrained and may result in α 's that are not within the range of 0 to 1. While it is difficult to give a

reasonable interpretation of a smoothing constant being negative or larger than 1, the numerical calculations of the component estimates and the forecasts do not cause any difficulties.

Since the model in (3.1) is linear in the components, a recursive least squares estimator with a forgetting factor can be applied for estimation of the component (see Section 2.2.1). This estimator is an alternative to Winters' procedure even though only one smoothing constant is included in the estimator, i.e. the components are equally smoothed. This is frequently a disadvantage because the components may not show equal stability. If this is the case the resulting forgetting factor turns out to be a bad compromise that causes one component to be smoothed too much and another to little. Winters' procedure which, on the contrary, attributes different smoothing constants to the various components, offers a lot more flexibility.

However, Parkum (1992) has proposed an extension of the recursive least squares method, the selective forgetting method. This method attributes different forgetting factors to different directions in the parameter space. The directions are parallel to the eigenvectors of the covariance matrix associated with the parameter estimates. The advantage of this method is that nearly constant or poorly excited parameter directions can be associated with forgetting factors close to 1 in order to obtain weak updating in these directions. On the other hand parameter directions with quick time-variation and sufficient excitation can be updated strongly by choosing smaller values of the corresponding forgetting factors.

Additional Comments

Diagnostic Checks. In order to validate Winters' forecast procedure with given values of the smoothing constants the resulting

one-step-ahead prediction errors (residuals) should be examined. For an optimal forecasting procedure the residual sequence will be white noise with the mean zero.

Madsen (1989) lists various ways of testing the white noise hypothesis statistically and visually. One way is to test whether the sample autocorrelation function of the residuals has significant values for non-zero lags. If more than 5% of the autocorrelations lie outside the 95% confidence interval ($\approx \pm 2/\sqrt{N}$) the white noise hypothesis is rejected.

Since exponential smoothing procedures do not automatically ensure that the average of the residuals is zero it should also be tested whether the mean value is zero or not. An approximate test of the hypothesis " $H_0 : \text{mean} = 0$ " at the 5% significance level is carried out by comparing the average of the residuals with $\pm 2s/\sqrt{N}$, where s is the mean square error. If the average lie outside this interval the hypothesis is rejected.

The Equivalent ARIMA Model. It is worth noticing that the predictor in (3.5) together with its updating equations, (3.2)-(3.4), constitutes a linear filter (predictor filter) through which observations of $y(t)$ are passed. The smoothing constants are the parameters of this filter. As the optimal one-step-ahead predictor for an ARIMA model is a linear filter too, there exists a one-to-one correspondence between Winters' one-step-ahead forecast procedure and an ARIMA model of suitable structure and parametrization. In order to find this correspondence it is assumed that the data is generated by an ARIMA model, and that the real model parameters (the smoothing constants) are known.

Since Winters' one-step-ahead predictor takes the form

$$\hat{y}(t+1|t) = \hat{\mu}(t) + \hat{\beta}(t) + \hat{p}(t+1-s),$$

the corresponding ARIMA model is implicitly described by

$$y(t) = \hat{\mu}(t-1) + \hat{\beta}(t-1) + \hat{p}(t-s) + e(t) \quad (3.13)$$

$$\hat{\mu}(t) = \alpha_1(y(t) - \hat{p}(t-s)) + (1 - \alpha_1)(\hat{\mu}(t-1) + \hat{\beta}(t-1)) \quad (3.14)$$

$$\hat{\beta}(t) = \alpha_2(\hat{\mu}(t) - \hat{\mu}(t-1)) + (1 - \alpha_2)\hat{\beta}(t-1) \quad (3.15)$$

$$\hat{p}(t) = \alpha_3(y(t) - \hat{\mu}(t)) + (1 - \alpha_3)\hat{p}(t-s), \quad (3.16)$$

where the fact that $y(t) = \hat{y}(t|t-1) + e(t)$ has been utilized to obtain (3.13) from the predictor. Equations (3.14)-(3.16) are Winters' updating equations. $e(t)$ denotes the noise term. The ARIMA model (or rather an IMA model) in an explicit form is derived by eliminating $\hat{\mu}(t)$, $\hat{\beta}(t)$ and $\hat{p}(t)$. At first, however, it is shown that $\hat{\mu}(t)$, $\hat{\beta}(t)$ and $\hat{p}(t)$ are governed by simple IMA models.

Equation (3.14) is rewritten by using (3.13):

$$\begin{aligned} \hat{\mu}(t) &= \alpha_1(y(t) - \hat{p}(t-s) - \hat{\mu}(t-1) - \hat{\beta}(t-1)) \\ &\quad + \hat{\mu}(t-1) + \hat{\beta}(t-1) \\ &= \alpha_1 e(t) + \hat{\mu}(t-1) + \hat{\beta}(t-1), \end{aligned} \quad (3.17)$$

or

$$\nabla \hat{\mu}(t) = \alpha_1 e(t) + \hat{\beta}(t-1), \quad (3.18)$$

where ∇ denotes the difference operator defined by $\nabla x(t) = x(t) - x(t-1)$. Equation (3.15) is rewritten by using (3.18):

$$\begin{aligned} \hat{\beta}(t) &= \alpha_2(\hat{\mu}(t) - \hat{\mu}(t-1) - \hat{\beta}(t-1)) + \hat{\beta}(t-1) \\ &= \alpha_1 \alpha_2 e(t) + \hat{\beta}(t-1), \end{aligned}$$

or

$$\nabla \hat{\beta}(t) = \alpha_1 \alpha_2 e(t). \quad (3.19)$$

This result is used in (3.18):

$$\nabla \hat{\mu}(t) = \alpha_1 e(t) + \frac{\alpha_1 \alpha_2}{\nabla} e(t-1),$$

or

$$\nabla^2 \hat{\mu}(t) = \alpha_1 e(t) + \alpha_1 (\alpha_2 - 1) e(t-1). \quad (3.20)$$

Finally (3.16) is rewritten by using (3.13) and (3.17):

$$\begin{aligned} \hat{p}(t) &= \alpha_3 (y(t) - \hat{\mu}(t) - \hat{p}(t-s)) + \hat{p}(t-s) \\ &= \alpha_3 (y(t) - \alpha_1 e(t) - \hat{\mu}(t-1) - \hat{\beta}(t-1) - \hat{p}(t-s)) \\ &\quad + \hat{p}(t-s) \\ &= \alpha_3 (1 - \alpha_1) e(t) + \hat{p}(t-s), \end{aligned}$$

or

$$\nabla_s \hat{p}(t) = \alpha_3 (1 - \alpha_1) e(t), \quad (3.21)$$

where ∇_s denotes the seasonal difference operator defined by $\nabla_s x(t) = x(t) - x(t-s)$. Now the IMA models in (3.19), (3.20) and (3.21) are used to eliminate $\hat{\mu}(t-1)$, $\hat{\beta}(t-1)$ and $\hat{p}(t-s)$ in (3.13). After some tedious calculations the following seasonal IMA model is reached:

$$\begin{aligned} \nabla_s \nabla^2 y(t) &= e(t) + (\alpha_1 (1 - \alpha_2) - 2) e(t-1) + (1 - \alpha_1) e(t-2) \\ &\quad + ((1 - \alpha_1) \alpha_3 - 1) e(t-s) \\ &\quad + \alpha_1 (2\alpha_3 - \alpha_2 - 1) e(t-s-1) \\ &\quad + (1 - \alpha_1) (\alpha_3 - 1) e(t-s-2). \end{aligned} \quad (3.22)$$

It can be concluded that Winters' one-step-ahead predictor is the optimal predictor (i.e. $\hat{y}(t|t-1) = E[y(t)|y(t-1), y(t-2), \dots]$) if, and only if, the data generating mechanism is given by this model.

Note that the parametrization of the model in (3.22) is very unusual. Consider the following IMA models which are parametrized in a more

common way

$$\nabla_s \nabla^2 y(t) = (1 + \psi_1 q^{-1})(1 + \psi_2 q^{-1})(1 + \psi_s q^{-s})e(t) \quad (3.23)$$

$$\nabla_s \nabla^2 y(t) = (1 + \psi_1 q^{-1} + \psi_2 q^{-2})(1 + \psi_s q^{-s})e(t) \quad (3.24)$$

$$\begin{aligned} \nabla_s \nabla^2 y(t) &= (1 + \psi_1 q^{-1})(1 + \psi_2 q^{-1} + \psi_s q^{-s} \\ &\quad + \psi_{s+1} q^{-(s+1)})e(t) \end{aligned} \quad (3.25)$$

$$\begin{aligned} \nabla_s \nabla^2 y(t) &= (1 + \psi_1 q^{-1} + \psi_2 q^{-2} + \psi_s q^{-s} \\ &\quad + \psi_{s+1} q^{-(s+1)} + \psi_{s+2} q^{-(s+2)})e(t). \end{aligned} \quad (3.26)$$

Here q^{-1} is the backward shift operator defined by $q^{-1}x(t) = x(t-1)$. A question is: Do any relationships exist between the parameters (ψ 's) of these models and the parameters (α 's) of the model in (3.22)? The answer is that only for some trivial cases of (3.26) that relationship exists. Thus the number of parameters is reduced from five to three if Winters' forecast procedure is used instead of an ordinary IMA model.

Note that Equations (3.19), (3.20) and (3.21) define the embedded models governing the variations of trend, level and seasonal profile.

Equations (3.13) to (3.16) correspond to the innovations form of a state-space model (innovations representation, see e.g. Ljung (1987)). This is easily established by writing (3.13), (3.17), (3.19) and (3.21) as matrix equations:

$$\begin{aligned} \hat{\mathbf{x}}(t) &= \mathbf{A}\hat{\mathbf{x}}(t-1) + \mathbf{K}e(t) \\ \mathbf{y}(t) &= \mathbf{C}\hat{\mathbf{x}}(t-1) + e(t), \end{aligned}$$

where

$$\hat{\mathbf{x}}(t) = (\hat{\mu}(t) \hat{\beta}(t) \hat{p}(t) \cdots \hat{p}(t-s+1))^T$$

$$\begin{aligned}
 A &= \begin{pmatrix} 1 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & 1 & & 0 & 0 \\ \vdots & \vdots & & \ddots & & \vdots \\ 0 & 0 & 0 & & 1 & 0 \end{pmatrix} \\
 K &= (\alpha_1 \quad \alpha_1\alpha_2 \quad \alpha_3(1-\alpha_1) \quad 0 \quad \cdots \quad 0)^T \\
 C &= (1 \ 1 \ 0 \ \cdots \ 0 \ 1).
 \end{aligned}$$

The vector K is known as the stationary Kalman gain and $e(t)$ as the innovation.

The IMA model (3.22) corresponds to Winters' one-step-ahead predictor but it does not in general apply to the k -step-ahead predictor,

$$\hat{y}(t+k|t) = \hat{\mu}(t) + \hat{\beta}(t)k + \hat{p}(t+k-s), \quad 1 \leq k \leq s. \quad (3.27)$$

In order to derive ARIMA models corresponding to the k -step-ahead predictor it should be utilized that the optimal k -step-ahead prediction error in an ARIMA process is a moving average of order $k-1$ (e.g. Ljung (1987)):

$$y(t) - \hat{y}(t|t-k) = e(t) + c_1e(t-1) + \cdots + c_{k-1}e(t-k+1). \quad (3.28)$$

Now insert (3.27) into (3.28):

$$\begin{aligned}
 y(t) &= \hat{\mu}(t-k) + \hat{\beta}(t-k)k + \hat{p}(t-s) \\
 &\quad + e(t) + c_1e(t-1) + \cdots + c_{k-1}e(t-k+1), \quad 1 \leq k \leq s.
 \end{aligned}$$

This is the k -step version of (3.13), and together with the updating equations, (3.14), (3.15) and (3.16), it implicitly describes the ARIMA model of interest. The explicit form of this model is achieved by eliminating $\hat{\mu}(t)$, $\hat{\beta}(t)$ and $\hat{p}(t)$, but this is a rather laborious task when $k > 1$.

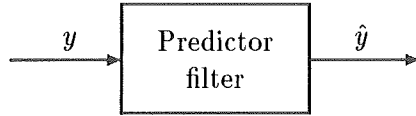


Figure 3.2: *The predictor filter.*

Note that since the moving average coefficients c_1, \dots, c_{k-1} can take any values, there is an infinite number of ARIMA models leading to the same k -step-ahead predictor as Winters' forecast procedure when $k > 1$.

Stability of Predictor Filters. At first consider simple exponential smoothing, where the updating equation of the level takes the form (see e.g. Abraham and Ledolter (1983))

$$\hat{\mu}(t) = (1 - \alpha)\hat{\mu}(t - 1) + \alpha y(t),$$

and the k -step-ahead prediction equation is

$$\hat{y}(t + k|t) = \hat{\mu}(t).$$

These equations can be seen as a state-space formulation of the linear filter giving the forecasts from the sequence of observations, $\{y(t)\}$. Ljung (1987) calls a filter of this kind a predictor filter (Figure 3.2). Since the pole of the filter is $1 - \alpha$, stability is ensured by the restriction $-1 \leq 1 - \alpha \leq 1$, i.e. $0 \leq \alpha \leq 2$. If this stability condition is not satisfied, the observations and the initial estimate, $\hat{\mu}(0)$, of the level influence the forecasts more the older they are, and this is not desirable. In other words it is required that the predictor filter has poles within or on the border of the unit disc in the complex plane.

In order to analyse the stability of the predictor filter corresponding to Winters' forecast procedure, the updating equations and the

prediction equation are transformed into a linear state-space form:

$$\begin{aligned}\hat{\mathbf{x}}(t) &= \tilde{\mathbf{A}}\hat{\mathbf{x}}(t-1) + \tilde{\mathbf{B}}y(t) \\ \hat{y}(t+k|t) &= \tilde{\mathbf{C}}\hat{\mathbf{x}}(t)\end{aligned}$$

where

$$\begin{aligned}\hat{\mathbf{x}}(t) &= (\hat{\mu}(t) \hat{\beta}(t) \hat{p}(t) \cdots \hat{p}(t-s+1))^T \\ \tilde{\mathbf{A}} &= \begin{pmatrix} 1 - \alpha_1 & 1 - \alpha_1 & 0 & \cdots & 0 & -\alpha_1 \\ -\alpha_1\alpha_2 & 1 - \alpha_1\alpha_2 & 0 & \cdots & 0 & -\alpha_1\alpha_2 \\ -\alpha_3(1 - \alpha_1) & -\alpha_3(1 - \alpha_1) & 0 & \cdots & 0 & 1 - \alpha_3(1 - \alpha_1) \\ 0 & 0 & 1 & & 0 & 0 \\ \vdots & \vdots & & \ddots & & \vdots \\ 0 & 0 & 0 & & 1 & 0 \end{pmatrix} \\ \tilde{\mathbf{B}} &= (\alpha_1 \ \alpha_1\alpha_2 \ \alpha_3(1 - \alpha_1) \ 0 \ \cdots \ 0)^T \\ \tilde{\mathbf{C}} &= (1 \ k \ 0 \ \cdots \ 0 \ \underbrace{1 \ 0 \ \cdots \ 0}_{((k-1) \text{ modulus } s)+1 \text{ elements}}).\end{aligned}$$

The stability condition for this filter is that the eigenvalues of $\tilde{\mathbf{A}}$ should lie within or on the border of the unit disc in the complex plane.

Note that $\tilde{\mathbf{A}}$ is independent of the prediction horizon, k . Therefore the relationship between the eigenvalues (= the poles of the predictor filter) and the smoothing constants is also independent of k . Furthermore it is well-known that the poles of the k -step-ahead predictor filter corresponding to an ARIMA model are equal to the zeros of the MA part (see e.g. Ljung (1987) and Holst (1977)). Thus it is concluded that the zeros of the MA parts of the equivalent ARIMA models mentioned on page 118 are independent of k . In other words, the right hand side of (3.22) applies to all the equivalent ARIMA models independent of k .

3.1.2 Alternative Procedures for Exponential Smoothing

In this section four forecast procedures which can be considered as modified versions of Winters' forecast procedure are proposed. The procedures were originally described in Madsen *et al.* (1990). All of them include a level and a seasonal component, but unlike Winters' procedure they do not consider a trend component. The reason for this is that a trend component does not improve the forecasts of air temperature significantly (see Section 3.1.3). Furthermore two of the procedures which have not been treated previously in the literature include a scaling factor for the seasonal component.

Two models are considered:

$$y(t) = \mu + p(t) + e(t) \quad (3.29)$$

$$y(t) = \mu + fp(t) + e(t) . \quad (3.30)$$

where $p(i) = p(i + s)$ (for $i = 1, 2, \dots$). The only difference between (3.29) and (3.30) is the scaling factor, f , which is motivated in the comments below. Note that f and the seasonal profile, $\{p(i)\}_{i=1, \dots, s}$, cannot be identified simultaneously in (3.30) since different combinations of f and $p(t)$ lead to the same model. However, in the smoothing procedures, $p(t)$ is updated independently of f implying that the identifiability problem is eliminated.

Before the four forecast procedures are stated some comments should be given:

- Each of the four forecast procedures consists of some recursive updating equations, some directions for the initialization of these equations and a prediction equation. The recursive equations describe how to update estimates of μ , $p(t)$ and f .

Procedures 1 and 3 are based upon model (3.29) and procedures 2 and 4 upon model (3.30).

- In contrast to Winters' forecast procedure, the observation included in the updating equations for the seasonal profile, $\hat{p}(t)$, is not adjusted for the current level of data. Thus the mean of the profile tends to become equal to the current level rather than to become zero. Therefore, $\hat{\mu}(t)$ should be considered as a level correction accounting for a possible difference between the mean of the seasonal profile and the level of the data. Experiments with the air temperature series have shown that it is beneficial to let the updating equation for the seasonal profile be independent of the $\hat{\mu}(t)$.
- The only difference between procedure 1 and 3 is the order in which the level correction, $\hat{\mu}(t)$, and the profile, $\hat{p}(t)$, are updated. There is a similar difference between procedure 2 and 4. The experimental results in Section 3.1.3 show that for the air temperature data the forecasting results are sensitive to the updating order.
- The seasonal profile is only updated once during a seasonal period. In Winters' forecast procedure and procedure 1 and 2 this means that the deseasonalizing in the updating equations for the level and the scaling factor is carried out using an estimate of the seasonal component which is s samples out of date. This weakness which was also pointed out by Harvey (1989) is a further reason for reordering the updating equations as done in procedure 3 and 4.
- Procedure 2 and 4 represent non-linear filters, and the forecasts are consequently sensitive to the zero point of the observations.
- The four forecast procedures stated above reflect the use they are intended for, namely to forecast air temperatures. It is well-known that the diurnal variations of air temperature are very

dependent on the net radiation (see Section 3.2). The difference between the maximum and minimum air temperature is much larger on bright days than on cloudy days. A change from a bright sky to cloudy weather and vice versa can take place in the course of few hours. The seasonal profile, which is included in the forecast procedures, consists of 24 values – one for each of the 24 hours of the day – but each value is only updated once a day. Consequently, the profile may adapt too slowly to the new situation if the weather changes. This is why models and forecast procedures with profile scaling factors are suggested. Since the scaling factor is updated at each sampling instant, the resulting profile, $\hat{f}(t)\hat{p}(t)$, will adapt more quickly to changes.

- In the initialization step all the N observations are used. As an alternative the initial values could have been based upon the first K ($< N$) observations, where K is chosen in consideration of the smoothing constants and N . In this way the forecasts obtained during the first samples would possibly be better.

Procedure 1

Updating equations:

$$\begin{aligned}\hat{\mu}(t) &= \alpha_1(y(t) - \hat{p}(t - s)) + (1 - \alpha_1)\hat{\mu}(t - 1) \\ \hat{p}(t) &= \alpha_2 y(t) + (1 - \alpha_2)\hat{p}(t - s)\end{aligned}$$

Initialization:

$$\begin{aligned}\hat{\mu}(0) &= 0 \\ \hat{p}(i) &= \frac{\sum_{j=1}^N I(i, j)y(j)}{\sum_{j=1}^N I(i, j)}, \quad i = 1, \dots, s\end{aligned}$$

Prediction equation:

$$\hat{y}(t + k|t) = \hat{\mu}(t) + \hat{p}(t + k - m),$$

where $m = (q + 1)s$ for $k = qs + 1, \dots, (q + 1)s$ ($q = 0, 1, 2, \dots$)

□

Procedure 2

Updating equations:

$$\hat{f}(t) = \begin{cases} \alpha_3 \frac{y(t) - \hat{\mu}(t-1)}{\hat{p}(t-s)} + (1 - \alpha_3)\hat{f}(t-1) & \text{if } \hat{p}(t-s) \neq 0 \\ \hat{f}(t-1) & \text{if } \hat{p}(t-s) = 0 \end{cases}$$

$$\hat{\mu}(t) = \alpha_1(y(t) - \hat{f}(t)\hat{p}(t-s)) + (1 - \alpha_1)\hat{\mu}(t-1)$$

$$\hat{p}(t) = \alpha_2 y(t) + (1 - \alpha_2)\hat{p}(t-s)$$

Initialization:

$$\hat{\mu}(0) = 0$$

$$\hat{p}(i) = \frac{\sum_{j=1}^N I(i, j)y(j)}{\sum_{j=1}^N I(i, j)}, \quad i = 1, \dots, s$$

$$\hat{f}(0) = 1$$

Prediction equation:

$$\hat{y}(t+k|t) = \hat{\mu}(t) + \hat{f}(t)\hat{p}(t+k-m),$$

where $m = (q + 1)s$ for $k = qs + 1, \dots, (q + 1)s$ ($q = 0, 1, 2, \dots$)

□

Procedure 3

Updating equations:

$$\hat{p}(t) = \alpha_2 y(t) + (1 - \alpha_2)\hat{p}(t-s)$$

$$\hat{\mu}(t) = \alpha_1(y(t) - \hat{p}(t)) + (1 - \alpha_1)\hat{\mu}(t-1)$$

Initialization:

$$\hat{\mu}(0) = 0$$

$$\hat{p}(i) = \frac{\sum_{j=1}^N I(i, j)y(j)}{\sum_{j=1}^N I(i, j)}, \quad i = 1, \dots, s$$

Prediction equation:

$$\hat{y}(t + k|t) = \hat{\mu}(t) + \hat{p}(t + k - m),$$

where $m = (q + 1)s$ for $k = qs + 1, \dots, (q + 1)s$ ($q = 0, 1, 2, \dots$)

□

Procedure 4

Updating equations:

$$\begin{aligned} \hat{p}(t) &= \alpha_2 y(t) + (1 - \alpha_2) \hat{p}(t - s) \\ \hat{f}(t) &= \begin{cases} \alpha_3 \frac{y(t) - \hat{\mu}(t - 1)}{\hat{p}(t)} + (1 - \alpha_3) \hat{f}(t - 1) & \text{if } \hat{p}(t) \neq 0 \\ \hat{f}(t - 1) & \text{if } \hat{p}(t) = 0 \end{cases} \\ \hat{\mu}(t) &= \alpha_1 (y(t) - \hat{f}(t) \hat{p}(t)) + (1 - \alpha_1) \hat{\mu}(t - 1) \end{aligned}$$

Initialization:

$$\begin{aligned} \hat{\mu}(0) &= 0 \\ \hat{p}(i) &= \frac{\sum_{j=1}^N I(i, j)y(j)}{\sum_{j=1}^N I(i, j)}, \quad i = 1, \dots, s \\ \hat{f}(0) &= 1 \end{aligned}$$

Prediction equation:

$$\hat{y}(t + k|t) = \hat{\mu}(t) + \hat{f}(t) \hat{p}(t + k - m),$$

where $m = (q + 1)s$ for $k = qs + 1, \dots, (q + 1)s$ ($q = 0, 1, 2, \dots$)

□

Equivalent Embedded ARIMA Models

Due to the linearity of procedures 1 and 3, ARIMA models corresponding to the one-step-ahead predictors can be derived. The models are derived in the same way as the ARIMA model in Section 3.1.1 (see page 115 ff.).

The equivalent embedded models describing the variations of the level and the seasonal profile for procedure 1 are

$$\begin{aligned}\nabla\hat{\mu}(t) &= \alpha_1 e(t) \\ \nabla_s \nabla\hat{p}(t) &= \alpha_2(1 + (\alpha_1 - 1)q^{-1})e(t),\end{aligned}$$

and the resulting ARIMA (IMA) model becomes

$$\nabla_s \nabla y(t) = (1 + (\alpha_1 - 1)q^{-1})(1 + (\alpha_2 - 1)q^{-s})e(t). \quad (3.31)$$

It turns out that this model is equivalent to the k -step-ahead predictor for all $k \geq 1$ (Bjerre (1992)).

Concerning the one-step-ahead predictor of procedure 3 it is found that the equivalent models of the level and the seasonal profile are

$$\begin{aligned}(1 + (\alpha_1\alpha_2 - 1)q^{-1})\hat{\mu}(t) &= \alpha_1(1 - \alpha_2)e(t) \\ (1 + (\alpha_1\alpha_2 - 1)q^{-1})\nabla_s\hat{p}(t) &= \alpha_2(1 + (\alpha_1 - 1)q^{-1})e(t).\end{aligned}$$

In this case a true ARIMA model is obtained

$$(1 - a_1q^{-1})\nabla_s y(t) = (1 + (\alpha_1 - 1)q^{-1})(1 + (\alpha_2 - 1)q^{-s})e(t), \quad (3.32)$$

where

$$a_1 = 1 - \alpha_1\alpha_2. \quad (3.33)$$

Bjerre (1992) showed that (3.32) also holds for $k > 1$ if a_1 is found as an solution to

$$a_1^k - (1 - \alpha_1)a_1^{k-1} - \alpha_1(1 - \alpha_2) = 0$$

instead of using (3.33).

Choice of the Smoothing Constants

As regards the choice of smoothing constants, the comments given in Section 3.1.1 also apply here. The criterion used for choice of smoothing constants corresponds to the one in (3.11),

$$\min_{\alpha(k)} \text{SSE}(\alpha(k)) = \sum_{t=k}^N [\varepsilon(t|t-k)]^2, \quad k = 1, 2, \dots, \quad (3.34)$$

where $\alpha(k) = (\alpha_1(k) \alpha_2(k))^T$ for procedures 1 and 3, and $\alpha = (\alpha_1(k) \alpha_2(k) \alpha_3(k))^T$ for procedures 2 and 4. $\varepsilon(t|t-k)$ is the k -step-ahead forecast error ($= y(t) - \hat{y}(t|t-k)$).

3.1.3 Prediction of Air Temperature – Results and Discussion

The Data

Winters' forecast procedure and procedures 1 to 4 have been tested with air temperature observations measured in Esbjerg. The observations are hourly averages, and each average has been computed from 12 instantaneous measurements during one hour. The hourly observations have been collected during the period from 6 July 1989 4:00p.m. to 22 August 1989 2:00p.m. Hence, the number of observations is $N = 1127$ with $s = 24$ observations within each seasonal period. Figure 3.3 shows the observations. A 24 hour seasonal variation appears from the data. Since the data are from a summer period the temperature of the air is at its minimum between 4a.m. and 7a.m. and at its maximum between 2p.m. and 5p.m. In certain periods the mean level of the series seems to follow a positive or negative trend (e.g. for $60 \leq t \leq 240$), and sometimes a jump in the mean level can be observed (e.g. for $t \approx 585$). Also the amplitude of the variations

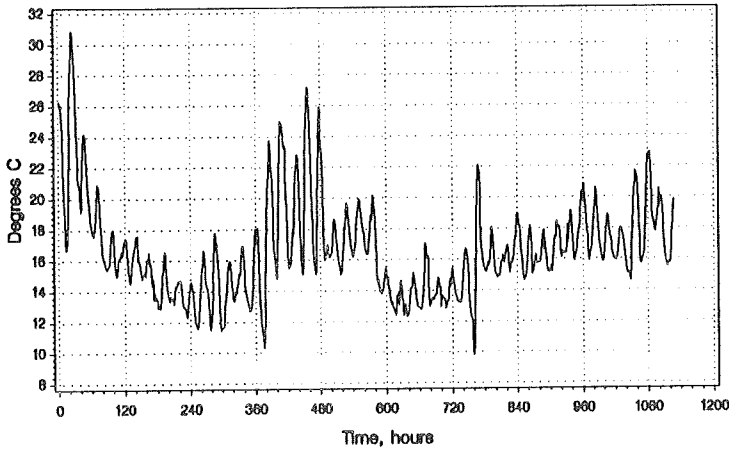


Figure 3.3: *Hourly averaged temperatures of the air during the period from 6 July 1989 4:00p.m. to 22 August 22 1989 2:00p.m.*

is exposed to changes, both slow (e.g. for $220 \leq t \leq 300$) and sudden ones (e.g. for $t \approx 485$). The changes of level and amplitude are strongly related to changes in the weather. The largest amplitudes are seen on bright days while cloudy weather implies modest amplitudes. A sudden change in the level very often occurs while a cold or a warm front is passing.

Before the prediction results for the various forecast procedures are presented, prediction results for four very simple models (random walk type models) are shown (Table 3.2). Furthermore an IMA model has been identified and estimated:

$$\begin{aligned}
 \nabla_{24}\nabla y(t) &= (1 + 0.51q^{-1} + 0.19q^{-2} + 0.14q^{-3})(1 - 0.12q^{-6}) \\
 &\quad (1 - 0.13q^{-8} - 0.07q^{-9} - 0.10q^{-10} - 0.08q^{-11}) \\
 &\quad (1 + 0.08q^{-21})(1 + 0.08q^{-23} - 0.83q^{-24} + 0.09q^{-25}) \\
 &\quad (1 - 0.14q^{-48})e(t).
 \end{aligned}
 \tag{3.35}$$

Table 3.1: Verbal description of the predictors resulting from some simple models.

Model	Forecasts
$y(t) = K + e(t)$	The forecasts are equal to the average, \hat{K} , of the entire series and independent of time and the prediction horizon: $\hat{y}(t+k t) = \hat{K}$.
$\nabla y(t) = e(t)$	The forecasts are equal to the latest observation and independent of the prediction horizon: $\hat{y}(t+k t) = y(t)$.
$\nabla_{24}y(t) = e(t)$	In order to forecast the temperature at a future hour, the observation at the same hour of the preceding day is used: $\hat{y}(t+k t) = y(t+k-24)$.
$\nabla_{24}\nabla y(t) = e(t)$	It is assumed that the gradients at future hours are the same as gradients at the same hours of the preceding day: $\hat{y}(t+k t) = y(t) + y(t+k-24) - y(t-24)$.

Although this model seems to be rather sophisticated, all 13 parameters are statistically significant (t- and F-test). The model structure has been found by estimating various models and checking their residual autocorrelation structure in order to find a model reducing the observations to white noise. The parameters have been estimated by minimizing the sum of squared one-step-ahead prediction errors (unconditional least squares criterion) using the ARIMA procedure of the SAS System.

The predictors corresponding to the simple models are very easily described in verbal terms. This is done in Table 3.1.

From Table 3.2 it is seen that the IMA model provide better forecasts than any of the simple models, but this is not surprising since 13

Table 3.2: Standard deviations of the k -step-ahead prediction errors for four simple models and an IMA model.

k	MODELS				
	Random walk models				IMA model
	$y(t) =$ $K + e(t)$ $(\hat{K} = 16.65)$	$\nabla y(t) =$ $e(t)$	$\nabla_{24}y(t) =$ $e(t)$	$\nabla_{24}\nabla y(t) =$ $e(t)$	$\nabla_{24}\nabla y(t) =$ $\psi(q^{-1})e(t)$
	$\hat{\sigma}_e$	$\hat{\sigma}_e$	$\hat{\sigma}_e$	$\hat{\sigma}_e$	$\hat{\sigma}_e$
1	3.167	0.621	2.138	0.614	0.407
2	3.167	1.145	2.138	1.040	0.728
3	3.167	1.620	2.138	1.387	0.991
4	3.167	2.055	2.138	1.679	1.221
5	3.167	2.446	2.138	1.920	1.408
6	3.167	2.791	2.138	2.123	1.569
7	3.167	3.085	2.138	2.287	1.695
8	3.167	3.328	2.138	2.421	1.796
9	3.167	3.516	2.138	2.520	1.872
10	3.167	3.650	2.138	2.591	1.928
11	3.167	3.729	2.138	2.633	1.967
12	3.167	3.753	2.138	2.655	1.994
13	3.167	3.726	2.138	2.666	2.013
14	3.167	3.650	2.138	2.670	2.024
15	3.167	3.533	2.138	2.679	2.032
16	3.167	3.380	2.138	2.697	2.040
17	3.167	3.197	2.138	2.725	2.044
18	3.167	2.995	2.138	2.766	2.052
19	3.167	2.782	2.138	2.817	2.062
20	3.167	2.572	2.138	2.875	2.075
21	3.167	2.381	2.138	2.940	2.093
22	3.167	2.233	2.138	3.012	2.117
23	3.167	2.145	2.138	3.082	2.146
24	3.167	2.138	2.138	3.136	2.184

parameters has been estimated in the IMA model. Concerning the traditional random walk model ($\nabla y(t) = e(t)$) it is worth noting that relatively large standard errors occur when predicting about 12 hours ahead. The reason is that the 12-step-ahead predictor expects same temperature 12 hours ahead as now, i.e. day temperatures at night and night temperatures in the daytime. The same phenomenon is also seen for the neighbouring predictors, although less dominating.

One way to evaluate the performance of the exponential smoothing procedures is to compare the standard deviations of their forecast errors with the corresponding standard deviations for the simple models and the IMA model. It can, of course, be expected that exponential smoothing provides more accurate forecasts than the simple models since exponential smoothing represents more complex models with their parameters (smoothing constants) being optimized for each prediction horizon (compare the simple models with (3.31) and (3.32)).

Optimal Smoothing Constants and Prediction Results

For each forecast procedure, the optimal values of the smoothing constants are found for prediction horizons of 1, 2, ..., 24 hours by use of the criterion in (3.34). The investigations include Winters' procedure which consists of the equations in (3.2), (3.3), (3.4), (3.5) and (3.8).

Tables 3.3 to 3.5 show the minimization results: the optimal smoothing constants and the standard errors (i.e. the standard deviations of the k -step-ahead prediction errors) for the various prediction horizons. Furthermore, the standard errors are shown in Figure 3.4 to permit a visual comparison of the results.

Normally smoothing constants are assumed to be in the interval $[0, 1]$.

Table 3.3: *Results from using Winters' forecast procedure: optimal smoothing constants and standard errors for $k = 1, \dots, 24$.*

Prediction horizon	Smoothing constants			Standard errors
	α_1	α_2	α_3	
k				$\hat{\sigma}$
1	1.464	-0.0028	0.0118	0.409
2	1.545	-0.0029	0.0047	0.748
3	1.621	-0.0029	0.0060	1.034
4	1.612	-0.0029	0.0058	1.294
5	1.648	-0.0029	0.0052	1.516
6	1.623	-0.0030	0.0043	1.720
7	0.086	-0.0025	-0.0268	1.830
8	0.082	-0.0025	-0.0268	1.868
9	0.081	-0.0025	-0.0267	1.896
10	0.082	-0.0025	-0.0268	1.914
11	0.085	-0.0025	-0.0270	1.926
12	0.091	-0.0025	-0.0274	1.931
13	0.098	-0.0025	-0.0279	1.930
14	0.108	-0.0024	-0.0287	1.926
15	0.119	-0.0024	-0.0298	1.921
16	0.134	-0.0024	-0.0312	1.914
17	0.151	-0.0024	-0.0328	1.910
18	0.171	-0.0024	-0.0350	1.908
19	0.196	-0.0024	-0.0379	1.912
20	0.226	-0.0024	-0.0393	1.925
21	0.265	-0.0024	-0.0392	1.947
22	0.329	-0.0024	-0.0390	1.979
23	0.485	-0.0024	-0.0370	2.018
24	1.136	-0.0025	0.1645	2.054

Table 3.4: Results from using the forecast procedures 1 and 2: optimal smoothing constants and standard errors for $k = 1, \dots, 24$.

Pred. Hor.	Procedure 1			Procedure 2			
	Smoothing constants		Std. err.	Smoothing constants			Std. err.
k	α_1	α_2	$\hat{\sigma}$	α_1	α_2	α_3	$\hat{\sigma}$
1	1.465	-0.0258	0.403	0.557	-0.0252	2.008	0.396
2	1.546	-0.0264	0.738	0.501	-0.0266	2.050	0.720
3	1.617	-0.0273	1.022	0.471	-0.0282	2.130	0.992
4	1.610	-0.0284	1.280	0.471	-0.0307	2.110	1.239
5	1.644	-0.0286	1.502	0.459	-0.0317	2.155	1.450
6	1.621	-0.0282	1.707	0.454	-0.0314	2.098	1.645
7	0.088	-0.0295	1.838	0.450	-0.0310	2.077	1.803
8	0.084	-0.0299	1.878	-0.024	-0.0319	0.119	1.845
9	0.082	-0.0303	1.906	-0.024	-0.0325	0.116	1.875
10	0.083	-0.0308	1.925	-0.025	-0.0331	0.118	1.895
11	0.086	-0.0313	1.937	-0.028	-0.0336	0.122	1.906
12	0.091	-0.0318	1.942	-0.030	-0.0342	0.129	1.912
13	0.098	-0.0324	1.942	-0.033	-0.0348	0.139	1.913
14	0.107	-0.0330	1.939	-0.037	-0.0354	0.150	1.910
15	0.118	-0.0337	1.933	-0.040	-0.0360	0.162	1.907
16	0.131	-0.0345	1.927	-0.043	-0.0367	0.178	1.904
17	0.147	-0.0356	1.923	-0.050	-0.0377	0.197	1.903
18	0.165	-0.0372	1.922	-0.063	-0.0392	0.223	1.905
19	0.186	-0.0393	1.926	-0.090	-0.0411	0.262	1.911
20	0.213	-0.0415	1.937	-0.146	-0.0428	0.322	1.924
21	0.249	-0.0439	1.958	-0.372	-0.0438	0.464	1.943
22	0.309	-0.0466	1.988	(*)	(*)	(*)	(*)
23	0.457	-0.0506	2.026	(*)	(*)	(*)	(*)
24	1.058	-0.0603	2.054	(*)	(*)	(*)	(*)

(*): The minimization algorithm did not converge.

Table 3.5: Results from using the forecast procedures 3 and 4: optimal smoothing constants and standard errors for $k = 1, \dots, 24$.

Pred. Hor.	Procedure 3			Procedure 4			
	Smoothing constants		Std. err.	Smoothing constants			Std. err.
k	α_1	α_2	$\hat{\sigma}$	α_1	α_2	α_3	$\hat{\sigma}$
1	1.462	0.0107	0.410	0.622	0.0068	2.166	0.403
2	1.547	0.0310	0.749	0.631	0.0255	2.400	0.728
3	1.626	0.0633	1.033	0.629	0.0513	2.608	0.998
4	1.621	0.1218	1.286	0.621	0.0825	2.560	1.242
5	1.658	0.2232	1.487	0.595	0.1213	2.554	1.445
6	1.638	0.3272	1.648	0.579	0.1802	2.438	1.623
7	1.637	0.4092	1.765	0.559	0.2712	2.381	1.758
8	0.159	0.3886	1.893	0.554	0.3764	2.370	1.860
9	0.149	0.3970	1.921	0.076	0.3148	0.073	1.916
10	0.148	0.4055	1.940	0.080	0.3339	0.071	1.936
11	0.152	0.4105	1.953	0.077	0.3346	0.079	1.948
12	0.158	0.4105	1.962	0.064	0.3161	0.097	1.955
13	0.167	0.4040	1.967	0.021	0.2408	0.136	1.958
14	0.177	0.3893	1.970	0.005	-0.0273	0.111	1.942
15	0.117	-0.0282	1.958	0.010	-0.0271	0.117	1.940
16	0.130	-0.0284	1.954	0.019	-0.0268	0.123	1.939
17	0.145	-0.0286	1.951	0.031	-0.0264	0.126	1.940
18	0.163	-0.0290	1.952	0.050	-0.0263	0.126	1.944
19	0.183	-0.0299	1.958	0.079	-0.0272	0.119	1.954
20	0.208	-0.0310	1.973	0.126	-0.0291	0.097	1.971
21	0.240	-0.0323	1.997	0.189	-0.0318	0.064	1.997
22	0.288	-0.0333	2.032	0.254	-0.0335	0.044	2.032
23	0.392	-0.0328	2.077	0.351	-0.0332	0.059	2.077
24	1.546	0.6881	2.118	(*)	(*)	(*)	(*)

(*): The minimization algorithm did not converge.

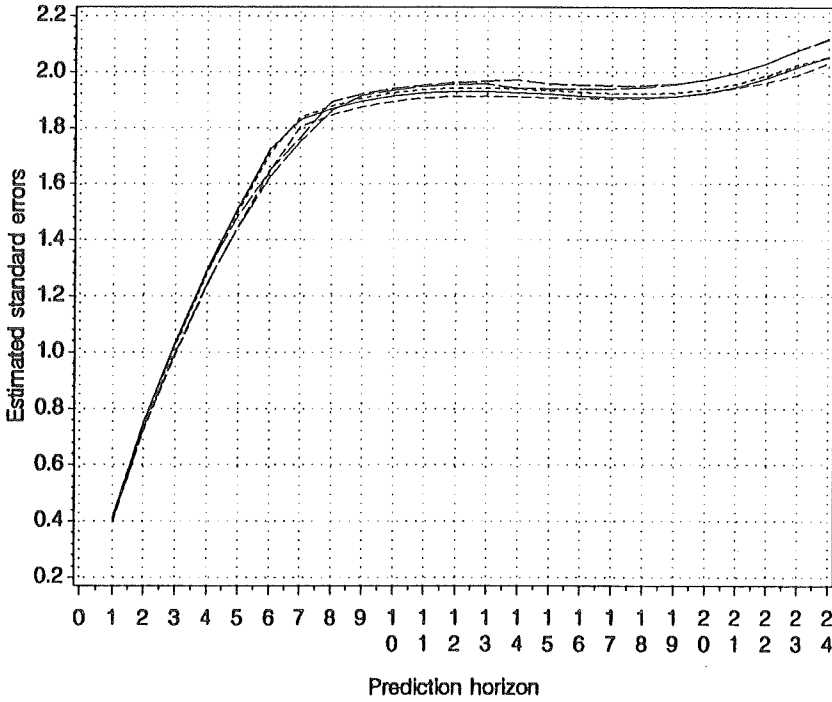


Figure 3.4: Standard errors versus the prediction horizon.

Winters' procedure: ——— . Procedure 1:
 Procedure 2: - - - - . Procedure 3: - · - · .
 Procedure 4: — — — .

The optimal smoothing constants in Tables 3.3 to 3.5 conflict with this assumption, but this does not necessarily mean that the predictor filters are unstable. For the linear procedures (Winters' procedure, procedure 1 and procedure 3) an analysis of the stability can be performed by means of the technique briefly described in Section 3.1.1 – that is by analyzing the eigenvalues of the transition matrix in a state-space formulation of the predictor filter. By doing so it is found that some of the estimated sets of smoothing constants result in eigenvalues slightly outside the unit disc in the complex plane. This does not cause any problems as long as a time series consisting of only 1127 observations is being forecasted. An exponential divergence of the components will occur but is so weak that the 1127th forecast is not seriously affected. On the other hand the problems will be very clear if the predictor filter is used for the observations in another very long series. One way to solve the problems is to employ longer series in the optimization phase. However, the source of the problems might be the way of initializing the forecast procedures. All the observations are used for the initialization but it would perhaps be more natural to use only some of the first observations for this purpose – especially because large variations of the trend and the seasonal pattern are seen in the observations. Though, experiments show that using fewer observations for the initialization gives rise to increasing standard deviation of prediction errors. An alternative way of obtaining the initial values would be to estimate them together with the smoothing constants by minimization of a prediction error criterion like the one in (3.34). A drawback of this approach is that the number of values to be estimated increases significantly – from two or three to 27 or 29 depending on whether it is a procedure with two or three smoothing constants.

In general, the optimal smoothing constants seem to be the largest for short prediction horizons and closer to zero for long horizons. Hence the smoothing effect becomes stronger as the prediction horizon is

prolonged. This tendency is not surprising since the most recent observation gives much information about the near future whereas information from earlier observations have to be given more weight when the distant future is to be predicted.

The optimal α_1 's in Winters' procedure (see Table 3.3) are about 1.5 to 1.6 for $k \leq 6$. From $k = 6$ to $k = 7$, α_1 drops considerably (from 1.623 to 0.086). The sign of α_3 is altered simultaneously. These changes reflect the fact that the function defined in (3.34) has more than one local minimum. Actually, for $k = 6$ and $(\alpha_1, \alpha_2, \alpha_3) \approx (0.09, -0.003, -0.03)$ a local minimum is found which is not much larger than the one shown in Table 3.3. The existence of local minima indicates that the predictor has to choose between two conflicting interests: a very quick and a slowly updating mechanism for the level component. For $k \leq 6$, the best results are experienced when the former interest is satisfied while the latter should be fulfilled for $k \geq 7$. Notice that the sudden drop of α_1 gives rise to a break of the standard error curve in Figure 3.4 about $k = 7$.

As it appears from Tables 3.4 and 3.5 similar drops in the smoothing constants are seen for procedures 1, 2, 3 and 4. The same comments as given to Winters' forecast procedure apply to these procedures.

The IMA model in (3.35) has a certain resemblance to the IMA model in (3.22) which is equivalent to Winters' one-step-ahead forecast procedure. Therefore, there is reason to believe that the parameters of the former model would also show sudden changes if they were optimized for different prediction horizons.

As to the standard errors listed in the tables and shown in Figure 3.4 it is seen that procedure 2 provides more accurate predictions than the other procedures for most prediction horizons. But for $k = 5, 6, 7$ procedure 4 provides slightly better results. Consequently it can be concluded that, among the procedures with three smoothing con-

starts one which includes a scaling factor should be preferred to one which includes a trend component (like in Winters' procedure). For most prediction horizons the relative difference between the largest and the smallest standard error is about 4%. This means that none of the procedures are really superior or inferior to the others. When the results are compared with those in Table 3.2, it is found that the forecast accuracy of the exponential smoothing procedures is between 5% and 90% better than for the four simple models (depending on forecast horizon). As to the IMA model it seems to perform better than Winters' procedure and procedures 1 and 3 if the forecast horizon is less than about 10 hours. For horizons longer than 10 hours, the IMA model becomes inferior. When the IMA model is compared with procedures 2 and 4, it is found that the IMA model is only superior for horizons between 3 and 9 (10) hours. That the IMA model is the best for certain horizons may in part be due to the fact that the IMA model has 13 parameters while the exponential smoothing procedures only have 2 or 3 parameters (the smoothing constants). However, comparing the IMA model with the smoothing procedures for horizons longer than one hour is not quite fair since the IMA model uses the same set of parameters for all horizons while the smoothing procedures use individual sets of parameters for different horizons.

Figure 3.5 shows the one-step-ahead forecasts made by Winters' procedure together with the real observations. The two curves are almost indistinguishable. Figures 3.6 and 3.7 show the corresponding plots for six and seven-step-ahead forecasts. These are chosen as typical examples. In these cases it is easy to distinguish observations and forecasts. Note that Figures 3.6 and 3.7 represent quickly and slowly updating, respectively (cf. the above discussion). Figures 3.8 and 3.9 show magnified versions of the same two plots. In both cases reasonable forecasts are obtained for $240 \leq t \leq 368$. For $368 < t \leq 480$, however, the six and seven-step-ahead forecasts are extremely different. The six-step-ahead predictor tries to track variations by varying

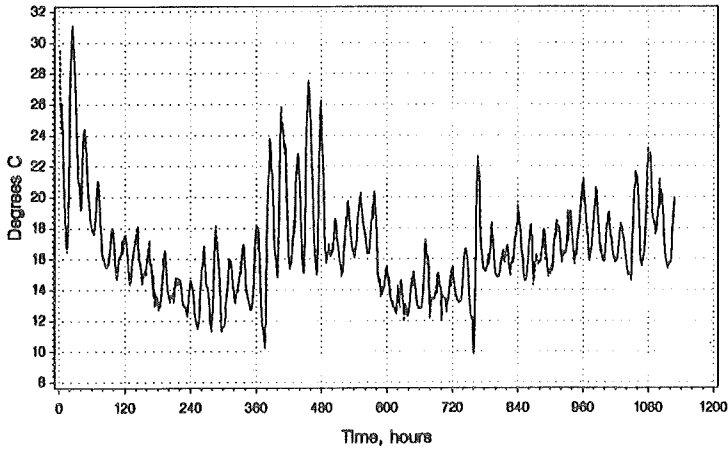


Figure 3.5: *One-step-ahead forecasting with Winters' forecast procedure (dashed curve) and real observations (solid curve)*

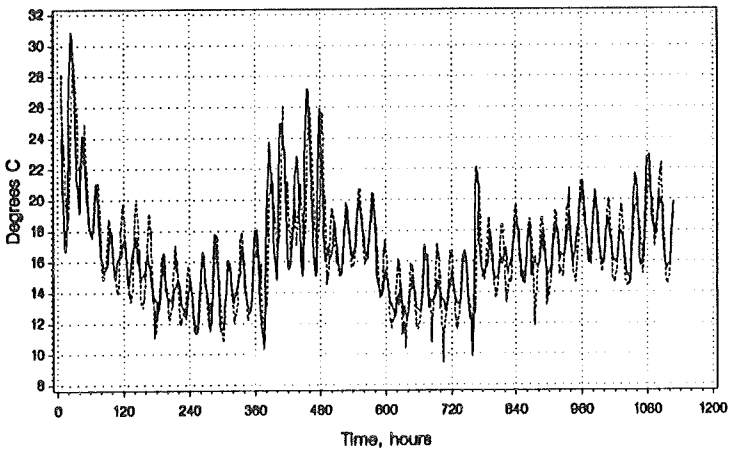


Figure 3.6: *Six-step-ahead forecasting with Winters' forecast procedure (dashed curve) and real observations (solid curve)*

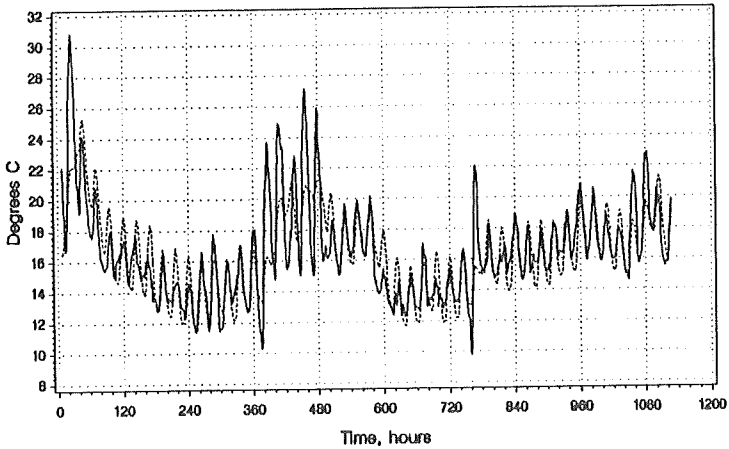


Figure 3.7: *Seven-step-ahead forecasting with Winters' forecast procedure (dashed curve) and real observations (solid curve)*

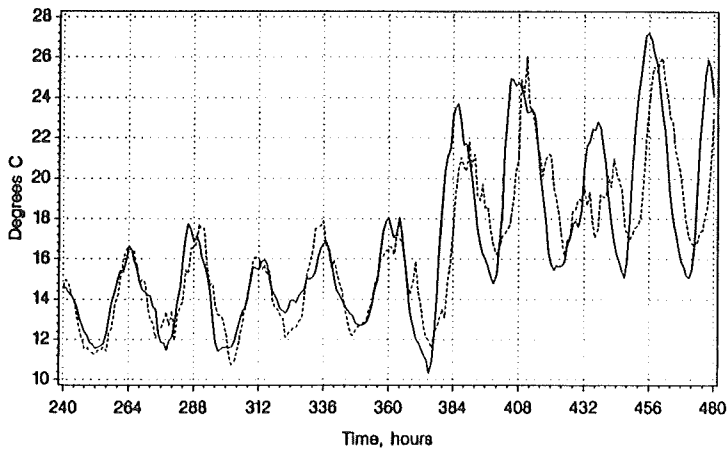


Figure 3.8: *Six-step-ahead forecasting with Winters' forecast procedure (dashed curve) and real observations (solid curve)*

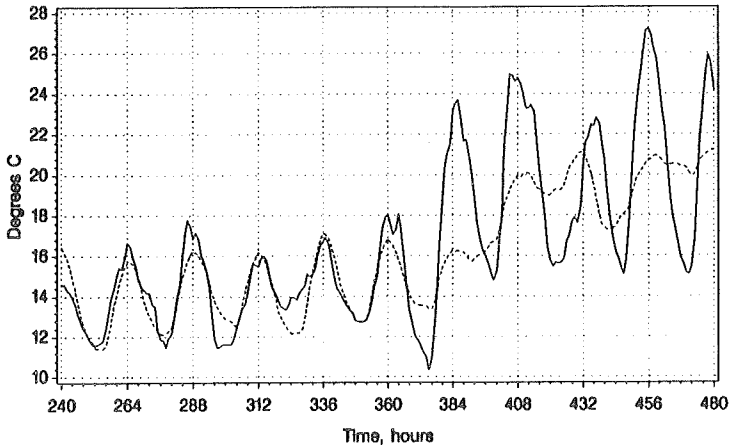


Figure 3.9: *Seven-step-ahead forecasting with Winters' forecast procedure (dashed curve) and real observations (solid curve)*

the level (not the seasonal profile). Unfortunately this leads to a 6 hour time lag between the observations and the forecasts. For the seven-step-ahead predictor it has been found to be more profitable to intensify the smoothing of the level. In this way the time lag phenomenon is avoided, but the amplitude of the forecasts is, on the other hand, much smaller than the amplitude of the real observations.

To illustrate in more details how Winters' six and seven-step-ahead predictors behave, Figures 3.10 and 3.11 show the estimated level, trend and seasonal components versus time. Most of the time, the trend component, $\hat{\beta}(t)$, is very close to zero. For that reason it has been multiplied by 500 in the figures. Actually, the six and seven-step-ahead forecasts are not influenced very much by $\hat{\beta}(t)$ (less than $0.2\text{ }^{\circ}\text{C}$). Therefore a simpler model without the trend component would perform perfectly well. Notice that the estimated seasonal profiles for the six and seven-step-ahead predictors are almost identical and stable in time. In contrast to that, the estimated levels are

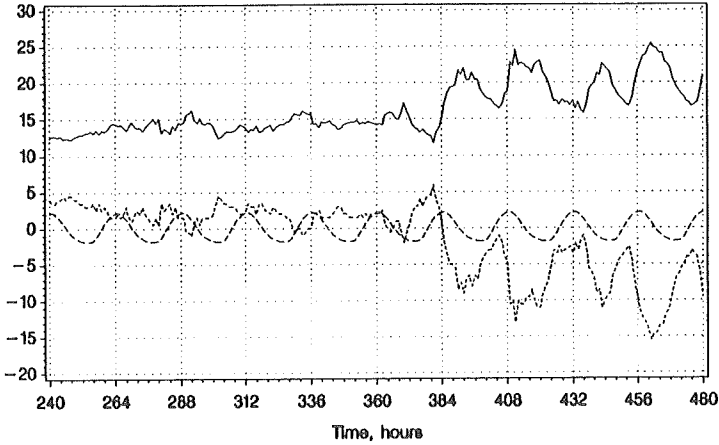


Figure 3.10: *Estimated components in Winters' six-step-ahead predictor:* $\hat{\mu}(t - 6)$: — . $500 \times \hat{\beta}(t - 6)$: - - - - . $\hat{p}(t - 24)$: - · - · - · .

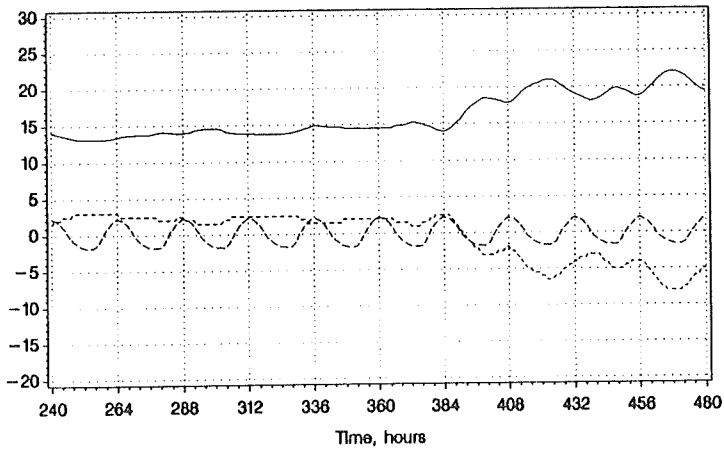


Figure 3.11: *Estimated components in Winters' seven-step-ahead predictor:* $\hat{\mu}(t - 7)$: — . $500 \times \hat{\beta}(t - 7)$: - - - - . $\hat{p}(t - 24)$: - · - · - · .

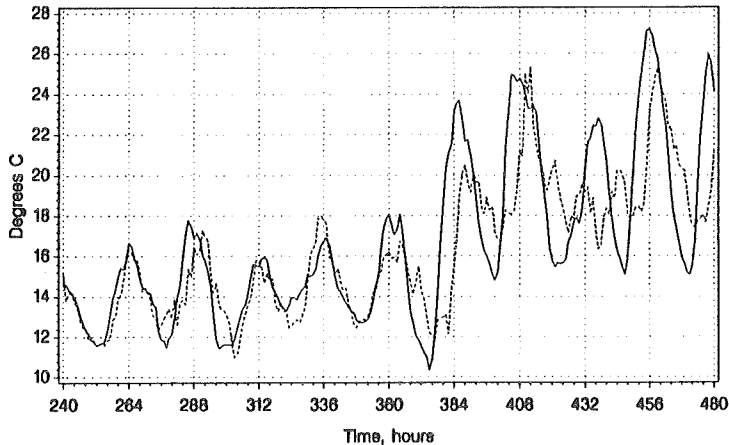


Figure 3.12: *Seven-step-ahead forecasting with procedure 2 (dashed curve) and real observations (solid curve)*

completely distinct as regards smoothness (cf. the above discussion).

So far it is mostly Winters' forecast procedure which has been discussed. A great deal of this discussion also applies to the other procedures. It is, however, reasonable briefly to mention procedure 2 which gave the best results. Figures 3.12 and 3.13 show plots for the seven-step-ahead predictor corresponding to those in Figures 3.9 and 3.11, respectively. Note that in Figure 3.13 the scaling factor multiplied with the seasonal component has been plotted instead of the scaling factor itself. The seven-step-ahead predictor of procedure 2 belongs to the quickly updating ones (cf. Winters' six-step-ahead predictor). The scaling factor shows very large variations, and for $376 \leq t \leq 480$ the level correction shows considerable variation too. Until $t \approx 386$ the scaling factor is less than one, but the growing amplitude of the real observations implies that subsequent values become larger than one. The increasing scaling factor tends to increase the mean of the forecasts, but this tendency is partly compensated by

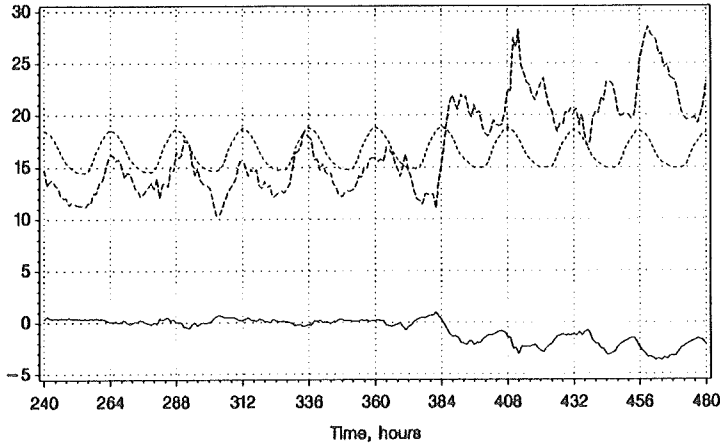


Figure 3.13: *Estimated components in the seven-step-ahead predictor of procedure 2: $\hat{\mu}(t-7)$: ——— . $\hat{p}(t-24)$: - - - - - . $\hat{f}(t-7)\hat{p}(t-24)$: - · - · - .*

the negative level correction. A certain increase in the mean level is in fact desirable since the mean level of the real observations increases as well.

The estimated seasonal profiles in Figures 3.10, 3.11 and 3.13 are almost stable in time. The reason is that the smoothing constants used for the updating of the seasonal components are very small (see α_3 in Table 3.3 and α_2 in Table 3.4) combined with the fact that each of the 24 elements in the seasonal profile is only updated once a day, i.e. about $1127/24 \approx 47$ times during the entire series. It is seen from Tables 3.3 and 3.4 that the strong smoothing of the seasonal profile takes place for almost all horizons in Winters' forecast procedure and procedure 1 and 2. One of the reasons why it is optimal to smooth that strong is that all the observations is used for estimation of the initial components in the profile. Experiments show that less than 100 of the first observation should be used for the initialization in

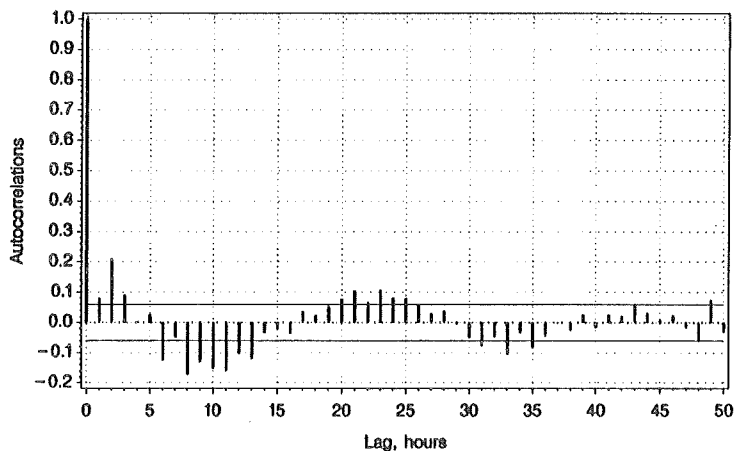


Figure 3.14: *Sample autocorrelations for the one-step-prediction errors produced by Winters' forecast procedure.*

order to make a stronger updating of the seasonal profile optimal. However, using that few observations also imply that the standard deviation of the prediction errors increases considerably.

Diagnostic Checks

In order to test whether the one-step-ahead predictors of the various forecast procedures are able to filter the sequence of observations to white noise, the simple diagnostic checks described in Section 3.1.1 are applied. Figure 3.14 shows the sample autocorrelations for the residuals obtained with Winters' one-step-ahead predictor. As it appears from the figure, the autocorrelations exceed the 95% confidence limits at several lags. Consequently, it can be concluded that the residual sequence is not white noise. The average of the one-step-ahead forecast errors is $-0.001\text{ }^{\circ}\text{C}$ which corresponds to approximately 9% of its standard deviation. This means that the one-step-ahead

forecasts apparently are unbiased.

The diagnostic checks applied to the forecast procedures 1 to 4 lead to the same conclusion: The one-step-ahead forecast errors are correlated, but unbiased. The sample autocorrelation functions of the residual sequences from procedures 1 to 4 and Winters' forecast procedure are very similar to each other (cf. Figure 3.14).

Remarks on Detection of Sudden Changes

The exponential smoothing procedures considered in the previous sections base their forecasts entirely on historical observations of the air temperature. As the past air temperatures do not give any advance information of a sudden change in the future weather situation, such forecast procedures can only detect the changes when they are observed at the location. Therefore the efficiency of a forecast procedure depends very much on its ability to adapt to the new situation. Exponential smoothing is generally well suited for processes exposed to slow changes, but not to processes with abrupt changes. The introduction of a scaling factor in the smoothing equations results in some increase of the adaptability to sudden changes. The drawback of the scaling factor idea in the previously introduced form is, however, that the forecasts depend on the zero point of the observations. This becomes apparent if the air temperature varies around zero. Then the estimated profile also varies around zero, implying that division by a small figure in the updating equation for the scaling factor very often occur (see e.g. procedure 2). This may lead to undesirable instabilities.

Combining Forecasts

Recent research shows that the scaling effect is obtainable in a more robust manner. The method used to obtain this effect is closely related to techniques by which the forecasts are calculated by combining forecasts produced by various methods (see, e.g., Makridakis *et al.* (1984), Newbold and Granger (1974)). Assume that k -step-ahead forecasts $\hat{y}^{(1)}(t+k|t), \hat{y}^{(2)}(t+k|t), \dots, \hat{y}^{(m)}(t+k|t)$ from m different forecast methods are available. Then one way to obtain a combined forecast is to compute a simple average,

$$\hat{y}(t+k|t) = \frac{1}{m} \sum_{i=1}^m \hat{y}^{(i)}(t+k|t).$$

This method does not pay attention to the mutual correlations between the m forecasts, but a weighted average can be used to cope with this problem (Makridakis *et al.* (1984)),

$$\hat{y}(t+k|t) = \sum_{i=1}^m w_i \hat{y}^{(i)}(t+k|t), \quad \sum_{i=1}^m w_i = 1,$$

where the weights, w_i ($i = 1, \dots, m$), are based on the sample covariance matrix of the prediction errors of the m methods. Once the weights are determined, they are used throughout the data series.

In the following it is illustrated how to make the weights in the linear combination of forecasts adaptive. Consider the following updating equations for a level and a seasonal profile:

$$\begin{aligned} \bar{\mu}(t) &= \alpha_1 y(t) + (1 - \alpha_1) \bar{\mu}(t-1) \\ \bar{p}(t) &= \alpha_2 y(t) + (1 - \alpha_2) \bar{p}(t-s). \end{aligned} \tag{3.36}$$

Note that these equations are decoupled versions of the corresponding equations in procedures 1 and 3 (simple exponential smoothing, Abraham and Ledolter (1983)). By using the statistics $\bar{\mu}(t)$ and

$\bar{p}(t)$, two different k -step-ahead forecasts are obtained independently: $\hat{y}_1(t+k|t) = \bar{\mu}(t)$ and $\hat{y}_2(t+k|t) = \bar{p}(t+k-s)$ ($k = 1, \dots, s$). These forecasts are used as regressors in a linear regression model,

$$y(t+k) = A\bar{\mu}(t) + B\bar{p}(t+k-s) + e(t), \quad k = 1, \dots, s, \quad (3.37)$$

where A and B are unknown coefficients. In order to estimate the coefficients, recursive least squares estimation with a forgetting factor λ is applied. By doing so, the resulting k -step-ahead prediction equation becomes

$$\hat{y}(t+k|t) = \hat{A}(t)\bar{\mu}(t) + \hat{B}(t)\bar{p}(t+k-s), \quad k = 1, \dots, s. \quad (3.38)$$

The smoothing parameters, α_1 , α_2 and λ , are found individually for each prediction horizon by minimization of the SSE criterion. Note that $\hat{B}(t)$ replaces the scaling factors in the forecasting procedures 2 and 4. $\hat{A}(t)$ ensures that the mean level of the forecasts is adjusted when the estimate of $\hat{B}(t)$ changes. The updating equations in (3.36) are combined through the regression equation (3.37). Since the coefficients are estimated by recursive least squares, the correlation between the level and the seasonal profile is considered. This very important quality was not present in the previously described procedures.

The predictor in (3.38) has been used to forecast the air temperature observations. The standard deviations of the forecast errors (standard errors) are shown in Figure 3.15. The predictor provides less accurate forecasts than the other forecast procedures for prediction horizons shorter than 5 hours. For horizons longer than 5 hours, however, the predictor offers significant improvements. This result is due to the regression coefficients being updated very strongly ($0.8 \leq \lambda \leq 0.9$). The updating of $\bar{\mu}(t)$ and $\bar{p}(t)$, on the other hand, is rather weak (small values of α_1 and α_2). Thus the level and the seasonal profile are relatively constant in time, and Equation (3.38) becomes an example of adaptive forecasting rather than combined forecasting.

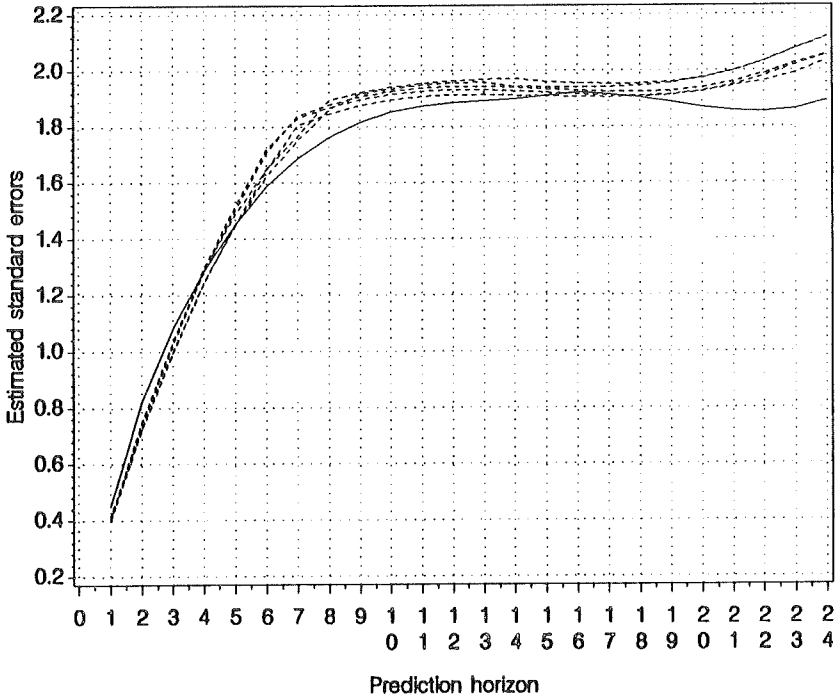


Figure 3.15: Standard errors versus the prediction horizon. Solid curve: Predictor in (3.38). Dashed curves: Previously described forecast procedures (cf. 3.4).

However, the approach combining exponential smoothing and recursive least squares regression can be considered a special case of a more general approach. A regression model similar to (3.37) can be used to form an adaptive combination of forecasts produced by several different models (e.g. exponential smoothing, ARIMA models, Box-Jenkins models, non-linear models etc.). The current estimates of the regression coefficients express how much weight should be attributed to the various forecasts when computing the final aggregate forecast. This approach is probably well suited for prediction of many types of processes which is governed by models with embedded variation of the parameters.

3.1.4 Conclusion

In this section various exponential smoothing procedures have been used to predict the air temperature. First Winters' seasonal forecast procedure was considered, and next four alternative forecast procedures being modifications of Winters' procedure were proposed. Two of the alternative procedures were non-linear and have not been seen previously in the literature. A seasonal IMA model was also identified, and finally an approach for combining forecasts was suggested.

A crucial advantage of the exponential smoothing procedures is that they are very simple as regards implementation and computations. Winters' procedure is based on a model containing three components: a level, a trend and a seasonal profile. In order to follow the time-variations of these components, they are updated at each sampling instant through the use of three recursive equations with associated smoothing constants. The four alternative procedures include a level and a seasonal component, and in addition to this, two of them include a scaling factor for the seasonal profile. Winters' procedure and the two alternative procedures without scaling factors are linear in

the components, and equivalent seasonal IMA models can be derived from these procedures. This derivation shows that the linear forecast procedures can be viewed as forecast models with embedded models of the level, (the trend,) and the seasonal component. The embedded models are IMA models with time-invariant parameters.

The smoothing constants were found separately for each individual prediction horizon. Since the forecast procedures only contain two or three parameters (smoothing constants), it is important that they be chosen optimally, i.e. by minimizing the prediction error variance. But in practice the use of individual sets of smoothing constants and smoothed components for each prediction horizon may not be desirable (e.g. due to the memory requirements). In such cases it might be more convenient to use a single set of smoothing constants by optimizing a multi-step criterion and thus obtain a forecast procedure that applies to all desired horizons.

It is important that the actual choice of smoothing constants defines a stable predictor filter. For simple exponential smoothing the stability is ensured if the smoothing constant is between zero and two, i.e. if the pole of the predictor filter is between minus one and one. For forecast procedures with more than one smoothing constant and coupled updating equations, however, it is not quite as easy to establish the stability condition since the poles of the predictor filter are complex functions of the smoothing constants.

Since sudden changes of both level and amplitude of the air temperature can occur, it is important that the forecast procedures adapt very fast to the changes. Procedures including a scaling factor therefore gave the best results.

In case of short prediction horizons (one to nine hours ahead), a traditional seasonal ARIMA model (in this case a seasonal IMA model) provided better results than all the smoothing procedures. When

predicting further ahead, the most accurate forecasts were obtained by the smoothing procedures.

For prediction horizons above five hours the most promising results were obtained by a method which combined forecasts from two independent exponential smoothing procedures: a simple exponential smoothing procedure and a simple seasonal exponential smoothing procedure. By this method the resulting forecasts are obtained as a weighted sum of the two separate forecasts, and the weights are estimated adaptively using a recursive least squares algorithm with a forgetting factor. The combination of forecasts in this way results in a quick adaption to sudden changes in the level and the diurnal amplitude of the air temperature. Though, the optimal smoothing of the level and the seasonal profile turns out to be very strong indicating that the resulting forecasts should be considered as adaptive forecasts rather than combined forecasts.

3.2 Models Relating the Variation in the Air Temperature to Other Climate Variables

This section presents linear stochastic models in continuous time relating the variations of the air temperature to other climatic variables – primarily the net radiation. The models are parametrized through physical interpretable parameters which are so-called embedded parameters. In the model building phase and the parameter estimation phase, ideas are taken from Madsen (1985) and Madsen, Thyregod and Holst (1987) as a starting point.

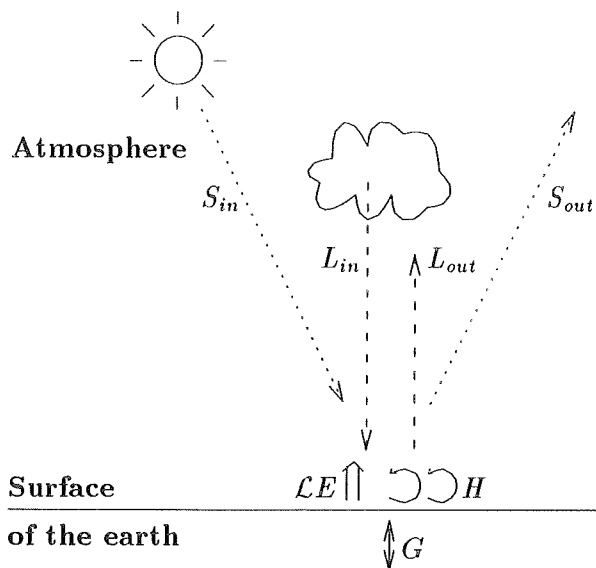


Figure 3.16: *Illustration of the energy balance at the surface of the ground.*

3.2.1 Formulation of Models

The energy source for the climatic system is the incident solar radiation. At the surface of the ground, the sum of the incoming short-wave radiation (S_{in}) from the sun and the long-wave radiation (L_{in}) from the clouds and the atmosphere makes up the positive contribution to the net radiation, R_n . The reflected part (S_{out}) of the incoming short-wave radiation together with emitted long-wave radiation (L_{out}) from the ground make up the negative contribution to the net radiation (see the illustration in Figure 3.16). Consequently

the net radiation can be expressed as

$$R_n = S_{in} + L_{in} - S_{out} - L_{out} .$$

Thus the net radiation is the part of the incoming radiation which is not returned to space/atmosphere as outgoing radiation. During night the short-wave radiation vanishes which implies that R_n most frequently becomes negative.

Since no accumulation of energy takes place at the surface of the ground, the net radiation energy is balanced by a net heat flux away from the surface. In this way the surface energy balance is maintained. The net heat flux is a sum of three different heat fluxes: the latent heat flux (flux of aqueous vapour), $\mathcal{L}E$, the sensible heat flux, H , and the soil heat flux, G . The following equation states the energy balance:

$$R_n = \mathcal{L}E + H + G . \quad (3.39)$$

The sensible and the latent heat fluxes are called turbulent fluxes.

Modelling Heat Fluxes Using the Resistance Method

This subsection about the resistance method summarizes some results from Berkowicz and Prahm (1982).

The sensible heat flux is closely related to the vertical gradient of the air temperature, T . Correspondingly, the latent heat flux is related to the vertical gradient of the humidity, q . The humidity is here expressed by the aqueous vapour pressure. The following gradient-transfer approximations are most widely used in atmospheric research:

$$H = -K_T(z) \frac{\partial T}{\partial z} \quad (3.40)$$

$$\mathcal{L}E = -K_q(z) \frac{\partial q}{\partial z} , \quad (3.41)$$

where $K_T(z)$ and $K_q(z)$ are diffusivities and z the height. The diffusivities depend on the flux which is considered and on the properties of the medium.

In the lowest part of the atmospheric boundary layer (the surface layer) the turbulent flux remains constant. Under these conditions, the so-called resistance method can be applied. The term “resistance method” refers to the analogy with Ohm’s first law,

$$\text{Current} = \frac{\text{Potential difference}}{\text{Resistance}} .$$

The resistance equations equivalent to (3.40) and (3.41) are

$$H = \frac{T(z_0) - T(z_1)}{r_h} \rho c_p \tag{3.42}$$

and

$$\mathcal{L}E = \frac{q(z_0) - q(z_1)}{r_w} \frac{\rho c_p}{\gamma} , \tag{3.43}$$

where ρ , c_p and γ are the air density, the specific heat of the air at a constant pressure and the psychrometric constant, respectively. Furthermore, r_h and r_w are aerodynamic resistances to fluxes of heat and aqueous vapour from height z_0 to height z_1 . By definition z_0 is the height of roughness where the wind speed is zero. It is assumed that both the temperature and humidity are measured at the same height, z_1 . There are simple relationships between the diffusivities in (3.40)/(3.41) and the constants in (3.42)/(3.43):

$$\begin{aligned} r_h &= \rho c_p \int_{z_0}^{z_1} \frac{1}{K_T(z)} dz \\ r_w &= \frac{\rho c_p}{\gamma} \int_{z_0}^{z_1} \frac{1}{K_q(z)} dz . \end{aligned}$$

The transport of aqueous vapour from the surface to the nearby surroundings complies with

$$\mathcal{L}E = \frac{Dq(z_0)}{r_s} \frac{\rho c_p}{\gamma} . \tag{3.44}$$

Here $Dq(z_0) = q_s(T(z_0)) - q(z_0)$ is the humidity deficit of which $q_s(T(z_0))$ is the saturated vapour pressure at the temperature $T(z_0)$ and $q(z_0)$ is the actual vapour pressure at the surface. The surface resistance, r_s , depends on the surface conditions which are characterized by vegetation and the water content of the soil. The surface resistance is not influenced by atmospheric turbulence.

Turbulent diffusion governs the transport of heat as well as aqueous vapour. Therefore the assumption

$$r_h = r_w \equiv r_a \quad (3.45)$$

is reasonable. The aerodynamic resistance r_a depends on wind speed and the roughness of the surface. Berkowicz and Prahm (1982) give a scheme for calculation of r_a . From Equations (3.39) and (3.42)-(3.45) the Penman-Monteith formula for the latent heat flux is achieved:

$$\mathcal{L}E = \frac{(R_n - G)\Delta r_a \gamma^{-1} + [q_s(T(z_1)) - q(z_1)]\rho c_p \gamma^{-1}}{r_s + (1 + \Delta \gamma^{-1})r_a}, \quad (3.46)$$

where

$$\Delta = \frac{\partial q_s(T)}{\partial T}. \quad (3.47)$$

Combining (3.39) and (3.46), the corresponding expression of the sensible heat flux is found as

$$H = \frac{(R_n - G)(r_a + r_s) + [q_s(T(z_1)) - q(z_1)]\rho c_p \gamma^{-1}}{r_s + (1 + \Delta \gamma^{-1})r_a}, \quad (3.48)$$

At this stage it would be appropriate to add a few comments to the expressions of the latent and the sensible heat fluxes in (3.46) and (3.48). The expressions are reached by considering static conditions. Consequently they have static nature, and do not describe the dynamic properties of the heat fluxes. Furthermore the static considerations are based on several approximations. The resistances r_a and r_s are generally time-varying. The aerodynamic resistance, r_a ,

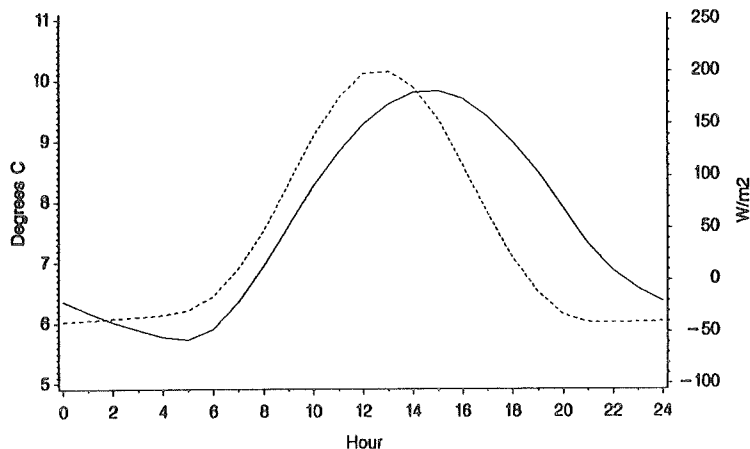


Figure 3.17: The mean diurnal variations of air temperature (solid curve) and net radiation (dashed curve).

is affected by a time-varying wind speed while the surface resistance, r_s , is determined by varying surface conditions (wetness and vegetation). In spite of (3.46) being only a static relationship, it will be used as an approximation in some of the dynamic models formulated later in this chapter.

Dynamic Relation between Air Temperature and Net Radiation

Madsen (1985) argues that a model of the variations of air temperature should be dynamic and have the net radiation as input. The arguments include a comparison of the mean diurnal and annual variations of the air temperature with the corresponding mean variations of the net radiation. Figures 3.17 and 3.18 show reproductions of graphs used in the comparison. The data used for the graphs are

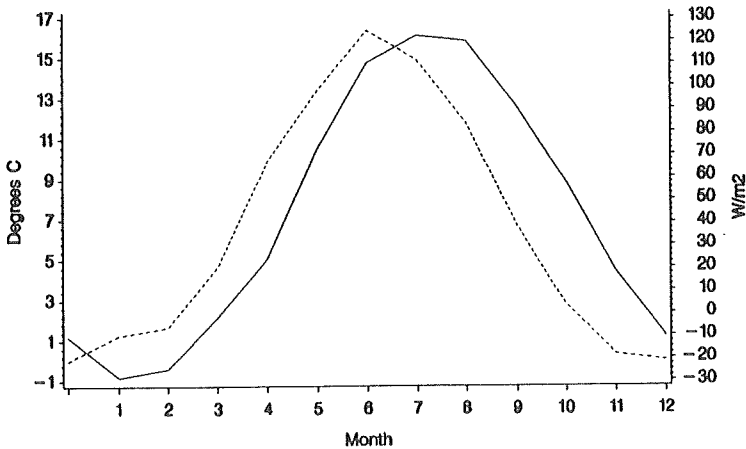


Figure 3.18: *The mean annual variation of air temperature (solid curve) and net radiation (dashed curve).*

69384 hourly observations of air temperature and the net radiation (see the brief description of the data in Section 3.2.2). Average values of the air temperature and the net radiation are calculated for each hour of the day (Figure 3.17). Corresponding average values are calculated for each month of the year (Figure 3.18). In both cases, the variations of the air temperature follow the variations of the net radiation with some time-delay. This indicates that the air temperature is driven by the net radiation through a dynamic system. Madsen (1985) draws attention to the fact that a second order linear, dynamic (possibly time-varying) system with net radiation as input would explain a great deal of the variation of the air temperature. The phase shifts indicate that, as a first approximation, the system should contain time constants of about 2 hours and 50 days for the diurnal and the annual variations, respectively (Madsen (1985)). The diurnal variations of the air temperature, which are approximately governed by a first order linear, dynamic system (time constant \approx 2 hours), exhibit a very typical asymmetric shape: The rise of the

temperature to its maximum happens significantly faster than the fall to its minimum. This is, however, the expected response from a first order system, driven by an input (net radiation) which is almost is the upper half of a sine wave plus an offset.

Assuming that the system governing the air temperature is properly approximated by a second order dynamic model, it is most likely that the dynamics can be modelled as two main heat reservoirs in the earth-atmosphere system. In the real system, however, the heat capacities are distributed in the land masses, the oceans and the air. This means that a detailed model would involve partial differential equations in time and the three spatial dimensions. An approximation which is frequently adopted for distributed dynamic systems is to lump the capacities into a finite number of point-capacities connected by resistances.

Madsen (1985) suggests that the smallest time constant is determined by the capacity of air and the resistances through which heat is exchanged with the surroundings. The arguments for this suggestion is based on the fact that the heat conductivity and the heat capacity of the soil are very low implying that the surface of the ground responds rather quickly to variations of the net radiation. The capacity of air is small and the exchange of heat between the surface and the air takes place through a small resistance. Hence changes of the net radiation affect the air temperature with a relatively small time constant (cf. Figure 3.17).

The largest time constant is attributed to the heat capacity of the sea. The sea constitutes an enormous heat reservoir to which heat is transferred through a small resistance (the surface layer of the sea is directly heated by the sun, and the energy is distributed very effectively by convection). Thus the diurnal variations of the sea temperature are negligible while the long-term effects of sun-heating imply delayed annual variation of the sea temperature. Because an

exchange of heat is going on between the sea and the air over land, the temperature of this air shows annual variations too (cf. Figure 3.18).

A Simple Linear, Dynamic Model for the Air Temperature

The above considerations lead to a simple linear, dynamic model of the relationship between the net radiation and the air temperature (over land). The dynamics are assumed to be time-invariant. The model is illustrated in Figure 3.19 together with an equivalent electric network. The currents and potentials in this network are heat fluxes and temperatures, respectively. The following symbols are used to denote temperatures, heat fluxes, capacities and resistances:

- T_a' Air temperature 2 m above the surface of the ground [$^{\circ}\text{C}$].
- T_{∞}' Constant temperature (probably somewhere in the atmosphere) [$^{\circ}\text{C}$].
- T_w' Temperature of the sea [$^{\circ}\text{C}$].
- R_n' Net radiation [W/m^2].
- c_a Heat capacity of the air [$\text{Wh}/^{\circ}\text{Cm}^2$].
- c_w Heat capacity of the sea [$\text{Wh}/^{\circ}\text{Cm}^2$].
- r_{wa} Resistance to the exchange of heat between the sea and the air over land [$^{\circ}\text{Cm}^2/\text{W}$].
- $r_{a\infty}$ Resistance to the exchange of heat between the air and the constant temperature T_{∞}' [$^{\circ}\text{Cm}^2/\text{W}$].

The long-term mean of the net radiation is balanced against a corresponding long-term mean of the heat flux through the resistance $r_{a\infty}$.

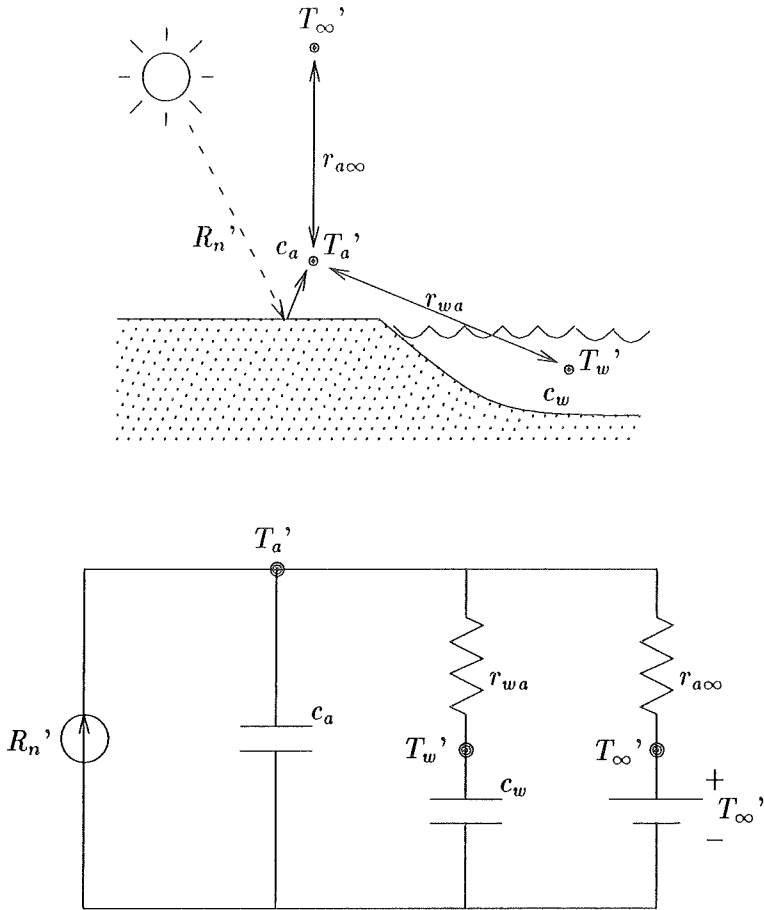


Figure 3.19: A simple linear dynamic model of the relationship between the net radiation and the air temperature.

The differential equations corresponding to the model is derived from Figure 3.19:

$$c_w \frac{dT_w'}{dt} = \frac{T_a' - T_w'}{r_{wa}} \quad (3.49)$$

$$c_a \frac{dT_a'}{dt} = \frac{T_{\infty}' - T_a'}{r_{a\infty}} + \frac{T_w' - T_a'}{r_{wa}} + R_n' \quad (3.50)$$

In this model it is assumed that the system can be characterized by four embedded parameters, namely the capacities c_a and c_w and the resistances r_{wa} and $r_{a\infty}$.

Linear Stochastic Models in Continuous Time

The continuous time model described by (3.49) and (3.50) is an approximation. Therefore noise terms are introduced in the model in order to account for the deviations from the real system. A two-dimensional stochastic differential equation is then obtained:

$$dT = ATdt + BR_n dt + dw(t), \quad (3.51)$$

where

$$\begin{aligned} T &= (T_w \ T_a)^T \\ A &= \begin{pmatrix} -\frac{1}{c_w r_{wa}} & \frac{1}{c_w r_{wa}} \\ \frac{1}{c_a r_{wa}} & -\frac{1}{c_a} \left(\frac{1}{r_{wa}} + \frac{1}{r_{a\infty}} \right) \end{pmatrix} \\ B &= (0 \ 1/c_a)^T \\ w(t) &= (w_1(t) \ w_2(t))^T. \end{aligned}$$

The two-dimensional stochastic process $w(t)$ is assumed to be a Wiener-process with the incremental covariance $R_1' dt$. Temperatures and fluxes are corrected by their long-term averages: $T_w = T_w' - \bar{T}_w$, $T_a = T_a' - \bar{T}_a$, $R_n = R_n' - \bar{R}_n$.

It is assumed that discrete time measurements of the air temperature and the net radiation are available. The measurements are inevitably superposed by noise (measurement errors). The error, $e(t)$, of the recorded air temperature, T_r' , is modelled as

$$T_r(t) = T_r'(t) - \bar{T}_r = CT(t) + e(t), \quad (3.52)$$

where $C = (0 \ 1)$ and $T_r(t)$ is the recorded temperature corrected by its long-term average. The measurement errors, $\{e(t)\}$, are assumed to be a sequence of independent $N(0, R_2)$ stochastic variables (Gaussian white noise with the mean zero). Furthermore, it is assumed that $e(t)$ and $w(t)$ are mutually independent. Since the model is linear, $e(t)$ is Gaussian, and $w(t)$ is a Wiener-process, the discrete time output, $T_r(t)$, is Gaussian too.

Equations (3.51) and (3.52) constitute the simplest input-output model dealt with in the following. The model belongs to a large class of models which will be termed *linear stochastic models in continuous time*. In the general case, the dimensions of the state vector, T , the input vector, U , ($U = R_n$ in (3.51)) and the observation vector, $Y(t)$, ($Y(t) = T_r(t)$ in (3.52)) correspond to the structure of the physical system considered. The dimensions of the matrices A , B , C , R_1' and R_2 are chosen accordingly. The general model becomes

$$dT = ATdt + BUdt + dw(t) \quad (3.53)$$

$$Y(t) = CT(t) + DU(t) + e(t), \quad (3.54)$$

where T , U and Y have mean value zero¹. The term $DU(t)$ where D is a matrix of suitable dimensions, has been included in (3.54) with the purpose of obtaining an entirely general formulation of the model ($D = 0$ in (3.52)). As an attempt to explore more of the

¹Corrections for the long-term averages are assumed.

physics behind the variations of the air temperature, this general model structure has been used as a basis.

As indicated in Figure 3.16 on page 153, only a certain part of the net radiation contributes directly to the heating of the air – namely the sensible heat flux. The latent heat flux only affects the air temperature if condensation of aqueous vapour takes place, and then the influence is mostly indirect. The soil heat flux influences the air temperature through a dynamic system, i.e. the soil. Therefore it would be desirable to decompose the net radiation into its individual heat flux components so that the net radiation in the model shown in Figure 3.19 can be replaced by the sensible heat flux as the driving “force”.

Since measurements of the soil heat flux are part of the data set used in the estimation phase later, a model with $R_n - G$ as input instead of R_n is suggested. According to (3.39), $R_n - G = \mathcal{L}E + H$ which implies that the sum of the latent and the sensible heat fluxes becomes input in this case – and not merely the sensible heat flux as originally intended. However, some improvement may still be brought about. The model is defined by

$$\mathbf{T} = (T_w \ T_a)^T \quad (3.55)$$

$$\mathbf{U} = (R_n \ G)^T \quad (3.56)$$

$$\mathbf{Y}(t) = T_r(t) \quad (3.57)$$

$$\mathbf{w}(t) = (w_1(t) \ w_2(t))^T \quad (3.58)$$

$$\mathbf{e}(t) = e(t) \quad (3.59)$$

$$\mathbf{A} = \begin{pmatrix} -\frac{1}{c_w r_{wa}} & \frac{1}{c_w r_{wa}} \\ \frac{1}{c_a r_{wa}} & -\frac{1}{c_a} \left(\frac{1}{r_{wa}} + \frac{1}{r_{a\infty}} \right) \end{pmatrix} \quad (3.60)$$

$$\mathbf{B} = \begin{pmatrix} 0 & 0 \\ 1/c_a & -1/c_a \end{pmatrix} \quad (3.61)$$

$$\mathbf{C} = (0 \ 1) \quad (3.62)$$

$$D = (0 \ 0) . \tag{3.63}$$

The only way this model differs from the model defined by (3.51) and (3.52) is that $R_n - G$ is used as input instead of R_n . Therefore, the values of the parameters, c_a , c_w , r_{wa} and $r_{a\infty}$, are also different.

With the object of creating a model which is closer to the physical system, the expression of $\mathcal{L}E$ in (3.46) is considered.

$$\mathcal{L}E = f(\Delta) ((R_n - G)\Delta r_a \gamma^{-1} + [q_s(T(z_1)) - q(z_1)]\rho c_p \gamma^{-1}) ,$$

where $f(\Delta) = 1/(r_s + (1 + \Delta\gamma^{-1})r_a)$ and $\Delta = \partial q_s(T)/\partial T$. This expression is not linear in the parameters r_a , r_s , γ , ρ and c_p . Therefore a polynomial series expansion of the function $f(\Delta)$ is introduced

$$f(\Delta) = \sum_{i=0}^{\infty} f_i \Delta^i . \tag{3.64}$$

If only the first (constant) term of this series expansion is used as an initial approximation, $\mathcal{L}E$ can be expressed as

$$\mathcal{L}E = \alpha_0(R_n - G)\Delta + \beta_0 Dq(z_1) , \tag{3.65}$$

where $Dq(z_1) = q_s(T(z_1)) - q(z_1)$ is the humidity deficit at the height z_1 . The original parameters in the expression of $\mathcal{L}E$ have been combined in two parameters α_0 and β_0 . If the two first terms of the series expansion (3.64) are included - i.e. by using $f(\Delta) \approx f_0 + f_1\Delta$ - the approximate expression of $\mathcal{L}E$ becomes

$$\mathcal{L}E = \alpha_0(R_n - G)\Delta + \alpha_1(R_n - G)\Delta^2 + \beta_0 Dq(z_1) + \beta_1 Dq(z_1)\Delta . \tag{3.66}$$

The data set contains observations of R_n , G , $q(z_1)$ and $q_s(T(z_1))$. Furthermore, Δ can be calculated from the observed air temperatures:

The saturated vapour pressure is a function of the air temperature. An empirical relationship is (Hansen *et al.* (1981)):

$$q_s(T_a) = e^{k_1 - \frac{k_2}{T_a + T_0}},$$

of which $k_1 = 26.042$, $k_2 = 5362.7$ and $T_0 = 273.2$ (T_a is stated in °C which implies that $T_a + T_0$ is in °K.). This gives

$$\Delta(T_a) = \left. \frac{\partial q_s(T)}{\partial T} \right|_{T=T_a} = \frac{k_2}{(T_a + T_0)^2} e^{k_1 - \frac{k_2}{T_a + T_0}} = \frac{k_2}{(T_a + T_0)^2} q_s(T_a),$$

To facilitate the notation, some new quantities are introduced

$$I_0 = (R_n - G)\Delta$$

$$I_1 = (R_n - G)\Delta^2$$

$$J_0 = Dq(z_1)$$

$$J_1 = Dq(z_1)\Delta.$$

Accordingly (3.65) and (3.66) are rewritten as

$$\mathcal{L}E = \alpha_0 I_0 + \beta_0 J_0 \quad (3.67)$$

$$\mathcal{L}E = \alpha_0 I_0 + \alpha_1 I_1 + \beta_0 J_0 + \beta_1 J_1. \quad (3.68)$$

The important thing about these approximate expressions of $\mathcal{L}E$, is that they are linear in the parameters (α_1 , β_1 , α_2 and β_2), and that I_0 , I_1 , J_0 and J_1 can be calculated from existing observations in the data set.

The above description of the latent heat flux can be used as embedded models of the energy balance at the surface of the ground in an extended model. This model and the equivalent electric network are shown in Figure 3.20. Along with this model a number of new quantities have been introduced (some of the previously assigned designations have been changed):

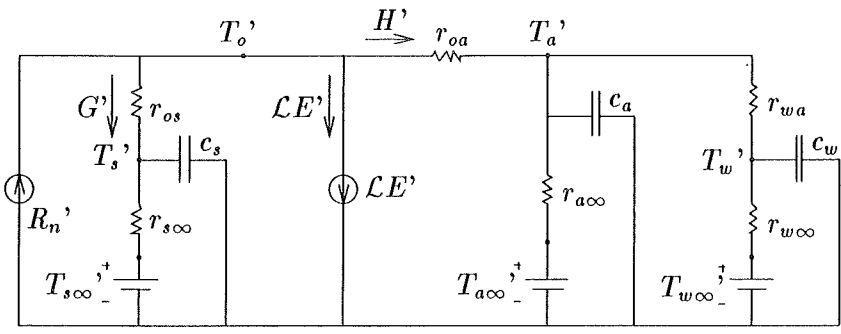
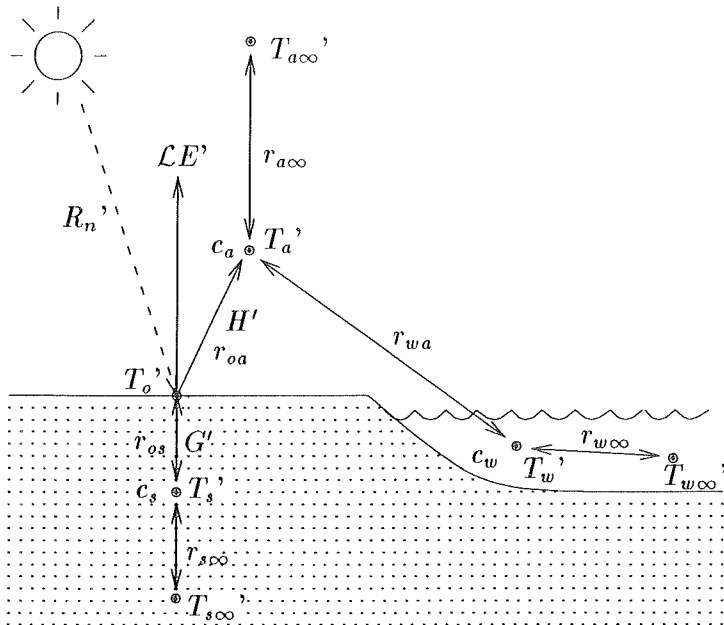


Figure 3.20: A linear dynamic model of the air temperature based on a partitioning of the net radiation in the latent, the sensible and the soil heat flux.

- T_s' Temperature of the upper part of the soil [$^{\circ}\text{C}$].
- $T_{s\infty}'$ Constant temperature of some deeper layers of the soil [$^{\circ}\text{C}$].
- T_o' Temperature of the surface of the ground [$^{\circ}\text{C}$].
- $T_{a\infty}'$ Constant temperature (probably somewhere in the atmosphere) [$^{\circ}\text{C}$].
- $T_{w\infty}'$ Constant temperature of some deep ocean [$^{\circ}\text{C}$].
- H' Sensible heat flux [W/m^2].
- $\mathcal{L}E'$ Latent heat flux [W/m^2].
- G' Soil heat flux [W/m^2].
- c_s Heat capacity of the upper part of the soil [$\text{Wh}/^{\circ}\text{Cm}^2$].
- r_{os} Resistance to the exchange of heat between the surface of the ground and the soil below [$^{\circ}\text{Cm}^2/\text{W}$].
- $r_{s\infty}$ Resistance to the exchange of heat between the upper part of the soil and the deeper layers of the soil [$^{\circ}\text{Cm}^2/\text{W}$].
- r_{oa} Resistance to the exchange of heat between the surface of the ground and the air [$^{\circ}\text{Cm}^2/\text{W}$].
- $r_{w\infty}$ Resistance to the exchange of heat between the sea and some deep ocean [$^{\circ}\text{Cm}^2/\text{W}$].

As it appears from the figure, the net radiation is partitioned into the latent, the sensible and the soil heat flux, in accordance with (3.39). The soil heat flux enters a first order dynamic system representing the soil. The temperature and the capacity of the upper part of the soil is T_s' and c_s , respectively. An exchange of heat with the deeper layers of the soil occur through the resistance $r_{s\infty}$. The amount (and direction) of heat transferred between the surface of the ground and the upper part of the soil is determined by the difference in temperature $T_o' - T_s'$. The latent heat flux is governed by (3.67) or (3.68) (both approximations will be used in the models). The surface of the ground and the air exchange heat through the resistance r_{oa} .

The rate of this exchange, i.e. the sensible heat flux, is determined by the difference in temperature $T_o' - T_a'$. The sensible heat flux enters a dynamic RC-network which is very similar to the one in Figure 3.19 apart from the capacity and the constant temperature that are introduced into the sub-network representing the sea. Notice that a possible loss (gain) of energy to (from) both the atmosphere and the ocean takes place through the resistances $r_{a\infty}$ and $r_{w\infty}$.

By means of Kirchhoff's law used for the nodes of the network represented by T_o' , T_s' , T_a' and T_w' , the equations describing the model are reached (all heat fluxes and temperatures are corrected for their long-term averages):

$$R_n = \frac{T_o - T_s}{r_{os}} + \mathcal{L}E + \frac{T_o - T_a}{r_{oa}} \quad (3.69)$$

$$c_s \frac{dT_s}{dt} = \frac{T_o - T_s}{r_{os}} - \frac{T_s}{r_{s\infty}} \quad (3.70)$$

$$c_a \frac{dT_a}{dt} = \frac{T_o - T_a}{r_{oa}} + \frac{T_w - T_a}{r_{wa}} - \frac{T_a}{r_{a\infty}} \quad (3.71)$$

$$c_w \frac{dT_w}{dt} = \frac{T_a - T_w}{r_{wa}} - \frac{T_w}{r_{w\infty}} \quad (3.72)$$

In (3.69) $\mathcal{L}E$ should be replaced by one of the expressions from (3.67) and (3.68). T_s , T_a and T_w are considered as the state variables of the model described by (3.69)-(3.72). As measurements of the air temperature and the soil heat flux are part of the data set, both quantities are considered as output variables of the model. The air temperature is one of the state variables while the soil heat flux is expressed indirectly by the state variables as

$$G = \frac{T_o - T_s}{r_{os}} \quad (3.73)$$

By solving (3.69) with respect to T_o and using the result to eliminate T_o in (3.70), (3.71) and (3.73), the model in its final form is readily reached. If the approximation in (3.67) is adopted, the model turns

into:

$$\mathbf{T} = (T_s \ T_a \ T_w)^T \quad (3.74)$$

$$\mathbf{U} = (R_n \ I_0 \ J_0)^T \quad (3.75)$$

$$\mathbf{Y}(t) = (T_r(t) \ G_r(t))^T \quad (3.76)$$

$$\mathbf{w}(t) = (w_1(t) \ w_2(t) \ w_3(t))^T \quad (3.77)$$

$$\mathbf{e}(t) = (e_1(t) \ e_2(t))^T \quad (3.78)$$

$$\mathbf{A} = \begin{pmatrix} -\frac{1}{c_s} \left(x_1 + \frac{1}{r_{s\infty}} \right) & \frac{x_1}{c_s} & 0 \\ \frac{x_1}{c_a} & -\frac{1}{c_a} \left(x_1 + \frac{1}{r_{wa}} + \frac{1}{r_{a\infty}} \right) & \frac{1}{c_a r_{wa}} \\ 0 & \frac{1}{c_w r_{wa}} & -\frac{1}{c_w} \left(\frac{1}{r_{wa}} + \frac{1}{r_{w\infty}} \right) \end{pmatrix} \quad (3.79)$$

$$\mathbf{B} = \begin{pmatrix} x_4 & -\alpha_0 x_4 & -\beta_0 x_4 \\ x_5 & -\alpha_0 x_5 & -\beta_0 x_5 \\ 0 & 0 & 0 \end{pmatrix} \quad (3.80)$$

$$\mathbf{C} = \begin{pmatrix} 0 & 1 & 0 \\ -x_1 & x_1 & 0 \end{pmatrix} \quad (3.81)$$

$$\mathbf{D} = \begin{pmatrix} 0 & 0 & 0 \\ x_2 & -\alpha_0 x_2 & -\beta_0 x_2 \end{pmatrix}, \quad (3.82)$$

where $G_r(t)$ is the recorded soil heat flux corrected for the long-term average. To simplify the notation, five auxiliary variables have been introduced:

$$\begin{aligned} x_1 &= \frac{1}{r_{oa} + r_{os}}, & x_2 &= r_{oa} x_1, & x_3 &= r_{os} x_1, \\ x_4 &= \frac{x_2}{c_s}, & x_5 &= \frac{x_3}{c_a}. \end{aligned}$$

If the approximation in (3.68) is used instead, the following changes should be observed:

$$\mathbf{U} = (R_n \ I_0 \ I_1 \ J_0 \ J_1)^T \quad (3.83)$$

$$\mathbf{B} = \begin{pmatrix} x_4 & -\alpha_0 x_4 & -\alpha_1 x_4 & -\beta_0 x_4 & -\beta_1 x_4 \\ x_5 & -\alpha_0 x_5 & -\alpha_1 x_5 & -\beta_0 x_5 & -\beta_1 x_5 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix} \quad (3.84)$$

$$\mathbf{D} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ x_2 & -\alpha_0 x_2 & -\alpha_1 x_2 & -\beta_0 x_2 & -\beta_1 x_2 \end{pmatrix}. \quad (3.85)$$

Estimation of the Parameters

The parameters of the models above formulated are estimated using the maximum likelihood method. The estimation is based on hourly observations of the input and output of the models. In short, the maximum likelihood method finds the values of the model parameters that maximize the joint probability density of all the observed output values. The function of the parameters which is maximized is called a likelihood function.

In the following it is briefly described how to derive the maximum likelihood estimator corresponding to the model in (3.53) and (3.54). A more detailed description can be found elsewhere (e.g. Melgaard and Madsen (1992) or Madsen and Melgaard (1991)).

Since the input and output are observed at discrete times ($t = 1, 2, \dots, N$), the continuous time model has to be sampled in order to compute the likelihood function. This sampling transformation is obtained by integrating the differential equation (3.53) through each sampling interval (1 hour). The result is a stochastic difference equation (in state-space form),

$$\mathbf{T}(t+1) = \Phi \mathbf{T}(t) + \Gamma \mathbf{U}(t) + \mathbf{v}(t). \quad (3.86)$$

Here $\mathbf{v}(t)$ is the sampled version of the Wiener-process, $\mathbf{w}(t)$, i.e. Gaussian white noise with the mean zero and the covariance \mathbf{R}_1 . It

is assumed that the input remains constant ($= U(t)$) through the sampling interval $[t, t + 1]$.

Let θ be a vector of the unknown parameters in the model – i.e. the physical parameters included in the definitions of the matrices A , B , C and D plus the covariance matrices R_1^t and R_2 . Furthermore, let $\mathcal{Y}(t)$ be a vector containing the observed output up to and including time t

$$\mathcal{Y}(t) = (Y(t) Y(t-1) \cdots Y(1) Y(0)),$$

and let $\hat{Y}(t|t-1)$ denote the conditional mean $E[Y(t)|\mathcal{Y}(t-1), \theta]$ and $R(t)$ the corresponding conditional covariance², $\text{Var}[Y(t)|\mathcal{Y}(t-1), \theta]$. Given the entire set of observed output values, $\mathcal{Y}(N)$, and the parameter vector, θ , the logarithm of the conditional likelihood function (conditioned on $Y(0)$) will turn out as

$$\log L(\mathcal{Y}(N); \theta) = -\frac{1}{2} \left(\sum_{t=1}^N \log \det R(t) + \sum_{t=1}^N \varepsilon^T(t) R^{-1}(t) \varepsilon(t) \right) + \text{const.}$$

where $\varepsilon(t) = Y(t) - \hat{Y}(t|t-1)$ ($t = 1, 2, \dots, N$) are the residuals of the output process.

Together with (3.54), (3.86) makes up a linear stochastic input-output model in a state-space formulation. Hence the conditional mean $\hat{Y}(t|t-1)$ and the conditional covariance $R(t)$, which are required in order to evaluate the log-likelihood function, are easily computed recursively by means of an ordinary Kalman filter. The Kalman filter is described in Appendix A.

²Since $Y(t)$ is Gaussian distributed, the conditional distribution $p(Y(t)|\mathcal{Y}(t-1), \theta)$ is Gaussian too. Consequently this distribution is completely characterized by its mean and covariance.

The hourly observations of input and output are not instantaneous measurements but hourly averages. Thus the observations recorded at time t are averages computed from instantaneous measurements recorded during the time interval $[t - 1, t[$. This means that the average observation of the output at a given time is already affected by the average input recorded at same time. In (3.86), however, it is assumed that the input affects the output after some delay. In order to solve this problem, (3.86) is modified as

$$\mathbf{T}(t + 1) = \Phi\mathbf{T}(t) + \Gamma\mathbf{U}(t + 1) + \mathbf{v}(t) . \quad (3.87)$$

The Kalman filter must be modified correspondingly.

The stability of the state-space models are ensured by the fact that they are based on RC-networks (R=resistance, C=capacity) which are known to define stable systems (provided that resistances and capacities are positive, of course).

It is impossible to estimate both covariances \mathbf{R}'_1 and \mathbf{R}_2 without imposing any restrictions on their structures. If it is assumed that \mathbf{R}'_1 is $n \times n$ and \mathbf{R}_2 is $s \times s$, then, taking into account that covariances are symmetric, the total number of different covariance elements is $\frac{1}{2}(n^2 + n) + \frac{1}{2}(s^2 + s)$. By considering the innovations form of the state-space model, however, it becomes clear that for identifiability the covariance elements must be parametrized through $ns + \frac{1}{2}(s^2 + s)$ free parameters at the most (see e.g. Ljung (1987)). A simple, yet frequently used way to fulfil this requirement is to assume that the covariances are diagonal.

Having found an estimate, $\hat{\boldsymbol{\theta}}$, of the parameters, an assessment of the accuracy is of great interest. This assessment is based on the fact that the maximum likelihood estimator is asymptotically $N(\boldsymbol{\theta}, \Sigma_{\boldsymbol{\theta}})$ distributed, where $\Sigma_{\boldsymbol{\theta}} = \mathbf{H}^{-1}$. \mathbf{H} is defined as

$$\mathbf{H} = -E \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^2} L(\mathcal{Y}(N); \boldsymbol{\theta}) \right] .$$

An estimate of Σ_θ is obtained by replacing the expectation by the observed value of the actual estimate of θ

$$\hat{V}[\hat{\theta}] = \hat{\Sigma}_\theta = - \left(\frac{\partial^2}{\partial \theta^2} L(\mathcal{Y}(N); \theta) \right) \Big|_{\theta = \hat{\theta}}$$

Since it is impossible to derive any analytical expression of the maximum likelihood estimator, the log-likelihood function has to be maximized by a numerical method. A computer software package, CTLSM (Continuous Time Linear Stochastic Modelling), is able to perform this maximization. A description of CTLSM is found in Melgaard and Madsen (1993).

3.2.2 Experimental Results

In this section the parameters of the above formulated dynamic models are estimated. Before the results are presented, a brief description of the data is given.

The Data Set

The data set consists of measurements made at a climate and water balance station, Højbakkegård, which belongs to the Royal Danish Veterinary and Agricultural University, Denmark. Højbakkegård is situated near Tåstrup, about 20 km west of Copenhagen. The climate station is freely exposed and surrounded by ordinary agricultural fields.

The output from the various measuring equipment (radiometers, thermometers etc.) was sampled every 10 minutes by a data-logger. From the sampled values hourly averages were computed. The data set

used here, which is a subset of an original data set, contains hourly averages of the net radiation, the air temperature, the soil heat flux, the actual vapour pressure and the saturated vapour pressure. The observations cover the period from 1 February 1966, to 31 December 1973, i.e. a period of 69384 hours.

The following is a brief description of the instantaneous measurements, the measuring equipments etc. See Hansen *et al.* (1981) for more details.

Net radiation [W/m²] The net radiation was measured by a polyethylene shielded radiometer placed one meter above the ground which was covered with short grass.

Air temperature [°C] The air temperature was measured by a resistance thermometer placed 2 m above the ground in a convective open box. The accuracy of the instantaneous measurements is ± 0.5 °C.

Soil heat flux [W/m²] The soil heat flux was measured by a heatflow-meter at a depth of 3 cm. The surface of the ground above the heatflow-meter was covered by short grass. The accuracy of the measurement is within the interval [-10%,+20%]. The measuring errors were primarily due to variations of the thermal conductivity of the surrounding soil.

Vapour pressure [Pa] The actual vapour pressure was computed from the dew point, which was measured by a LiCl dew point sensor placed in a convective open box 2 m above the ground. The

relationship between the vapour pressure, q , and the dew point, T_d , is

$$q = e^{26.042 - \frac{5362.7}{T_d}}, \quad (3.88)$$

where q is stated in Pa and T_d is in Kelvin. The accuracy is $\pm 4\%$.

Saturated vapour pressure [Pa] The saturated vapour pressure 2 m above the ground was computed from the air temperature at the same height. Equation (3.88) was used for this computation: q and T_d were replaced by the saturated vapour pressure and the air temperature, respectively.

Estimation Results

Four models are considered:

- **Model 1.** Described by Equations (3.51) and (3.52).
- **Model 2.** Described by Equations (3.55)-(3.63).
- **Model 3.** Described by Equations (3.74)-(3.82).
- **Model 4.** Described by Equations (3.74), (3.76)-(3.79), (3.81), (3.83), (3.84) and (3.85).

For each of these models the parameters included in the matrices A , B , C and D and the diagonal elements of the covariances R'_1 and R_2 have been estimated (elements outside the diagonals of the covariances have been fixed to zero). However, some of the estimates have turned out to be statistical insignificant. Therefore, the results given below have been found by re-estimation of the models with appropriate fixed values of insignificant parameters.

In the following the estimation results are presented, and in Section 3.2.3 they are discussed.

Model 1. Table 3.6 shows the estimation results of model 1. The following parameters have been estimated by maximization of the likelihood function: c_w , c_a , r_{wa} , $r_{a\infty}$, $\mathbf{R}'_{1,1,1}$ and $\mathbf{R}'_{1,2,2}$ ($\mathbf{R}'_{1,i,j}$ denotes the (i, j) th element of \mathbf{R}'_1). An attempt to estimate R_2 along with these parameters leads to an insignificant estimate. For that reason, R_2 is fixed at a value which is chosen from *a priori* knowledge of the measurement error. It is realistic to assume that the measurement errors connected to the hourly averaged values are within the interval $[-0.1^\circ\text{C}, 0.1^\circ\text{C}]$. Supposing that this interval corresponds to ± 3 times the standard deviation in a normal distribution, it is relevant to choose $R_2 = (0.1/3 \text{ }^\circ\text{C})^2 = 0.001111 \text{ }^\circ\text{C}^2$. However, the estimation results seem to be rather insensitive to the choice of R_2 (for $R_2 \in [1 \cdot 10^{-10}, 1 \cdot 10^{-3}]$ the variation of the other parameters is less than 0.5% and less than their standard deviations).

For purpose of comparison, Table 3.6 also shows the estimation results obtained by Madsen (1985) by the use of the same model structure and the same data. Due to limited computer capacity he considered the Kalman gain, $\mathbf{k}(t)$, to be constant, \mathbf{k} , throughout the time series, and estimated this “stationary” Kalman gain instead of the covariances \mathbf{R}'_1 and R_2 . Therefore Table 3.6 compares the stationary Kalman gains, $\mathbf{k}(\infty)$, instead of the covariances.

The CTLSM program also permits estimation of the stationary Kalman gain instead of the covariances. By doing so, the results found by Madsen (1985) can be reconstructed with close accuracy ($\pm 0.1\%$). Madsen (1985) excluded an unknown number of the first observations when evaluating the likelihood criterion, and this is probably the reason why minor deviations occur.

Table 3.6: *Estimation results of model 1. The figures in square brackets are the corresponding estimates found by Madsen (1985).*

Parameter	Estimate			Standard deviation	
c_w	382.0	[501.4]	Wh/°Cm ²	9.0	[25.9]
c_a	120.8	[130.9]	Wh/°Cm ²	0.9	[1.9]
r_{wa}	0.01776	[0.02169]	°Cm ² /W	0.00023	[0.00057]
$r_{a\infty}$	0.08651	[0.08096]	°Cm ² /W	0.00214	[0.00380]
$R'_{1,1,1}$	1.748		°C ² /h	0.037	
$R'_{1,2,2}$	0.06030		°C ² /h	0.00128	
R_2	0.001111		°C ²	Fixed	
$k_1(\infty)$	2.642	[2.328]			[0.045]
$k_2(\infty)$	0.994	[1.446]			[0.006]
λ_1	-0.6885	[-0.5218]	h ⁻¹		
λ_2	-0.0205	[-0.0166]	h ⁻¹		
τ_1	1.46	[1.92]	h		
τ_2	48.8	[60.1]	h		
Estimate of the one-step-ahead prediction error variance: 0.4460 ² [0.4686 ²] °C ²					

Comments:

- $R'_{1,i,j}$ denotes the (i, j) th element in R'_1 .
- Standard deviations of quantities which have not been estimated directly through the maximization of the likelihood function are not available.
- "Fixed" in the column "Standard deviation" means that the corresponding "Estimate" was fixed during the estimation.

By computing the eigenvalues, λ_i ($i = 1, 2$), of the estimated A -matrix in the continuous time model (3.53), estimates of the time constants in the dynamic system are obtainable as

$$\tau_i = -\frac{1}{\lambda_i}, \quad i = 1, 2.$$

The eigenvalues as well as the time constants are shown in Table 3.6.

Model 2. For model 2 the estimation results shown in Table 3.7 are achieved. As for model 1 the measurement error variance of the air temperature has been fixed: $R_2 = 0.001111 \text{ }^\circ\text{C}^2$.

Model 3. For model 3 the estimates shown in Table 3.8 are achieved. Note that fixed values ($1.0 \cdot 10^{-20}$) have been assigned to the covariance elements $R_{1,2,2}$ and $R_{2,2,2}$. The reason is that when estimated the estimates turn out to be small (less than $1.0 \cdot 10^{-10}$) and statistical insignificant. Replacing the estimates by fixed values does not have much influence on the rest of the parameter estimates: they remain unchanged with four significant digits.

Model 4. Estimation results of model 4 are shown in Table 3.9. The covariance elements $R_{1,2,2}$ and $R_{2,2,2}$ have been fixed due to statistical insignificance. The explanation given in connection with model 3 also applies here.

3.2.3 Discussion and Conclusion

It was expected that the estimates for model 1 would have been closer to the estimates found for the same model by Madsen (1985). Table

Table 3.7: Estimation results of model 2.

Parameter	Estimate		Standard deviation
c_w	306.4	Wh/°Cm ²	7.4
c_a	121.9	Wh/°Cm ²	1.0
r_{wa}	0.02021	°Cm ² /W	0.00028
$r_{a\infty}$	0.09496	°Cm ² /W	0.00230
$R'_{1,1,1}$	2.139	°C ² /h	0.047
$R'_{1,2,2}$	0.06456	°C ² /h	0.00111
R_2	0.001111	°C ²	Fixed
$k_1(\infty)$	2.841		
$k_2(\infty)$	0.995		
τ_1	1.58	h	
τ_2	45.3	h	
Estimate of the one-step-ahead prediction error variance: 0.4507 ² °C ²			

Comments:

- $R'_{1,i,j}$ denotes the (i, j) th element in R'_1 .
- Standard deviations of quantities which have not been estimated directly through the maximization of the likelihood function are not available.
- "Fixed" in the column "Standard deviation" means that the corresponding "Estimate" was fixed during the estimation.

Table 3.8: Estimation results of model 3.

Parameter	Estimate		Standard deviation
c_s	96.91	Wh/°Cm ²	1.58
c_a	59.16	Wh/°Cm ²	0.65
c_w	5497	Wh/°Cm ²	523
r_{os}	0.2219	°Cm ² /W	0.0010
$r_{s\infty}$	1.204	°Cm ² /W	0.061
r_{oa}	0.004021	°Cm ² /W	0.000077
$r_{a\infty}$	0.9651	W/°Cm ²	0.0097
r_{wa}	0.01444	°Cm ² /W	0.00016
$r_{w\infty}$	0.02345	°Cm ² /W	0.00204
α_0	0.005722	°C/Pa	0.000042
β_0	-0.4971	W/(m ² Pa)	0.0055
$R'_{1,1,1}$	0.4743	°C ² /h	0.0049
$R'_{1,2,2}$	1.0·10 ⁻²⁰	°C ² /h	Fixed
$R'_{1,3,3}$	0.3845	°C ² /h	0.0035
$R_{2,1,1}$	0.007602	°C ²	0.000219
$R_{2,2,2}$	1.0·10 ⁻²⁰	°C ²	Fixed
$k_{1,1}(\infty)$	0.934		
$k_{2,1}(\infty)$	0.934		
$k_{3,1}(\infty)$	1.607		
$k_{1,2}(\infty)$	-0.223	°Cm ² /W	
$k_{2,2}(\infty)$	0.003	°Cm ² /W	
$k_{3,2}(\infty)$	0.018	°Cm ² /W	
τ_1	0.783	h	
τ_2	19.2	h	
τ_3	127	h	
$\hat{\sigma}^2(\text{Pred. error for air. temp.}) = 0.3655^2 \text{ °C}^2$			

Comment: “Fixed” in the column “Standard deviation” means that the corresponding “Estimate” was fixed during the estimation.

Table 3.9: Estimation results of model 4.

Parameter	Estimate		Standard deviation
c_s	101.6	Wh/°Cm ²	0.5
c_a	26.89	Wh/°Cm ²	0.31
c_w	2278	Wh/°Cm ²	4
r_{os}	0.2144	°Cm ² /W	0.0010
$r_{s\infty}$	1.135	°Cm ² /W	0.034
r_{oa}	0.008543	°Cm ² /W	0.000176
$r_{a\infty}$	1.104	°Cm ² /W	0.001
r_{wa}	0.03227	°Cm ² /W	0.00038
$r_{w\infty}$	0.05664	°Cm ² /W	0.00004
α_0	0.01415	°C/Pa	0.00008
α_1	$-5.418 \cdot 10^{-5}$	°C ² /Pa ²	$5.2 \cdot 10^{-7}$
β_0	-0.3214	W/(m ² Pa)	0.0047
β_1	$7.301 \cdot 10^{-4}$	(W°C)/(m ² Pa ²)	$2.60 \cdot 10^{-5}$
$R'_{1,1,1}$	0.4573	°C ² /h	0.0052
$R'_{1,2,2}$	$1.0 \cdot 10^{-20}$	°C ² /h	Fixed
$R'_{1,3,3}$	0.4104	°C ² /h	0.0041
$R_{2,1,1}$	0.007359	°C ²	0.000246
$R_{2,2,2}$	$1.0 \cdot 10^{-20}$	°C ²	Fixed
$k_{1,1}(\infty)$	0.936		
$k_{2,1}(\infty)$	0.936		
$k_{3,1}(\infty)$	1.640		
$k_{1,2}(\infty)$	-0.220	°Cm ² /W	
$k_{2,2}(\infty)$	0.003	°Cm ² /W	
$k_{3,2}(\infty)$	0.031	°Cm ² /W	
τ_1	0.730	h	
τ_2	20.4	h	
τ_3	124	h	
$\hat{\sigma}^2(\text{Pred. error for air. temp.}) = 0.3614^2 \text{ °C}^2$			

Comment: "Fixed" in the column "Standard deviation" means that the corresponding "Estimate" was fixed during the estimation.

3.6 shows that the deviations are considerable. The largest differences are about 25% and occur in the estimates of c_w and τ_1 . However, the differences are probably due to different parametrization of the covariance structure of the noise. In model 1 the parametrization is done through the diagonal elements of the covariance matrices while it is done through the Kalman gain in the model presented by Madsen (1985). The reduction of the prediction error variance obtained by model 1 indicates that this model provides the best description of the covariance structure.

Model 1 leads to conclusions very similar to those drawn by Madsen (1985). Studying the phase shifts between the net radiation and the air temperature in Section 3.2.1 it was found that time constants of about 2 hours and 50 days should be expected. The small time constant in Table 3.6 ($\tau_1 = 1.45$ h) is quite close to the expected value. The large time constant ($\tau_2 = 48.8$ h) is, however, far from the expected 50 days. One reason for this can be that estimation of a large time constant in a system containing both small and large time constants is difficult. Therefore it can be concluded that the estimate of τ_2 does not, as expected, correspond to the heat capacity of the sea. The estimated value should rather be attributed to some deeper layers of the soil.

The prediction error variance for model 2 is larger than for model 1. Correspondingly, the value of the likelihood function is smaller. This indicates that model 2 does not describe the system quite as well as model 1, and it can be concluded that the use of the net radiation minus the soil heat flux as an input variable was no such good idea as expected. Notice, however, that the magnitudes of the estimates for the two models are the same.

The only difference between model 3 and model 4 is that the approximative expression of the latent heat flux used in model 4 is better than the one used in model 3. From Tables 3.8 and 3.9 it appears

that the different choice of approximation results in parameter differences of more than 100% for some of the parameters. Though, it also appears that the estimates of the time constants are not that far apart (about 7% at the most). The reduction of the prediction error variance from model 3 to model 4 is about 2%. Since model 4 seems to represent the best approximation no further comments shall be made as to the estimation results of model 3.

For model 4 the prediction error variance of the air temperature is $0.3614^2 \text{ } ^\circ\text{C}^2$. Hence, the model improves the one-hour-ahead forecasting performance considerably compared with model 1 and 2. The expected time constant of 2 hours presumably corresponds to a mixture of τ_1 and τ_2 since it lies in between them. This indicates that a second order model is more appropriate than a first order model for an approximation of the short term dynamics of the underlying distributed system governing the diurnal variations of the air temperature. But the large time constant of about 50 days does not have a counterpart among those in Table 3.9. As for model 1 the reason may be that it is difficult to estimate a large time constant in a system containing both large and small time constants. Furthermore, the estimate of a large time constant (large compared to the sampling interval) has a general tendency to be too small. Ljung (1987) points out that it is difficult to let one and the same model handle more than two or, at the most, three decades of the frequency range. If the system is stiff (as the present climatic system) so that it contains widely separated time constants of interest, he suggests to build two or more models, each covering a proper part of the frequency range and each sampled with a corresponding, suitable sampling interval. The choice of sampling interval represent one way of suppressing unwanted frequencies, another and more efficient way is to apply an appropriate prefilter.

The magnitudes of the estimated soil parameters for model 4 seem to be reasonable. A typical value of the heat capacity of soil is $\rho c = 700 \text{ Wh/m}^3\text{ } ^\circ\text{C}$ (Hansen *et al.* (1981)). Therefore the estimate of c_s

corresponds to an upper soil layer of thickness $\hat{c}_s/\rho c = 0.15$ m. Since the thermal conductivity of soil is typically about $\lambda_h = 1.3$ W/m²C, the estimate of r_{os} corresponds to a depth of $\hat{r}_{os}\lambda_h = 0.28$ m. Thus, the estimates of c_s and r_{os} are in reasonable accordance with figures in Hansen *et al.* (1981) which indicate that the diurnal temperature variation in the soil is measurable down to a depth of 0.25 - 0.30 m. For the observed air temperatures, the average heat capacity of atmospheric air is about 0.35 Wh/m³C. This means, the estimated value of c_a represents the heat capacity of air up to a height of about $27/0.35 = 77$ m. It is likely, however, that the estimated capacity includes part of the upper soil layer.

The autocorrelation function and the partial autocorrelation function of the residual sequence of the air temperature have been estimated for all four models, and the pattern is the same for all four instances. Due to the large number of observations, the autocorrelation function contain several significant correlations, and a white noise hypothesis can be rejected on any reasonable level. A great similarity between the autocorrelation function and the partial autocorrelation function indicates that there is not much information on the variations of the air temperature left in the residuals.

The cross-correlation functions between the residuals of the air temperature and the various input variables have also been estimated. Several significant cross-correlations are found for all input variables indicating that they could be used to explain more of the variation of the air temperature. The largest cross-correlations for model 4 (about 0.08 to 0.09, absolute value) occur for the input variable R_n at lags 10 to 12.

The above discussion suggests that model 4 provide the best description of the local climate system as regards the air temperature. However, the model is far from perfect. First of all, none of the estimated time constants are larger than 5.2 days although the real system con-

tains a time constant of about 50 days. Comparing model 4 with model 1 it is clearly seen that the distance between the smallest and the largest time constant becomes larger when the model order is increased from two to three. There is a general tendency that the time constants of a lumped model of a distributed system move apart when the model order is increased. Therefore, the "missing" time constant might be found if a model of suitably large order was specified. To check whether the resulting model contains too many parameters, a cross-validation technique can be applied (see e.g. Ljung (1987)).

A general problem of the continuous time state-space models estimated here is that there is no guarantee that the final estimates of resistances and capacities represent the particular parts of the real system that they were supposed to. If a lumped model structure of, say, third order is specified, the estimation method determines how the real distributed system should be lumped to make it optimal; and the result does not necessarily agree with the intended interpretation of the model. An effective way of solving this problem is to apply observations corresponding to all or at least more of the state variables in the model. For the present models, measurements of the sea temperature would be most interesting since the influence of the sea temperature is badly described by the models. Secondly, it would be beneficial to have observations of the temperature in the upper part of the soil.

All observations were corrected for their sample means before they were used for estimation of the parameters. This procedure is required if not all variables influencing the air temperature are taken into account. In the present case, the horizontal transport of energy is neglected, and this indicates a need for the correction by the sample mean. Alternatively, the original observations could have been used when estimating the parameters, thus obtaining models which focus more on the description of the static conditions. Notice that doing it that way would imply that the set of parameters to be estimated

would be extended, e.g. by T_{∞}' in model 1.

As mentioned earlier, some of the thermal resistances depend, among other things, on wind speed and surface conditions. Therefore there is a chance that the model could be improved if embedded models of the variations of these resistances were specified.

Chapter 4

Multi-Step and Embedded Model Based Control

DURING the two previous decades, self-tuning and adaptive control of dynamic systems have been subject to extensive research. Various techniques for adaptive identification/estimation have been combined with different types of controllers in order to obtain self-tuning and adaptive qualities (see, e.g., Davis and Vinter (1985) or Åström and Wittenmark (1989)). The research has contributed to a theoretical understanding of many aspects of adaptive control, which has proved to be useful in several practical applications.

This chapter deals with multi-step predictive control, and in a major part of the chapter, a new control method which permits variations of the embedded model parameters is considered. Below is first explained how the traditional adaptive control works, and next why it is of interest to study control methods permitting embedded parameter variations.

In the traditional case, an adaptive estimator is used to track the time-variations of the model parameters. At time, t , the estimator provides the controller with a local estimate, $\hat{\mathcal{D}}_{t|t}$, of the system dynamics based upon observations of input and output up to time t . If the controller is prediction based, as is the case with many controllers (e.g. minimal variance controllers), it needs predictions of the future output as a function of the present and perhaps the future input. This means that the future dynamics, \mathcal{D}_{t+j} ($j > 0$), of the system is as interesting as the present dynamics. Most frequently, however, the estimator only gives an estimate of the present dynamics, and the controller has to use this estimate as a forecast for the future as well – i.e. the controller assumes that $\hat{\mathcal{D}}_{t+j|t} = \hat{\mathcal{D}}_{t|t}$ ($j > 0$). If the time-variation of the system is not too quick and the input-output delay not too long, this is a reasonable solution.

For systems exhibiting quick time-variation compared to the sampling interval, it may be beneficial to specify explicit embedded models of the time-variation of the dynamics. Such embedded models would provide forecasts, $\hat{\mathcal{D}}_{t+j|t}$ ($j > 0$), of future changes in the system dynamics. In Section 2.2.2, as an example, the diurnal variation of the parameters in a transfer function model is described by trigonometrical functions which correspond to an embedded model of the diurnal variation of the time-delay. Through this model forecasts of future variations of the time-delay is available. Most traditional controllers, however, are not designed to use this information on future time-variation.

Predictive control is a rather general method that can handle embedded parameter variations. First, in Section 4.1, it is explained why multi-step predictive control is relevant seen from the district heating point of view. Next, in Section 4.2, a pseudo generalized predictive control that is called weighted predictive control is proposed. Section 4.3 reviews Generalized Predictive Control (GPC) in its traditional form as proposed by Clarke *et al.* (1987A). The subsequent Section

4.4 shows how to obtain controllers for embedded model structures through a generalized version of GPC, and finally simulation results obtained by use of GPC are given in Section 4.5.

4.1 Control Strategy for a District Heating System

In Chapter 2 the great importance of tracking time-delays in district heating systems was explained. In that chapter only time-delays and transfer functions between the district heating plant and individual sites in the distribution network were considered. However, in order to perform an optimal control of the supply temperature, dynamic models of the total demand for heat in the entire system are required as well. Various models belonging to this category are proposed by Madsen *et al.* (1990) and Madsen *et al.* (1992). Since a model of the total heat demand considers the entire distribution system simultaneously, it has to describe a lot of different time-delays – in principle one for each consumer of heat. The presence of several different delays in a single dynamic system makes many traditional controllers inapplicable since they assume a unambiguous delay. For such a system, it is appropriate to use multi-step predictive control, e.g. generalized predictive control.

As described in Chapter 1 a strategy for controlling the supply temperature in a district heating system may aim at keeping the supply temperature as low as possible while taking certain restrictions into consideration. A very important restriction is the limitation of the pumping capacity in the district heating network. That means that the total flow of water (mass flow) through the network is subject to an upper physical limit.

Normally, variations of the total heat demand can be met by varying the mass flow through the network or by varying the supply temperature. Since the purpose is to keep the supply temperature as low as possible, the heat demand is primarily met by variations of the mass flow. If the heat demand is low enough, the supply temperature is kept close its minimum (usually determined by hot tap water restrictions), while the mass flow is subject to variations. For a moderately high heat demand the mass flow occasionally reaches its maximum value during the day, and an increase of the temperature is needed. In case the quantity of heat required by the consumers becomes so large that the mass flow is close to its maximum throughout the day (which is usually the case for non-summer situations), then variations of the temperature is the only way to meet variations of the heat demand. This is illustrated in Figure 4.1.

Note that an increasing heat demand can be satisfied immediately by increasing the mass flow (if it is not at its maximum value already). If an increasing heat demand should be met by a higher temperature, the temperature must be raised some time before the demand actually increase. This is due to the transport time (time-delay) from the production plant to the consumers.

Suppose that the heat consumption varies as in situation 3 in Figure 4.1, which normally is the case in non-summer periods. This means that a strategy to control the supply temperature is needed to ensure that the consumers are supplied with enough heat at any time. The control strategy must, of course, take into account that the heat is supplied with some delay, and consequently it must make use of predictions of the heat demand. This is the major reason for using predictive control of the supply temperature in district heating systems.

The time-delay from the production plant to the various consumers varies in accordance with the siting of the individual consumer in

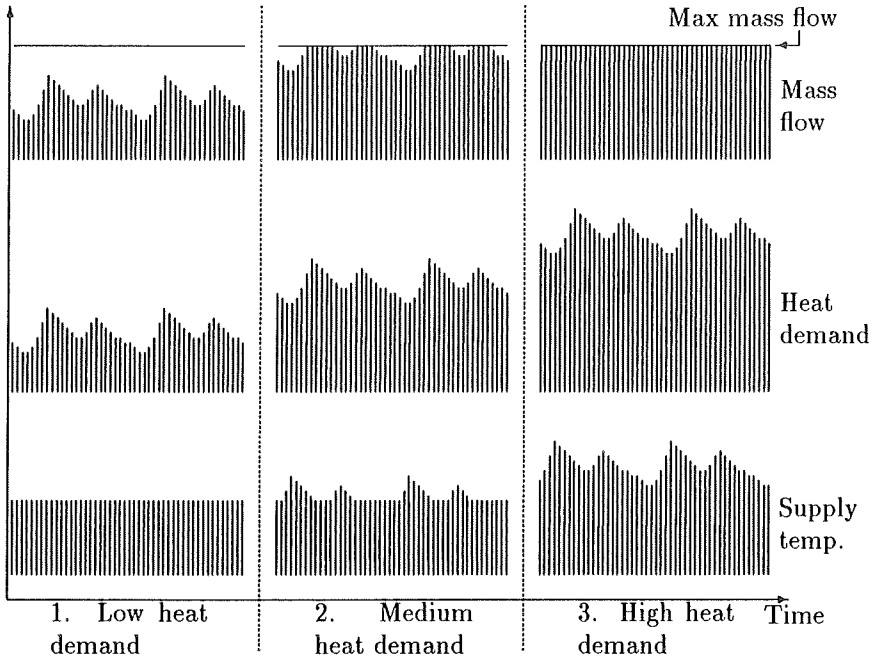


Figure 4.1: *Varying heat demand met by varying the mass flow and/or the supply temperature.*

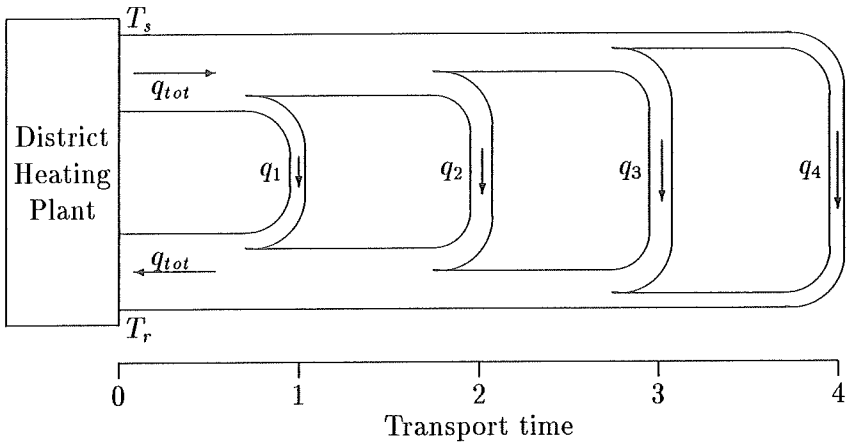


Figure 4.2: Flow to consumers put into groups with (almost) same time distance to the heat production plant.

the distribution network. Pipe dimensions and mass flow in the various sections of the distribution network determine the delay. On a discrete time scale, the total heat consumption may be divided into several subcategories. A certain amount of the heat produced at time t is consumed before time $t+1$, and another amount is consumed during the time interval $t+1$ to $t+2$ etc. (see Figure 4.2). Figure 4.3 illustrates what will happen if an impulse in the supply temperature, T_s , is produced at time $t=0$. Some of the consumers are affected by the impulse within the first time unit, i.e. during the time interval t to $t+1$. Assuming that the heat demand is not a function of the supply temperature, the mass flow, q_1 , to this group of consumers drops when the impulse arrives. Due to the dynamics of the system coming from the heat capacity of, e.g., the district heating pipes and the water itself, the influence of the temperature impulse on the mass flow at that group of consumers will decay according to certain time constants. As the impulse reaches the next groups of consumers, a

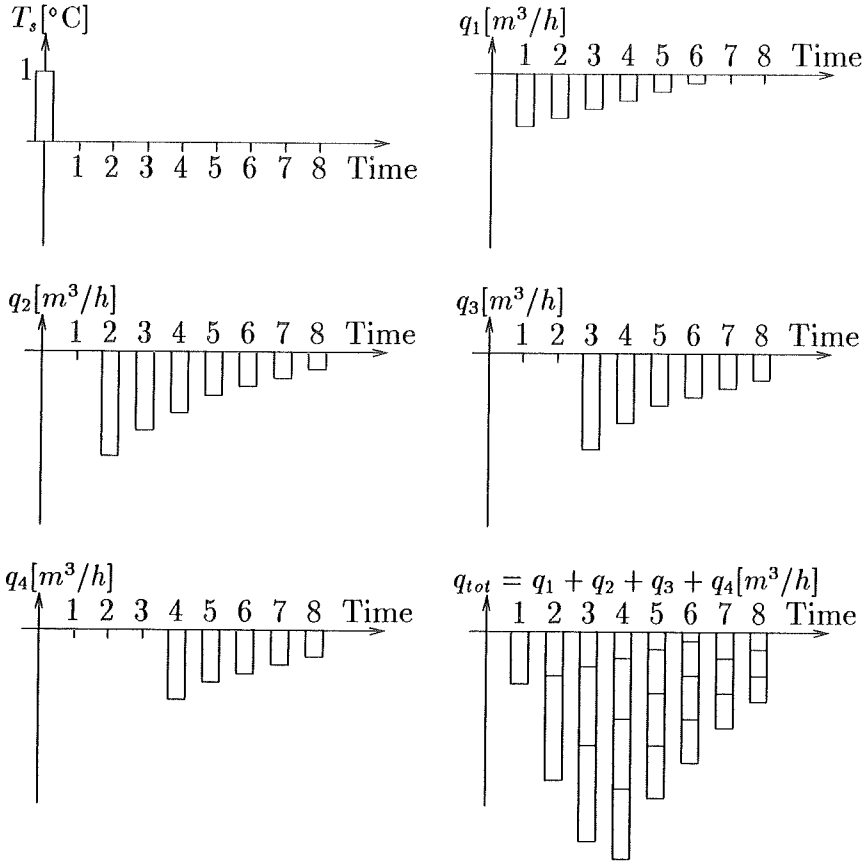


Figure 4.3: *Different transport times from the production plant to the various consumers. The effect of a supply temperature impulse is “smoothed” by the system.*

similar decrease in the mass flow is seen for these consumers. Since the total mass flow, q_{tot} , from the production plant is the sum of the mass flows to the various groups of consumers, this flow is also subject to a decrease and a subsequent increase. Due to different delays involved in this summation, the time-variation of the total mass flow is much smoother than it would have been in a hypothetical situation where all the heat consumption was concentrated at a single site.

The above discussion shows that a temporary change of the supply temperature affects the consumers after different delays. Therefore it is relevant to take heat load predictions of different horizons into consideration when controlling the supply temperature, and this, in turn, motivates the use of multi-step predictive control. Furthermore, there is a considerable diurnal and annual variation of the heat demand resulting in time-varying mass flows and time-delays. Therefore there is a need for models of the embedded parameter variations, and for multi-step predictive controllers which are able to utilize such models.

4.2 Weighted Predictive Control

In principle, a GPC assumes knowledge of the dynamic relationship between input and output of the system that is to be controlled. Expressed more directly, the GPC method assumes a model that specifies how the current control action affects the future output. In certain cases, however, such a model is not available. In this section a GPC-like controller for such insufficiently modelled systems is proposed. A simple example illustrates the idea.

4.2.1 A Simple Example

Suppose that the following model for two-step-ahead prediction of the output process $\{y(t)\}$ has been identified in the open-loop system

$$y(t) = ay(t-2) + e^*(t), \quad (4.1)$$

where $\{e^*(t)\}$ is a noise process with mean zero (possibly coloured). A two-step-ahead predictor (prediction of $y(t+2)$ given the observations $y(t), y(t-1), y(t-2), \dots$) is

$$\hat{y}^*(t+2|t) = ay(t). \quad (4.2)$$

Note that in case the noise process, $\{e^*(t)\}$, is coloured, this predictor is not optimal (the prediction error variance is not minimized).

Now assume that the real input-output system is governed by the (unknown) model

$$y(t) = ay(t-2) + bu(t-2) + e(t), \quad (4.3)$$

where $\{e(t)\}$ is a white noise process with the mean zero, and $\{u(t)\}$ is the input process also with the mean zero. The associated optimal two-step-ahead predictor (in the sense of a minimal prediction error variance) becomes

$$\hat{y}(t+2|t) = ay(t) + bu(t). \quad (4.4)$$

Since the predictor in (4.2) does not describe the relationship between the prediction and the current control action, $u(t)$, it cannot be directly applied for the design of a predictive controller. Note that (4.2) is still valid for prediction purposes although the prediction error variance is typically larger than for the input-output model in (4.3) since $e^*(t) = bu(t-2) + e(t)$.

Assume for a moment that the input-output model in (4.3) is known. Then a predictive controller can easily be obtained. If the purpose

of the controller is to minimize the variance of the control error, $y(t+2) - y^0(t+2)$, where $y^0(t+2)$ is a reference output value, the optimal control is found by solving the equation $\hat{y}(t+2|t) = y^0(t+2)$ with respect to $u(t)$ - i.e. (for $b \neq 0$):

$$ay(t) + bu(t) = y^0(t+2) \Leftrightarrow u(t) = -\alpha\delta(t+2). \quad (4.5)$$

Here $\alpha = 1/b$ is the “gain” of the controller and $\delta(t+2) = ay(t) - y^0(t+2)$ is the predicted value of the input-free control error, i.e., the predicted control error in case $u(t) = 0$. The controller in (4.5) is a special type of GPC, namely a minimal variance (MV) controller.

Now return to the assumption that the input-output model is unknown. It means that the controller gain, α , in (4.5) is unknown; but the input-free control error $\delta(t+2)$ can still be found using the predictor in (4.2) since $\hat{y}^*(t+2|t) = \hat{y}(t+2|t)|_{u(t)=0} = ay(t)$:

$$\delta(t+2) = \hat{y}^*(t+2|t) - y^0(t+2).$$

Therefore the controller in (4.5) can be expressed as

$$u(t) = -\alpha(\hat{y}^*(t+2|t) - y^0(t+2)).$$

However, the model (4.1) does not provide the gain α . Instead an estimate based on prior knowledge of the input-output system can be used. In order to find an optimal value of α it is necessary to carry out experiments with the real system letting the controller control the system with different values of α .

4.2.2 The General Case

In the general case it is assumed that models for prediction the output $N_1, N_1 + 1, \dots, N_2$ steps ahead are identified. The predictors are

assumed to be functions of past output and possibly future control actions and time:

$$\hat{y}(t+j|t) = f_j(y(t), y(t-1), \dots; u(t+j), u(t+j-1), \dots, u(t); t), \\ j = N_1, \dots, N_2.$$

The mean value of $u(t)$ is assumed to be zero. If the mean value of $u(t)$ is $\bar{u} \neq 0$ then $u(t)$ should be replaced by $u(t) - \bar{u}$. The reason why no past control actions occur as arguments to f_j is that they are implicitly represented by past output.

Suppose that the aim of the controller is to keep future output values, $y(t+N_1), \dots, y(t+N_2)$, close to the pre-specified reference values $y^0(t+N_1), \dots, y^0(t+N_2)$. If no future control actions are performed (i.e. $u(t+j) = 0, j \geq 0$), the predicted control errors become

$$\delta_j(t+j) = \\ f_j(y(t), y(t-1), \dots; 0, 0, \dots, 0; t) - y^0(t+j), \quad j = N_1, \dots, N_2.$$

The weighted predictive controller is then defined as

$$u(t) = -\alpha(t)^T \delta(t), \quad (4.6)$$

where $\alpha(t) = (\alpha_{N_1}(t), \dots, \alpha_{N_2}(t))^T$ is a vector of controller weights (gains), and $\delta(t) = (\delta_{N_1}(t+N_1), \dots, \delta_{N_2}(t+N_2))^T$ is a vector of predicted values of the input-free control errors. For time-invariant systems, a constant gain vector will be suitable.

As for the simple example above, the controller weights $\alpha(t)$ are chosen on the basis of prior knowledge about the system or by carrying out experiments. Values of N_1 and N_2 are chosen from similar considerations.

Comparing the control law in (4.6) with the real GPC laws derived in the subsequent sections, it is seen that they have the following in

common: the current control is a linear function of predicted values of the input-free control errors. (The term “input-free” covers cases like $u(t+j) = 0$ or $u(t+j) - u(t+j-1) = 0$ ($j \geq 0$) depending on the type of GPC controller.) Since the weighted predictive controller (4.6) is based upon general linear feedback from input-free control errors it is a general dynamic controller which covers both minimal variance control, linear quadratic control and generalized predictive control.

As the weighted predictive controller is based upon a kind of a “pseudo” GPC strategy, it inherits the robustness properties of the GPC methods (Sections 4.3 and 4.4) – provided that the controller gain, $\alpha(t)$, and the minimum and the maximum output horizons, N_1 and N_2 , are chosen appropriately.

4.2.3 Weighted Predictive Control of Supply Temperature

This section describes an example of application of weighted predictive control for a practical case, namely for the control of the supply temperature in a district heating system. The models and control method presented in the following are used at the combined heat and power plant Vestkraft in Esbjerg.

The control is based on models for prediction of heat load described by Madsen *et al.* (1990) and Sejling (1993). These models give rise to the following one-step-ahead predictor

$$\begin{aligned} \hat{p}(t+1|t) = & \alpha_1 p(t) + \alpha_2 p(t-1) + \alpha_{24} p(t-23) + \alpha_{25} p(t-24) + \alpha_{26} p(t-25) \\ & + \beta_{1,0} \nabla t_s(t+1) + \beta_{1,1} \nabla t_s(t) \\ & + \beta_{2,1} t_a(t) + \beta_{2,2} t_a(t-1) \end{aligned}$$

$$\begin{aligned}
& +\beta_{3,1}w(t) + \beta_{3,2}w(t-1) \\
& +\mu_1(t+1)I_{\{\text{workday}\}}(t+1) + \mu_2(t+1)I_{\{\text{weekend}\}}(t+1) + l,
\end{aligned}$$

where

- t is the time [h],
- $\hat{p}(t+1|t)$ is the one-step-ahead prediction of the heat load at time $t+1$ [J/s],
- $p(t)$ is the heat load [J/s],
- $t_s(t)$ is the supply temperature ($\nabla t_s(t) = t_s(t) - t_s(t-1)$) [$^{\circ}\text{C}$],
- $t_a(t)$ is the ambient air temperature [$^{\circ}\text{C}$],
- $w(t)$ is the wind speed [m/s],
- $\mu_1(t)$ is a diurnal heat load profile for workdays [J/s],
- $\mu_2(t)$ is a diurnal heat load profile for weekends [J/s],
- $I_{\{\text{workday}\}}(t)$ is an indicator which is 1 on workdays and 0 otherwise,
- $I_{\{\text{weekend}\}}(t)$ is equal to $1 - I_{\{\text{workday}\}}(t)$, and
- l is an adjusting parameter accounting for the variables having mean values different from zero [J/s].

Note that all variables are recorded as hourly averages and not as instantaneous values. That is why the prediction at time $t+1$ is a function of the supply temperature up to time $t+1$ and not only up to time t . The corresponding predictors for prediction horizons, j , between 2 and 22 are

$$\begin{aligned}
\hat{p}(t+j|t) = & \alpha_j p(t) + \alpha_{j+1} p(t-1) \\
& +\alpha_{24} p(t+j-24) + \alpha_{25} p(t+j-25) + \alpha_{26} p(t+j-26) \\
& +\beta_{1,0} \nabla t_s(t+j) + \beta_{1,1} \nabla t_s(t+j-1) \\
& +\beta_{1,2} \nabla t_s(t+j-2) + \beta_{1,3} \nabla t_s(t+j-3) \\
& +\beta_{2,j} t_a(t) + \beta_{2,j+1} t_a(t-1) + \beta_{2,j+2} t_a(t-2) + \beta_{2,j+3} t_a(t-3) \\
& +\beta_{3,j} w(t) + \beta_{3,j+1} w(t-1) + \beta_{3,j+2} w(t-2) + \beta_{3,j+3} w(t-3) \\
& +\mu_1(t+j)I_{\{\text{workday}\}}(t+j) + \mu_2(t+j)I_{\{\text{weekend}\}}(t+j) + l.
\end{aligned} \tag{4.7}$$

Some of the parameters in the model depend on j while others don't. Consider for instance the term $\alpha_{j+1}p(t-1)$. The subscript in α_{j+1} indicates that the time lag between the regressor $p(t-1)$ and the prediction $\hat{p}(t+j|t)$ is $j+1$. The same rule applies to the other parameter subscripts, and therefore only some of them depend on j . Actually, all of the parameter subscripts should have been provided with an extra j in order to indicate that a separate model with its own set of parameters is set up for each prediction horizon. However, in accordance with the notation used by Sejling (1993) these extra j 's are omitted.

The predictors are used as a basis for weighted predictive control of the supply temperature. Since the purpose of the controller is to perform an on-line control of the supply temperature in order to keep the mass flow, $q(t)$ [kg/s], close to but below a critical maximum value, q_{\max} , it would be more appropriate to have prediction models of the mass flow than models of the heat load. It will, however, be shown that load models can be used instead of flow models. In order to do that the following relationship between mass flow and heat load is required:

$$p(t) = c_w q(t)(t_s(t) - t_r(t)) , \quad (4.8)$$

where c_w [J/(kg°C)] is the heat capacity of water and $t_r(t)$ [°C] is the return temperature. Suppose that the future values of the mass flow at time $t+j$ should be kept below q_{\max} with probability π . If the actual time is t this means that

$$P\{q(t+j) \leq q_{\max}\} = \pi , \quad 0 < N_1 \leq j \leq N_2 , \quad (4.9)$$

where N_1 and N_2 are chosen so that $q(t+N_1), \dots, q(t+N_2)$ encompass future mass flow values significantly affected by the next control, $t_s(t+1)$. Actually, these equations give us the reference values, $q^0(t+j)$, of the mass flow implicitly. This will be shown later on. If (4.8) is used to substitute for $q(t+j)$ in (4.9) we obtain

$$P\{p(t+j) \leq c_w q_{\max}(t_s(t+j) - t_r(t+j))\} = \pi . \quad (4.10)$$

Since the return temperature j steps ahead is unknown, a predictor has to be introduced. As the variation of the return temperature is very slow and with a variance which is much less than the variance of the supply temperature it is sufficient in the present context to use the following random walk model

$$t_r(t) = t_r(t-1) + e(t),$$

where $\{e(t)\}$ is assumed to be white noise with mean zero. The corresponding j -step-ahead predictor is

$$\hat{t}_r(t+j|t) = t_r(t), \quad j = 1, 2, 3, \dots \quad (4.11)$$

The heat load and the return temperature j steps ahead can be written as a sum of the prediction and the prediction error:

$$\begin{aligned} p(t+j) &= \hat{p}(t+j|t) + \varepsilon_p(t+j|t) \\ t_r(t+j) &= \hat{t}_r(t+j|t) + \varepsilon_{t_r}(t+j|t), \end{aligned} \quad (4.12)$$

where the prediction errors, $\varepsilon_p(t+j|t)$ and $\varepsilon_{t_r}(t+j|t)$, are assumed to be mutually independent and to have mean zero and variances $\sigma_p^2(j)$ and $\sigma_{t_r}^2(j)$. Inserting (4.12) into (4.10) gives (after a few manipulations)

$$\begin{aligned} &P\{\varepsilon_p(t+j|t) + c_w q_{\max} \varepsilon_{t_r}(t+j|t) \\ &\leq c_w q_{\max} (t_s(t+j) - \hat{t}_r(t+j|t)) - \hat{p}(t+j|t)\} = \pi. \end{aligned}$$

If assuming that the prediction errors are normally distributed this can be rewritten as

$$\frac{c_w q_{\max} (t_s(t+j) - \hat{t}_r(t+j|t)) - \hat{p}(t+j|t)}{\sqrt{\sigma_p^2(j) + c_w^2 q_{\max}^2 \sigma_{t_r}^2(j)}} = u_\pi, \quad (4.13)$$

where u_π is the 100π % quantile in the standardized normal distribution. The next step is to substitute for $\hat{t}_r(t+j|t)$ and $\hat{p}(t+j|t)$ in

this equation. But it is necessary to rewrite (4.7) first so it fits the control strategy.

In the district heating system in Esbjerg/Varde it has been estimated that the heat produced at Vestkraft reaches the first consumers after approximately 4 hours. Therefore the control is based on predictions of the heat load 4, 5 and 6 steps ahead. This means that $j = 4, 5, 6$ are the relevant prediction horizons in (4.7), and that $N_1 = 4$ and $N_2 = 6$.

Now introduce a control horizon $N_u = 1$ (see Section 4.3). This means that the choice of the next control, $t_s(t+1)$, is subject to

$$t_s(t+j) = t_s(t+1), \quad j = 2, \dots, 6 \quad (4.14)$$

($\Rightarrow \nabla t_s(t+j) = 0$ for $j = 2, \dots, 6$). Under this restriction, (4.7) can be expressed as

$$\hat{p}(t+j|t) = \beta_j \nabla t_s(t+1) + \gamma_j(t) = \beta_j(t_s(t+1) - t_s(t)) + \gamma_j(t), \quad (4.15)$$

where

$$\beta_j = \begin{cases} \beta_{1,3} & \text{if } j = 4 \\ 0 & \text{if } j = 5, 6 \end{cases}$$

$$\gamma_j(t) = \alpha_j p(t) + \alpha_{j+1} p(t-1) \\ + \alpha_{24} p(t+j-24) + \alpha_{25} p(t+j-25) + \alpha_{26} p(t+j-26) \\ + \beta_{2,j} t_a(t) + \beta_{2,j+1} t_a(t-1) + \beta_{2,j+2} t_a(t-2) + \beta_{2,j+3} t_a(t-3) \\ + \beta_{3,j} w(t) + \beta_{3,j+1} w(t-1) + \beta_{3,j+2} w(t-2) + \beta_{3,j+3} w(t-3) \\ + \mu_1(t+j) I_{\{\text{workday}\}}(t+j) + \mu_2(t+j) I_{\{\text{weekend}\}}(t+j) + l.$$

Inserting (4.11), (4.14) and (4.15) into (4.13) it is readily seen that

$$\frac{c_w q_{\max}(t_s(t+1) - t_r(t)) - \beta_j(t_s(t+1) - t_s(t)) - \gamma_j(t)}{\sqrt{\sigma_p^2(j) + c_w^2 q_{\max}^2 \sigma_{t_r}^2(j)}} = u_\pi,$$

or if the equation is solved with respect to $t_s(t+1)$,

$$t_s(t+1) = \frac{S_j + K t_r(t) - \beta_j t_s(t) + \gamma_j(t)}{K - \beta_j}, \quad (4.16)$$

where

$$\begin{aligned} K &= c_w q_{\max} \\ S_j &= u_\pi \sqrt{\sigma_p^2(j) + K^2 \sigma_{t_r}^2(j)}. \end{aligned}$$

For each of the prediction horizons (4, 5 and 6 hours), (4.16) result in a value, $t_{s,j}(t+1)$, of $t_s(t+1)$, which is a solution to (4.9), but the final controller should, of course, provide us with a unique value of $t_s(t+1)$. Such a value is obtained by constructing a weighted average of $t_{s,4}(t+1)$, $t_{s,5}(t+1)$ and $t_{s,6}(t+1)$:

$$t_s(t+1) = \chi^T \mathbf{t}_s(t+1), \quad (4.17)$$

where

$$\chi = (\chi_4 \ \chi_5 \ \chi_6)^T, \quad \sum_{j=4}^6 \chi_j = 1,$$

and

$$\begin{aligned} \mathbf{t}_s(t+1) &= (t_{s,4}(t+1) \ t_{s,5}(t+1) \ t_{s,6}(t+1))^T \\ &= \begin{pmatrix} \frac{S_4 + K t_r(t) - \beta_{1,3} t_s(t) + \gamma_4(t)}{K - \beta_{1,3}} \\ \frac{S_5 + K t_r(t) + \gamma_5(t)}{K} \\ \frac{S_6 + K t_r(t) + \gamma_6(t)}{K} \end{pmatrix}. \end{aligned}$$

Equation (4.17) constitute the controller in its operational form. However, it is not directly comparable with the controller in (4.6) where the control action is a linear feedback from the predicted values of input-free control errors. In order to obtain this form of the

controller, the j -step-ahead predictors are expressed in terms of the mass flow and the differenced supply temperature. To match the symbols used previously put

$$\begin{aligned} y(t) &= q(t) \\ u(t) &= t_s(t) - t_r(t-1). \end{aligned} \quad (4.18)$$

By replacing the heat load, the mass flow and the return temperature in (4.8) by the corresponding j -step-ahead predictions, the following approximation is obtained:

$$\begin{aligned} \hat{y}(t+j|t) = \hat{q}(t+j|t) &\approx \frac{\hat{p}(t+j|t)}{c_w(t_s(t+j) - \hat{t}_r(t+j|t))} \\ &= \frac{\beta_j u(t+1) + \gamma_j(t)}{c_w(u(t+1) + t_\Delta(t))}, \end{aligned} \quad (4.19)$$

where (4.11), (4.14) and (4.15) have been used, and where $t_\Delta(t) = t_s(t) - t_r(t)$. This linearization is reasonable if $\sigma_p(j)$ and $\sigma_{t_r}(j)$ are sufficiently small. Equation (4.16), which implicitly determines the reference value of the flow at time $t+j$, is easily expressed in terms of $u(t+1)$ instead of $t_s(t+1)$ by using (4.18):

$$u(t+1) = \frac{S_j - K t_\Delta(t) + \gamma_j(t)}{K - \beta_j}.$$

By inserting this into (4.19), the future reference value of the flow is found:

$$y^0(t+j) = \frac{\beta_j(S_j - K t_\Delta(t)) + K \gamma_j(t)}{c_w(S_j + \gamma_j(t) - \beta_j t_\Delta(t))}.$$

It is now possible to find an expression of $\delta_j(t+j)$ (the predicted value of the input-free control error):

$$\begin{aligned} \delta_j(t+j) &= \hat{y}(t+j|t)|_{u(t+1)=u(t+2)=\dots=u(t+6)=0} - y^0(t+j) \\ &= \frac{\gamma_j(t)}{c_w t_\Delta(t)} - y^0(t+j) \\ &= \frac{[\gamma_j(t) - \beta_j t_\Delta(t)][S_j - K t_\Delta(t) + \gamma_j(t)]}{c_w t_\Delta(t)[S_j + \gamma_j(t) - \beta_j t_\Delta(t)]}, \end{aligned}$$

where $1 \leq j \leq 6$. The next step is to find an expression of $\alpha(t)$ so that the controller

$$u(t+1) = -\alpha(t)^T \delta(t), \quad (4.20)$$

of which $\delta(t)$ is defined as in Section 4.2.2, is equivalent to the controller in (4.17). Comparison leads to

$$\alpha(t) = \begin{pmatrix} -\chi_4 \frac{c_w t_{\Delta}(t)[S_4 + \gamma_4(t) - \beta_{1,3} t_{\Delta}(t)]}{[K - \beta_{1,3}][\gamma_4(t) - \beta_{1,3} t_{\Delta}(t)]} \\ -\chi_5 \frac{c_w t_{\Delta}(t)[S_5 + \gamma_5(t)]}{K \gamma_5(t)} \\ -\chi_6 \frac{c_w t_{\Delta}(t)[S_6 + \gamma_6(t)]}{K \gamma_6(t)} \end{pmatrix}. \quad (4.21)$$

If (4.20) is compared with (4.6), it becomes apparent that the only difference between the two controllers is that the former gives a value of $u(t+1)$ while the latter gives a value of $u(t)$. This difference is due to different sampling procedures. Equation (4.6) is based upon instantaneous values while (4.20) is based on hourly averages recorded at the end of the hourly intervals. This actually means that both controllers give the value of the input within the time interval t to $t+1$.

Previously it was argued that the gain vector, $\alpha(t)$, of a controller for a time-invariant system should be constant. However, the model in (4.7) is in fact time-invariant but the gain vector in (4.21) depends on time. This seeming contradiction is due to the fact that the weighted predictive controller (4.20) corresponds to the prediction model for the flow in (4.19) and not to the model in (4.7).

4.2.4 Results Obtained at Vestkraft in Esbjerg

As previously mentioned, the weighted predictive control method described in Section 4.2.3 has been implemented at Vestkraft in Esbjerg. In the present section, a few results obtained in connection with

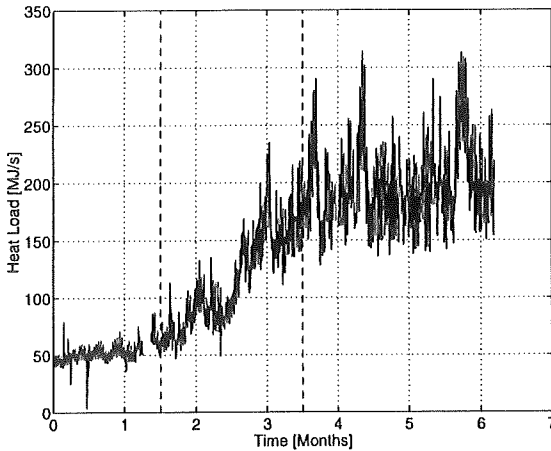


Figure 4.4: Heat load measurements (hourly averages). Period: 1 August 1991 to 5 February 1992.

this implementation are presented. The control parameters for (4.17) have been chosen as follows:

π [%]	q_{\max} [tons/h]	χ_4	χ_5	χ_6
99	4600	0.2	0.4	0.4

In order to obtain a reasonably smooth control signal (supply temperature), the χ -values have been chosen so that they are not too different. On the other hand, the average transport time from Vestkraft to the consumers is probably at least 6 hours. Therefore χ_5 and χ_6 are larger than χ_4 .

Figure 4.4 shows the heat load measurements recorded in a period from August 1991 to February 1992. The figure represents a period with summer load (from August to the middle of September, $0 \leq t \leq 1.5$), load in a transition period (from the middle of September

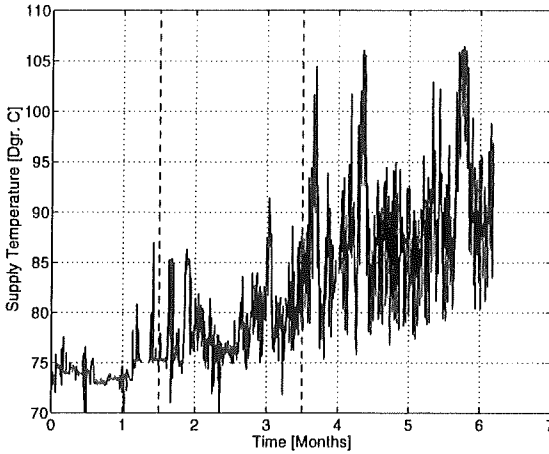


Figure 4.5: *Supply temperature (hourly averages). Period: 1 August 1991 to 5 February 1992.*

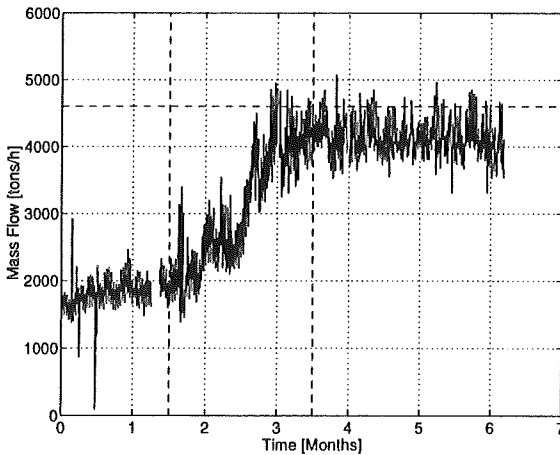


Figure 4.6: *Mass flow (hourly averages). Period: 1 August 1991 to 5 February 1992.*

to the middle of November, $1.5 \leq t \leq 3.5$) and a period with winter load (from the middle of November to February, $3.5 \leq t \leq 6.2$). The diurnal amplitude and the level of the load are clearly increasing during the period. This tendency is, of course, due to a greater part of the load being used for heating the buildings in the fall and winter seasons.

Figure 4.5 shows the supply temperature. Until the transition period begins ($t \approx 1.5$), the supply temperature is almost constant about 75°C . During the transition period, the level and the amplitude of the signal increase, but this is not due to the weighted predictive controller considered in this section. The increase is caused by other controllers of the implementation at Vestkraft which ensure that the water delivered to the consumers has a certain minimum temperature, which depends on the ambient air temperature (hot tap water restriction). At the end of the transition period ($t \approx 3$), the weighted predictive controller takes over. The reason is that the mass flow approaches its maximum value (4600 tons/h, see Figure 4.6). After this period the patterns of the heat load variation and the supply temperature variation are very much alike. The mass flow, however, does not show very large variations in this period. Thus the controller seems to perform what is was intended for: the level of the mass flow is kept close to its maximum value with small fluctuations only, while the supply temperature is subject to considerable variation in order to ensure that the varying heat demand is met. Analysis of data from this and other periods show that the weighted predictive controller is able to keep the average supply temperature in the winter season about 9°C lower than before the controller was brought into operation (see Appendix B).

4.3 Ordinary Generalized Predictive Control

The results presented in this section are well-known from the literature, and they are included because they are closely related to the new results presented in other parts of this chapter.

The generalized predictive controller was originally presented by Clarke *et al.* (1987A) (the main reference in this section), and the results was reproduced by Bitmead *et al.* (1990). Clarke *et al.* (1987A) show that GPC can be viewed as a generalization of well-known controllers (e.g., the ordinary minimal variance controller (MV) and the generalized minimal variance controller (GMV)). By comparative simulations they verify that the GPC is superior to widely accepted controllers (the PID, the GMV and algorithms based on pole-placement). The GPC is based on five fundamental ideas:

1. A special ARIMAX model¹ (rather than an ARMAX model) is adopted to describe the dynamic system considered:

$$A(q^{-1})y(t) = B(q^{-1})u(t-1) + \frac{C(q^{-1})}{\nabla}e(t), \quad (4.22)$$

where $A(q^{-1})$, $B(q^{-1})$ and $C(q^{-1})$ are polynomials in the backward shift operator q^{-1} :

$$\begin{aligned} A(q^{-1}) &= 1 + a_1q^{-1} + \cdots + a_nq^{-n} \\ B(q^{-1}) &= b_0 + b_1q^{-1} + \cdots + b_mq^{-m} \\ C(q^{-1}) &= 1 + c_1q^{-1} + \cdots + c_rq^{-r}, \end{aligned}$$

and ∇ is the differencing operator, $1 - q^{-1}$. If the system has a non-zero time-delay, k ($< m$), from input, $u(t)$, to output,

¹Clarke *et al.* (1987A) use the designation "CARIMA model" (Controlled Auto-Regressive and Integrated Moving-Average model).

$y(t)$, the leading coefficients, b_0, b_1, \dots, b_{k-1} , of the polynomial $B(q^{-1})$ are zero. $\{e(t)\}$ is a white noise process with the mean zero and the variance σ_e^2 . Clarke *et al.* (1987A) argue that the integrated noise term, $\frac{C(q^{-1})}{\nabla}e(t)$, in (4.22) is appropriate for many practical applications since the disturbances very often belong to a non-stationary stochastic process. Moreover they utilize that the integrator makes the derivation of the controller more straight-forward. But they do not discuss the consequences of including an integrating factor in the model for a system governed by an ordinary ARMAX model. Furthermore, the simulation studies given by Clarke *et al.* (1987A) are carried out without including any noise. The fact is that it may lead to identification of a “false” MA-parameter and very poor estimates of the other parameters if an ARIMAX model is assumed in the case of an ARMAX system.

2. The GPC takes the future sequence of output predictions over a finite horizon into account when deciding how to choose the current control value. The controller tries to minimize the control effort and the deviation between future output values and a predetermined reference signal at the same time. The strategy leads to a very robust controller with favourable stability properties.
3. The multi-step output predictors have been derived by means of a recursive version of the Diophantine equation for the relevant prediction horizons. The recursivity reduces the computational burden of the prediction evaluations. For models having embedded parameter variation, however, Diophantine equations *cannot* be employed. In this case, the multi-step predictors have to be derived within a more general framework.
4. Control increments rather than the control values themselves are weighted in the cost function. This is done to reduce the variation of the control signal.

5. A control horizon is chosen after which the control increments are fixed at zero. Future control actions before this horizon is reached are determined by a minimization of the cost function.

Clarke *et al.* (1987A) and Clarke *et al.* (1987B) restrict themselves to deal with ARIMAX models and Diophantine equations. In this presentation, however, the basic idea of multi-step predictive control is extended in order to develop controllers for embedded model structures, which also include time-varying ARMAX models and which do not require the first order integration introduced in (4.22). This is done in Section 4.4. In the present section, the ordinary GPC is described. The results will be used for comparison with the embedded model based GPC in Section 4.4.

First some attractive qualities of the GPC should be mentioned (Clarke *et al.* (1987A)):

1. Unlike MV controllers, the GPC is able to handle non-minimum phase systems.
2. Open-loop unstable systems and systems with badly damped poles are stabilized.
3. Systems with slowly time-varying or unknown time-delay are controlled without problems. The MV and the GMV controllers function rather badly in case the time-delay of the model does not coincide with the time-delay of the system.
4. Unlike most other controllers, the GPC is robust to an overestimated model order.

4.3.1 The Model and the Diophantine Equation

The GPC is basically designed to handle models with the C -polynomial set to 1 – i.e. an ARIX model:

$$A(q^{-1})y(t) = B(q^{-1})u(t-1) + \frac{1}{\nabla}e(t). \quad (4.23)$$

Since the derivation of the controller corresponding to this model illustrates the fundamental technique, the general model in (4.22) is not further considered.

To construct the j -step-ahead predictor of $y(t+j)$, the Diophantine identity,

$$1 = E_j(q^{-1})A(q^{-1})\nabla + q^{-j}F_j(q^{-1}), \quad (4.24)$$

is considered (simply obtained by polynomial division). The polynomials E_j and F_j , which are uniquely defined given A and the prediction horizon j , are of degrees $j-1$ and n , respectively. Multiplying (4.23) by $E_j\nabla q^j$ gives

$$E_j A \nabla y(t+j) = E_j B \nabla u(t+j-1) + E_j e(t+j).$$

In this equation the expression of $E_j A \nabla$ which can be obtained from the Diophantine equation (4.24) is inserted:

$$(1 - q^{-j}F_j)y(t+j) = E_j B \nabla u(t+j-1) + E_j e(t+j)$$

or

$$y(t+j) = E_j B \nabla u(t+j-1) + F_j y(t) + E_j e(t+j).$$

The optimal j -step-ahead predictor is found as the conditional expectation of $y(t+j)$ given output data up to time t . The predictor becomes a function of the input data up to time $t+j-1$:

$$\hat{y}(t+j|t) = G_j \nabla u(t+j-1) + F_j y(t), \quad (4.25)$$

where $G_j = E_j B$ (degree $j + m - 1$). The fact that the term $E_j e(t + j)$ is a sum of future noise components with the conditional expectation 0 has been used to obtain (4.25). It can be verified that the first j coefficients of the polynomial G_j are equal to the first j weights of the step response of the system.

The E_j and F_j polynomials thus play an important role in the derivation of the j -step-ahead predictor since they are used to obtain a separation into what is known and what is unknown at time t . Furthermore, the GPC is based on a set of predictors with the prediction horizon running from a minimum to a maximum value. Therefore, instead of solving the Diophantine equation a number of times for different values of j , the coefficients of E_j and F_j can conveniently be computed by recursion. If E_j and F_j ($j = 1, 2, 3, \dots$) are known, then

$$\begin{aligned} E_{j+1}(q^{-1}) &= E_j(q^{-1}) + q^{-j} f_{j0} \\ F_{j+1}(q^{-1}) &= q[F_j(q^{-1}) - f_{j0}A(q^{-1})\nabla] , \end{aligned} \quad (4.26)$$

where f_{j0} is the constant term of the F_j -polynomial. The recursions are initialized as

$$E_1(q^{-1}) = 1 , \quad F_1(q^{-1}) = q(1 - A\nabla) . \quad (4.27)$$

The proof of (4.27) is trivial (insert into (4.24)). (4.26) can be proved by an inductive proof:

Assume that E_j and F_j ($j = 1, 2, 3, \dots$) satisfy the Diophantine equation:

$$1 = E_j A\nabla + q^{-j} F_j .$$

Then it is shown that also E_{j+1} and F_{j+1} , given by (4.26), satisfy the Diophantine equation:

$$\begin{aligned} E_{j+1}A\nabla + q^{-(j+1)}F_{j+1} \\ = (E_j + q^{-j}f_{j0})A\nabla + q^{-(j+1)}q(F_j - f_{j0}A\nabla) \end{aligned}$$

$$\begin{aligned}
 &= E_j A \nabla + q^{-j} F_j \\
 &= 1
 \end{aligned}$$

□

If the model parameters are estimated adaptively, the E_j - and F_j -polynomials have to be recalculated at each sampling instant. In this self-tuning case the recursive solution of the Diophantine equation reduces the computational task significantly.

4.3.2 The GPC Cost Function

In practical cases the job of the controller is to keep the system output “close to” some predetermined reference signal, $\{y^0(t)\}$ which can be constant (set-point) or time-varying. The control error, i.e. the difference between the output and its reference, is limited by weighting it into the cost function defining the GPC. To prevent that a reduction of the control error variance imposes undesirably large variations of the control signal, the control increments ($\nabla u(t)$) are also weighted. At time t these requirements are met by the following optimization

$$\begin{aligned}
 &\min_{\nabla u(t+j-1)_{j=1,\dots,N_2}} J[N_1, N_2, \lambda_1, \dots, \lambda_{N_2}; t, \nabla u(t+j-1)_{j=1,\dots,N_2}] \\
 &= E_t \left[\sum_{j=N_1}^{N_2} [y(t+j) - y^0(t+j)]^2 + \sum_{j=1}^{N_2} \lambda_j [\nabla u(t+j-1)]^2 \right]
 \end{aligned} \tag{4.28}$$

subject to the model in (4.23) and

$$\nabla u(t+j-1) = 0, \quad j > N_u, \tag{4.29}$$

where

N_1 is the minimum cost (output) horizon,

N_2 is the maximum cost (output) horizon,
 N_u ($< N_2$) is the control horizon,
 λ_j (≥ 0) is a control-weighting sequence,

and $E_t[\cdot]$ denotes the conditional expectation of its argument conditioned on data up to time t . By eliminating $\nabla u(t + N_u), \dots, \nabla u(t + N_2 - 1)$ from (4.28) by using (4.29), the following optimization problem is obtained:

$$\begin{aligned} & \min_{\nabla u(t+j-1)_{j=1, \dots, N_u}} J[N_1, N_2, N_u, \lambda_1, \dots, \lambda_{N_u}; t, \nabla u(t+j-1)_{j=1, \dots, N_u}] \\ & = E_t \left[\sum_{j=N_1}^{N_2} [y(t+j) - y^0(t+j)]^2 + \sum_{j=1}^{N_u} \lambda_j [\nabla u(t+j-1)]^2 \right]. \end{aligned} \quad (4.30)$$

subject to the model in (4.23) and the constraint in (4.29).

The differenced control values $\nabla u(t), \nabla u(t+1), \dots, \nabla u(t+N_u-1)$ are “free” to be optimized, and the values $\nabla u(t+N_u), \dots, \nabla u(t+N_2-1)$ are fixed at zero (equivalent to an infinite control-weight). The optimization in (4.30) results in both the present and future control values, but only the present one is implemented. At the next sampling instant, the optimization is repeated. Figure 4.7 illustrates the GPC criterion.

The output and control horizons (N_1, N_2 and N_u) together with the control weights (λ_j) constitute the design parameters of the GPC controller. Clarke *et al.* (1987A) give a few guidelines for the choice of these parameters:

N_1 : *The minimum output horizon.* If the time-delay, k , from input to output is known exactly, N_1 is appropriately set to k . Choosing an N_1 which is lower than k implies superfluous calculations since the current control action, $u(t)$, does not affect output until time $t+k$. If k is unknown or time-varying, N_1 can be set on 1 or on a lower limit of k , if it is known. In this case the degree of the

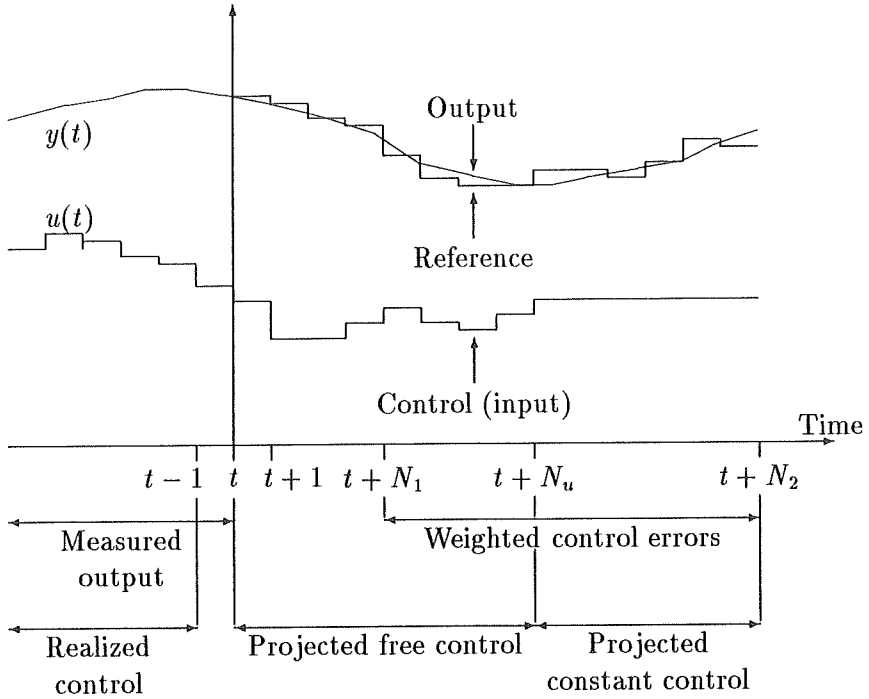


Figure 4.7: *The various signals in the GPC.*

estimated $B(q^{-1})$ -polynomial should be increased to encompass all possible values of k . Observe that N_1 should never be larger than the minimum time-delay.

N_2 : *The maximum output horizon.* For open-loop stable systems it is recommended to choose a value of N_2 close to the risetime of the system. The choice should ensure that all output values significantly affected by the current control are included in the cost function. For non-minimum phase systems this means that N_2 is likely to exceed the degree of $B(q^{-1})$. It is a fact, however, that the calculations associated with the optimization problem become more extensive with an increasing N_2 .

N_u : *The control horizon.* A proper value of this design parameter depends on the complexity (order) of the system. For simple systems $N_u = 1$ gives a reasonably good control. Increasing N_u leads to more active control and smaller control errors. As N_u is increased the changes in the optimal value of the current control action, $u(t)$, will vanish. As shown in Section 4.3.3 the solution of the optimization problem involves inversion of a $N_u \times N_u$ matrix. Therefore, it is of interest to keep N_u at a acceptable minimum where both the computational effort and the qualities associated with the resulting controller are considered. In case $N_u = 1$ is not sufficient for a particular system, a rule of thumb suggests that N_u be set equal to the number of poles near the stability boundary.

λ_j : *The sequence of control weights.* The control weights may, e.g. be set to zero. The restrictions in (4.29) implies that this choice works even if the system is non-minimum phase. This is unlike the GMV controller which requires a positive control weight in connection with non-minimum phase systems. If the control signal should be smoothed further, the λ -weights can be increased from zero.

4.3.3 Optimization of the GPC Cost Function

In this section the optimization problem defined in (4.30) is solved. The computations are made more compact by introducing matrix notation. At first it is shown how to bring all the j -step-ahead predictors together in one linear matrix equation.

As appear from (4.25), the j -step-ahead predictors are linear functions of future control actions (including $u(t)$). Hence a j -step predictor can be written as a sum of two terms: A variable term that depends linearly on the next j control actions ($u(t), \dots, u(t+j-1)$) and a constant term that depends entirely on past control actions and output measurements ($u(t-1), u(t-2), \dots$ and $y(t), y(t-1), \dots$). In order to obtain this separation in future and past, Equation (4.25) is rewritten as

$$\hat{y}(t+j|t) = G'_j(q^{-1})\nabla u(t+j-1) + z_j(t), \quad (4.31)$$

where $G'_j(q^{-1})$ and $z_j(t)$ are defined by

$$\begin{aligned} G'_j(q^{-1}) &= g_{j,0} + g_{j,1}q^{-1} + \dots + g_{j,j-1}q^{-(j-1)} \\ z_j(t) &= G''_j(q^{-1})\nabla u(t+j-1) + F_j(q^{-1})y(t), \end{aligned} \quad (4.32)$$

and $G''_j(q^{-1})$ is defined by

$$\begin{aligned} G'_j(q^{-1}) + G''_j(q^{-1}) &= G(q^{-1}) \\ &= g_{j,0} + g_{j,1}q^{-1} + \dots + g_{j,j+m-1}q^{-(j+m-1)}. \end{aligned}$$

Note that the polynomial G_j appearing in (4.25) has simply been divided into a sum of two terms, G'_j and G''_j , which are used as filters for future and past control actions, respectively. As mentioned earlier, the first j coefficients of G_j (i.e. the coefficients of G'_j) are equal to the first j weights of the step response. So if the step response weight

at lag i is denoted g_i , it is found that $g_{ji} = g_i$ for $i = 0, 1, \dots, j - 1$ independent of the prediction horizon and (cf. (4.32))

$$G'_j(q^{-1}) = g_0 + g_1q^{-1} + \dots + g_{j-1}q^{-(j-1)} .$$

Defining

$$\hat{\mathbf{y}}(t) = (\hat{y}(t + N_1|t), \hat{y}(t + N_1 + 1|t), \dots, \hat{y}(t + N_2|t))^T$$

$$\tilde{\mathbf{u}}(t) = (\nabla u(t), \nabla u(t + 1), \dots, \nabla u(t + N_u - 1))^T$$

$$\mathbf{z}(t) = (z_{N_1}(t), z_{N_1+1}(t), \dots, z_{N_2}(t))^T$$

$$\mathbf{G} =$$

$$\begin{pmatrix} g_{N_1-1} & g_{N_1-2} & \cdots & g_0 & 0 & \cdots & 0 & 0 \\ g_{N_1} & g_{N_1-1} & \cdots & g_1 & g_0 & \cdots & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots & \ddots & \vdots & \vdots \\ g_{N_u-2} & g_{N_u-3} & \cdots & g_{N_u-N_1-1} & g_{N_u-N_1-2} & \cdots & g_0 & 0 \\ g_{N_u-1} & g_{N_u-2} & \cdots & g_{N_u-N_1} & g_{N_u-N_1-1} & \cdots & g_1 & g_0 \\ \vdots & \vdots & & \vdots & \vdots & & \vdots & \vdots \\ g_{N_2-2} & g_{N_2-3} & \cdots & g_{N_2-N_1-1} & g_{N_2-N_1-2} & \cdots & g_{N_2-N_u} & g_{N_2-N_u-1} \\ g_{N_2-1} & g_{N_2-2} & \cdots & g_{N_2-N_1} & g_{N_2-N_1-1} & \cdots & g_{N_2-N_u+1} & g_{N_2-N_u} \end{pmatrix}$$

and using the restrictions defined in (4.29), the predictors in (4.31) with j running from N_1 to N_2 can be written in one linear matrix equation:

$$\hat{\mathbf{y}}(t) = \mathbf{G}\tilde{\mathbf{u}}(t) + \mathbf{z}(t) . \tag{4.33}$$

Note that \mathbf{G} is formed by extracting rows N_1 to N_2 and columns 1 to N_u of the following lower-triangular matrix of dimension $N_2 \times N_2$:

$$\begin{pmatrix} g_0 & 0 & \cdots & 0 & 0 \\ g_1 & g_0 & & 0 & 0 \\ \vdots & \vdots & \ddots & & \vdots \\ g_{N_2-2} & g_{N_2-3} & \cdots & g_0 & 0 \\ g_{N_2-1} & g_{N_2-2} & \cdots & g_1 & g_0 \end{pmatrix} .$$

By forming vectors of future output values and future reference values,

$$\begin{aligned}\mathbf{y}(t) &= (y(t + N_1), y(t + N_1 + 1), \dots, y(t + N_2))^T \\ \mathbf{y}^0(t) &= (y^0(t + N_1), y^0(t + N_1 + 1), \dots, y^0(t + N_2))^T,\end{aligned}$$

and a control-weighting matrix which is diagonal,

$$\Lambda = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & & 0 \\ \vdots & & \ddots & \\ 0 & 0 & & \lambda_{N_u} \end{pmatrix}$$

the cost function in (4.30) can now be written as

$$J(t) = E_t [(\mathbf{y}(t) - \mathbf{y}^0(t))^T (\mathbf{y}(t) - \mathbf{y}^0(t)) + \tilde{\mathbf{u}}(t)^T \Lambda \tilde{\mathbf{u}}(t)] \quad (4.34)$$

Furthermore, a vector of prediction errors is introduced:

$$\boldsymbol{\varepsilon}(t) = (\varepsilon(t + N_1|t), \varepsilon(t + N_1 + 1|t), \dots, \varepsilon(t + N_2|t))^T,$$

where $\varepsilon(t + j|t) = y(t + j) - \hat{y}(t + j|t)$. Then $\mathbf{y}(t) = \mathbf{G}\tilde{\mathbf{u}}(t) + \mathbf{z}(t) + \boldsymbol{\varepsilon}(t)$ which is substituted in (4.34):

$$\begin{aligned}J(t) &= \\ &E_t [(\mathbf{G}\tilde{\mathbf{u}}(t) + \mathbf{z}(t) + \boldsymbol{\varepsilon}(t) - \mathbf{y}^0(t))^T (\mathbf{G}\tilde{\mathbf{u}}(t) + \mathbf{z}(t) + \boldsymbol{\varepsilon}(t) - \mathbf{y}^0(t)) \\ &\quad + \tilde{\mathbf{u}}(t)^T \Lambda \tilde{\mathbf{u}}(t)].\end{aligned}$$

Since the j -step-ahead predictions and the corresponding prediction errors are uncorrelated, the expectation can easily be evaluated. The result becomes

$$\begin{aligned}J(t) &= (\mathbf{G}\tilde{\mathbf{u}}(t) + \mathbf{z}(t) - \mathbf{y}^0(t))^T (\mathbf{G}\tilde{\mathbf{u}}(t) + \mathbf{z}(t) - \mathbf{y}^0(t)) \\ &\quad + \tilde{\mathbf{u}}(t)^T \Lambda \tilde{\mathbf{u}}(t) + \text{constant term},\end{aligned}$$

where “constant term” is the sum of the j -step-ahead prediction error variances ($j = N_1, \dots, N_2$). Provided that $\mathbf{\Lambda}$ is chosen so that $\mathbf{G}^T \mathbf{G} + \mathbf{\Lambda}$ is positive definite, the minimum of $J(t)$ is found by setting the derivative of $J(t)$ with respect to $\tilde{\mathbf{u}}(t)$ equal to zero. We find that

$$\tilde{\mathbf{u}}(t) = (\mathbf{G}^T \mathbf{G} + \mathbf{\Lambda})^{-1} \mathbf{G}^T (\mathbf{y}^0(t) - \mathbf{z}(t)).$$

If the first row of $(\mathbf{G}^T \mathbf{G} + \mathbf{\Lambda})^{-1} \mathbf{G}^T$ is denoted \mathbf{l}^T then the current control action is

$$u(t) = u(t-1) + \mathbf{l}^T (\mathbf{y}^0(t) - \mathbf{z}(t)). \quad (4.35)$$

From (4.33) it is seen that the prediction of future output values is $\mathbf{z}(t)$ if no control action is taken ($\tilde{\mathbf{u}}(t) = \mathbf{0}$) – i.e. the future control signal is kept constant (“passive” control). Therefore the elements in the vector $\mathbf{y}^0(t) - \mathbf{z}(t)$ are predicted values of the input-free control errors, and $\mathbf{l}^T (\mathbf{y}^0(t) - \mathbf{z}(t))$ may be seen as an aggregate control error. These considerations lead to an interpretation of the GPC controller in (4.35): the optimal value of the present control is a result of an integration of aggregate input-free control errors. Clarke *et al.* (1987A) point out that the integral action ensures offset-free control.

If an off-line estimation of the model parameters is performed, the control gain vector, \mathbf{l} , can be computed before any control takes place. In case of recursive adaptive estimation, \mathbf{G} takes a new value at each sampling instant due to changes in the estimated step response function. Consequently \mathbf{l} must be recomputed at each sampling instant, and this, in turn, implies that the inversion of the matrix $\mathbf{G}^T \mathbf{G} + \mathbf{\Lambda}$ must be repeated.

4.4 Generalized Predictive Control for Embedded Models

This section describes how the GPC strategy can be extended to handle systems governed by models with embedded parameter variations. Compared to the presentation in Section 4.3, a more general point of view is used in the development of the controllers. Actually it is found that the GPC method presented by Clarke *et al.* (1987A) is a special case of the method stated here. The main difference is that time-variation of the model parameters is permitted, but a generalization of the cost function and the associated constraints is also introduced.

4.4.1 Model Structure and Output Prediction

The ARMAX Model

The control considered here is called XGPC, eXtended Generalized Predictive Control, and it is based upon the assumption that the j -step-ahead prediction of the system output can be expressed as a linear function of present and future controls (input values). The ARMAX model,

$$A_t(q^{-1})y(t) = B_t(q^{-1})u(t) + C_t(q^{-1})e(t), \quad (4.36)$$

belongs to a transfer function model structure which supports this assumption. In (4.36) $\{e(t)\}$ is white noise with the mean zero, and A , B and C are time-varying polynomials in q^{-1} :

$$\begin{aligned} A_t(q^{-1}) &= 1 + a_{1,t}q^{-1} + \cdots + a_{n,t}q^{-n} \\ B_t(q^{-1}) &= b_{1,t}q^{-1} + \cdots + b_{m,t}q^{-m} \\ C_t(q^{-1}) &= 1 + c_{1,t}q^{-1} + \cdots + c_{r,t}q^{-r}. \end{aligned}$$

The time-variation of the parameters in these polynomials is assumed to be described by known functions of time. Note that if both A and B have a root equal to 1, C is 1, and the model parameters are time-invariant, the ARIX model used for the formulation of the traditional GPC in Section 4.3.1 is obtained.

The j -Step-Ahead Predictor

The optimal j -step-ahead predictor, $\hat{y}(t+j|t)$, corresponding to (4.36) is found as the conditional expectation of $y(t+j)$ conditioned on observations of output up to and including time t . Due to the time-variation of the parameters the predictor *cannot* be obtained by means of a Diophantine equation. Instead (Madsen (1989))

$$\hat{y}(t+j|t) = -\sum_{i=1}^n a_{i,t+j} \hat{y}(t+j-i|t) + \sum_{i=1}^m b_{i,t+j} u(t+j-i) + \sum_{i=1}^r c_{i,t+j} \hat{e}(t+j-i|t), \quad j \geq 1 \quad (4.37)$$

$$\hat{y}(t+j|t) = y(t+j), \quad j < 1, \quad (4.38)$$

where

$$\hat{e}(s|t) = \begin{cases} e(s) = y(s) - \hat{y}(s|s-1) & \text{if } s \leq t \\ 0 & \text{if } s > t \end{cases}, \quad (4.39)$$

can be used. Thus the one-step-ahead prediction can partially be computed from known data ($y(t), y(t-1), \dots$ and $u(t-1), u(t-2), \dots$), but the prediction becomes a function of the unrealized control action $u(t)$. If $y(t+1)$ was known, the two-step-ahead prediction could be computed as a one-step-ahead prediction from time origin $t+1$ and would then become a function of $u(t)$ and $u(t+1)$. As $y(t+1)$ is unknown, it is replaced by $\hat{y}(t+1|t)$, and the unknown $e(t+1)$ is replaced by 0 (conditional expectations). In general, a j -step-ahead

prediction can be computed as a one-step-ahead prediction from time origin $t + j - 1$ where unknown output values are replaced by their predicted counterparts and unknown noise components are fixed to 0. This is illustrated by an example:

Example 4.1 Consider the ARX model ($n = 1, m = 2$)

$$y(t) + a_{1,t}y(t-1) = b_{1,t}u(t-1) + b_{2,t}u(t-2) + e(t).$$

The one-step-ahead predictor is a function of $u(t)$ and is found as

$$\begin{aligned}\hat{y}(t+1|t) &= E[y(t+1)|\mathcal{Y}(t)] \\ &= E[-a_{1,t+1}y(t) + b_{1,t+1}u(t) + b_{2,t+1}u(t-1) + e(t+1)|\mathcal{Y}(t)] \\ &= -a_{1,t+1}y(t) + b_{1,t+1}u(t) + b_{2,t+1}u(t-1),\end{aligned}$$

where $\mathcal{Y}(t) = (y(t), y(t-1), \dots, u(t), u(t-1), \dots)$. The two-step-ahead predictor becomes

$$\begin{aligned}\hat{y}(t+2|t) &= E[y(t+2)|\mathcal{Y}(t)] \\ &= E[-a_{1,t+2}y(t+1) + b_{1,t+2}u(t+1) + b_{2,t+2}u(t) \\ &\quad + e(t+2)|\mathcal{Y}(t)] \\ &= -a_{1,t+2}\hat{y}(t+1|t) + b_{1,t+2}u(t+1) + b_{2,t+2}u(t).\end{aligned}$$

If similar calculations are made for prediction horizons 3, 4, 5, ..., the predictor in (4.37) is seen to hold generally - i.e.:

$$\hat{y}(t+j|t) = -a_{1,t+j}\hat{y}(t+j-1|t) + b_{1,t+j}u(t+j-1) + b_{2,t+j}u(t+j-2).$$

□

The Predictor as a Linear Function of Future Controls

From (4.37) it is found that $\hat{y}(t+j|t)$ is linear in the input values up to time $t+j-1$ and output values up to time t . This is seen by

eliminating the \hat{y} 's and the \hat{e} 's in (4.37) recursively using (4.38) and (4.39). If $v_j(t)$ denotes the sum of all terms on the right-hand side of (4.37) involving known data ($y(t), y(t-1), \dots$ and $u(t-1), u(t-2), \dots$), the j -step-ahead predictor can be expressed as

$$\hat{y}(t+j|t) = \sum_{i=1}^j \bar{h}_{j,i,t} u(t+j-i) + v_j(t), \quad j = 1, 2, 3, \dots, \quad (4.40)$$

where $\bar{h}_{j,i,t}$ ($i = 1, 2, 3, \dots, j$) are the coefficients describing the linear dependence on future controls. These coefficients and $v_j(t)$ can, of course, be found by eliminating the \hat{y} 's and the \hat{e} 's in (4.37) recursively, but there is an easier way.

Write the model (4.36) as

$$y(t) = H_t(q^{-1})u(t) + G_t(q^{-1})e(t), \quad (4.41)$$

where the time-varying transfer functions

$$H_t(q^{-1}) = \sum_{i=0}^{\infty} h_{i,t} q^{-i}$$

and

$$G_t(q^{-1}) = \sum_{i=0}^{\infty} g_{i,t} q^{-i}$$

correspond to $B(q^{-1})/A(q^{-1})$ and $C(q^{-1})/A(q^{-1})$ in the time-invariant case. The coefficients $h_{0,t}, h_{1,t}, h_{2,t}, \dots$ ($h_{0,t} = 0$) are clearly the weights of the time-varying impulse response function describing the dynamic relationship between input and output² – i.e. $h_{i,t}$ is the extra contribution to $y(t)$ if $u(t-i)$ is increased by one. From

²Note that the coefficients of $H_t(q^{-1})$ cannot be computed by ordinary polynomial division of $B_t(q^{-1})$ by $A_t(q^{-1})$ since the q^{-1} operator affects the t index of the coefficients of A and B . One way to compute the $h_{i,t}$ weights is to pass impulses through the system $A_t(q^{-1})y(t) = B_t(q^{-1})u(t)$. This technique is described below.

(4.41) the j -step ahead output prediction is easily found (conditional expectation):

$$\hat{y}(t+j|t) = \sum_{i=0}^{\infty} h_{i,t+j} u(t+j-i) + \sum_{i=j}^{\infty} g_{i,t+j} e(t+j-i), \quad j = 1, 2, 3, \dots$$

Comparison with (4.40) shows that

$$\begin{aligned} \bar{h}_{j,i,t} &= h_{i,t+j}, \quad i = 1, \dots, j \\ v_j(t) &= \sum_{i=j+1}^{\infty} h_{i,t+j} u(t+j-i) + \sum_{i=j}^{\infty} g_{i,t+j} e(t+j-i), \end{aligned}$$

where $j = 1, 2, 3, \dots$. Thus

$$\hat{y}(t+j|t) = \sum_{i=1}^j h_{i,t+j} u(t+j-i) + v_j(t). \quad (4.42)$$

Equation (4.42) shows that the $h_{i,t+j}$ -values and $v_j(t)$ can be obtained as follows:

1. To find $v_j(t)$ ($j = 1, 2, 3, \dots$): Let $u(t+j-i) = 0$ for $i = 1, 2, \dots, j$ and compute $\hat{y}(t+j|t)$ using (4.37). This "prediction" becomes $v_j(t)$.
2. To find $h_{i,t+j}$ ($j = 1, 2, \dots$ and $i = 1, \dots, j$): Assume that

$$u(t+j-l) = \begin{cases} 1 & \text{if } l = i \\ 0 & \text{otherwise} \end{cases},$$

and $y(t-l) = 0$ for $l = 0, 1, 2, \dots$ (corresponding to $v_j(t) = 0$). Under these assumptions compute $\hat{y}(t+j|t)$ using (4.37). This "prediction" becomes $h_{i,t+j}$.

In step 1 we simply compute $v_j(t)$ as the prediction of $y(t+j)$ in case that the controller is passive from time t onwards. In step 2 we

compute $h_{i,t+j}$ as the response at time $t + j$ to a unit input pulse at time $t + j - i$.

An example illustrates the technique.

Example 4.2 *Reconsider the ARX model from Example 4.1. We will find the 1, 2, 3 and 4 step predictors as functions of future controls on the form (4.42). At first $v_1(t), \dots, v_4(t)$ are computed using (4.37) with $u(t) = u(t + 1) = u(t + 2) = u(t + 3) = 0$:*

$$\begin{aligned} v_1(t) &= \hat{y}(t + 1|t) = -a_{1,t+1}y(t) + b_{2,t+1}u(t - 1) \\ v_2(t) &= \hat{y}(t + 2|t) = -a_{1,t+2}\hat{y}(t + 1|t) = -a_{1,t+2}v_1(t) \\ v_3(t) &= \hat{y}(t + 3|t) = -a_{1,t+3}\hat{y}(t + 2|t) = -a_{1,t+3}v_2(t) \\ v_4(t) &= \hat{y}(t + 4|t) = -a_{1,t+4}\hat{y}(t + 3|t) = -a_{1,t+4}v_3(t). \end{aligned}$$

Note that $v_1(t), \dots, v_3(t)$ appear as intermediate results during the computation of $v_4(t)$. Furthermore, the computations are very simple as $v_j(t)$ is equal to $-a_{1,t+j}v_{j-1}(t)$ for $j > 1$.

Now $h_{1,t+1}, h_{2,t+2}, h_{3,t+3}$ and $h_{4,t+4}$ are found using (4.37) with $u(t) = 1, u(i) = 0$ for $i \neq t$ and $y(i) = 0$ for $i \leq t$:

$$\begin{aligned} h_{1,t+1} &= \hat{y}(t + 1|t) &&= b_{1,t+1} \\ h_{2,t+2} &= \hat{y}(t + 2|t) = -a_{1,t+2}\hat{y}(t + 1|t) + b_{2,t+2} &&= -a_{1,t+2}h_{1,t+1} \\ &&&+ b_{2,t+2} \\ h_{3,t+3} &= \hat{y}(t + 3|t) = -a_{1,t+3}\hat{y}(t + 2|t) &&= -a_{1,t+3}h_{2,t+2} \\ h_{4,t+4} &= \hat{y}(t + 4|t) = -a_{1,t+4}\hat{y}(t + 3|t) &&= -a_{1,t+4}h_{3,t+3}. \end{aligned}$$

Note that $h_{j,t+j} = -a_{1,t+j}h_{j-1,t+j-1}$ for $j > 2$.

If $u(t + 1) = 1, u(i) = 0$ for $i \neq t + 1$ and $y(i) = 0$ for $i \leq t$ we

obtain:

$$\begin{aligned} \hat{y}(t+1|t) &= 0 \\ h_{1,t+2} &= \hat{y}(t+2|t) = b_{1,t+2} \\ h_{2,t+3} &= \hat{y}(t+3|t) = -a_{1,t+3}h_{1,t+2} + b_{2,t+3} \\ h_{3,t+4} &= \hat{y}(t+4|t) = -a_{1,t+4}h_{2,t+3} , \end{aligned}$$

and if $u(t+2) = 1$, $u(i) = 0$ for $i \neq t+2$ and $y(i) = 0$ for $i \leq t$:

$$\begin{aligned} \hat{y}(t+1|t) &= 0 \\ \hat{y}(t+2|t) &= 0 \\ h_{1,t+3} &= \hat{y}(t+3|t) = b_{1,t+3} \\ h_{2,t+4} &= \hat{y}(t+4|t) = -a_{1,t+4}h_{1,t+3} + b_{2,t+4} , \end{aligned}$$

and finally if $u(t+3) = 1$, $u(i) = 0$ for $i \neq t+3$ and $y(i) = 0$ for $i \leq t$:

$$\begin{aligned} \hat{y}(t+1|t) &= 0 \\ \hat{y}(t+2|t) &= 0 \\ \hat{y}(t+3|t) &= 0 \\ h_{1,t+4} &= \hat{y}(t+4|t) = b_{1,t+4} . \end{aligned}$$

If the above results are inserted into (4.42), the output predictions up to 4 steps ahead can be written as functions of future controls. Furthermore, it can be concluded that the impulse response function of the system at time t becomes

$$h_{i,t} = \begin{cases} 0 & \text{if } i = 0 \\ b_{1,t} & \text{if } i = 1 \\ -a_{1,t}h_{1,t-1} + b_{2,t} & \text{if } i = 2 \\ -a_{1,t}h_{i-1,t-1} & \text{if } i \geq 3 \end{cases}$$

□

Now introduce a maximum prediction horizon N (≥ 1) which corresponds to N_2 in the traditional GPC (see Section 4.3.2). Considering (4.42) it is found that the j -step-ahead predictions, j running from 1 up to N , can be written as a linear matrix expression:

$$\hat{\mathbf{y}}(t) = \mathbf{H}(t)\mathbf{u}(t) + \mathbf{v}(t), \quad (4.43)$$

where

$$\begin{aligned} \hat{\mathbf{y}}(t) &= (\hat{y}(t+1|t), \dots, \hat{y}(t+N|t))^T \\ \mathbf{u}(t) &= (u(t), \dots, u(t+N-1))^T \\ \mathbf{v}(t) &= (v_1(t), \dots, v_N(t))^T \\ \mathbf{H}(t) &= \begin{pmatrix} h_{1,t+1} & 0 & \cdots & 0 & 0 \\ h_{2,t+2} & h_{1,t+2} & \cdots & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots \\ h_{N-1,t+N-1} & h_{N-2,t+N-1} & \cdots & h_{1,t+N-1} & 0 \\ h_{N,t+N} & h_{N-1,t+N} & \cdots & h_{2,t+N} & h_{1,t+N} \end{pmatrix}. \end{aligned}$$

The matrix $\mathbf{H}(t)$ corresponds to \mathbf{G} in Section 4.3.3. However, the rows of $\mathbf{H}(t)$ are time-varying impulse responses while the rows of \mathbf{G} are time-invariant step responses. Furthermore, \mathbf{G} is a result of truncating $N_1 - 1$ rows and $N_2 - N_u$ columns in a $N_2 \times N_2$ matrix while $\mathbf{H}(t)$ has full dimensions $N \times N$. For time-invariant systems the following equation determines the relationship between the weights, h_i ($i = 0, 1, 2, \dots$), in the impulse response function and the weights, g_j ($j = 0, 1, 2, \dots$), in the step response function

$$\sum_{i=0}^{\infty} h_i q^{-i} = \nabla \sum_{j=0}^{\infty} g_j q^{-j-1}.$$

In (4.43) it has been assumed that the same model is applied for all prediction horizons. Actually, this need not be the case. If different models are used, the j th row of $\mathbf{H}(t)$ and the j th element of $\mathbf{v}(t)$

belong to a special model designed for j -step-ahead prediction. Making use of an individual model for each horizon is often relevant if a non-linear system is approximated by linear models, i.e. if the optimal linearized predictor for the system depends on the prediction horizon.

4.4.2 The Cost Function

The control strategy is based on minimization of a cost function very similar to, yet more general than the ordinary GPC cost function in (4.28):

$$\begin{aligned} \min_{\mathbf{u}(t)} J(\Gamma(t), \Lambda(t), \omega(t); t, \mathbf{u}(t)) \\ = E_t[(\mathbf{y}(t) - \mathbf{y}^0(t))^T \Gamma(t) (\mathbf{y}(t) - \mathbf{y}^0(t)) \\ + \mathbf{u}(t)^T \Lambda(t) \mathbf{u}(t) + 2\omega(t)^T \mathbf{u}(t)] , \end{aligned} \quad (4.44)$$

where $\mathbf{y}(t)$ is a vector of future output values and $\mathbf{y}^0(t)$ is a vector of future set-points:

$$\begin{aligned} \mathbf{y}(t) &= (y(t+1), \dots, y(t+N))^T \\ \mathbf{y}^0(t) &= (y^0(t+1), \dots, y^0(t+N))^T \end{aligned}$$

and

$\Gamma(t)$ is a positive semidefinite³ and symmetric matrix which weights the control errors,

$\Lambda(t)$ is a positive semidefinite³ and symmetric matrix which weights the squared control values, and

$\omega(t)$ is a vector weighting the control values linearly.

³The matrix $\mathbf{H}(t)^T \Gamma(t) \mathbf{H}(t) + \Lambda(t)$ must be non-singular.

Choice of the Design Parameters of the Controller

Apart from the fact that the $\Gamma(t)$ -matrix should be symmetric and positive semidefinite, no further restrictions apply to it. However, it may, for instance, be the identity matrix as in the ordinary GPC presented by Clarke *et al.* (1987A). Another possibility is to choose the inverse covariance matrix of the prediction errors, i.e. $\Gamma(t) = \text{Var}[\varepsilon(t)]^{-1}$, where $\varepsilon(t) = \mathbf{y}(t) - \hat{\mathbf{y}}(t)$. This choice means that less attention is paid to the control error if the corresponding prediction is very uncertain. This seems reasonable since the current control cannot do very much about a prediction error which is primarily a result of stochastic disturbances.

In order to evaluate $\text{Var}[\varepsilon(t)]$, the relationship between the j -step-ahead prediction errors and future values of $e(t)$ is required. The following example shows how to find such a relationship and the covariance matrix.

Example 4.3 *The AR model corresponding to the ARX model in Example 4.1 is*

$$y(t) + a_{1,t}y(t-1) = e(t).$$

The deterministic terms of the ARX model involving $u(t)$ have been suppressed since they do not have any influence on the prediction error, $\varepsilon(t+j|t) = y(t+j) - \hat{y}(t+j|t)$ ($j \geq 1$). The prediction errors are easily found from the AR model:

$$\begin{aligned} \varepsilon(t+1|t) &= e(t+1) \\ \varepsilon(t+2|t) &= e(t+2) - a_{1,t+2}\varepsilon(t+1|t) \\ \varepsilon(t+3|t) &= e(t+3) - a_{1,t+3}\varepsilon(t+2|t) \\ \varepsilon(t+4|t) &= e(t+4) - a_{1,t+4}\varepsilon(t+3|t) \\ &\vdots \qquad \qquad \qquad \vdots \end{aligned}$$

Now element (i, l) in $\text{Var}[\varepsilon(t)]$ is obtained as

$$\text{Cov}[\varepsilon(t+i|t), \varepsilon(t+l|t)].$$

For $N = 3$ the covariance matrix $\text{Var}[\varepsilon(t)]$ becomes

$$\sigma_e^2 \begin{pmatrix} 1 & -a_{1,t+2} & a_{1,t+3}a_{1,t+2} \\ -a_{1,t+2} & 1 + a_{1,t+2}^2 & -a_{1,t+3}(1 + a_{1,t+2}^2) \\ a_{1,t+3}a_{1,t+2} & -a_{1,t+3}(1 + a_{1,t+2}^2) & 1 + a_{1,t+3}^2(1 + a_{1,t+2}^2) \end{pmatrix}.$$

□

If the control strategy should aim at keeping the control values close to zero, $\Lambda(t) = \lambda I$ ($\lambda > 0$) and $\omega = \mathbf{0}$ would be appropriate. If, on the contrary, the control value $u(t)$ should be kept close to some pre-specified value $u^0(t)$ and the control increments should be simultaneously damped, $\Lambda(t)$ and $\omega(t)$ may be chosen such that

$$\begin{aligned} & \mathbf{u}(t)^T \Lambda(t) \mathbf{u}(t) + 2\omega(t)^T \mathbf{u}(t) + \text{const.} = \\ & \sum_{j=1}^N \left(\lambda_{1,j} [u(t+j-1) - u^0(t+j-1)]^2 + \lambda_{2,j} [\nabla u(t+j-1)]^2 \right) \end{aligned}$$

where $\lambda_{1,j}$ and $\lambda_{2,j}$ ($1 \leq j \leq N$) are non-negative real numbers and “const.” is a term which is independent of $\mathbf{u}(t)$ but dependent on $u^0(t)$. This leads to

$$\Lambda(t) = \begin{pmatrix} \lambda_{1,1} + \lambda_{2,1} + \lambda_{2,2} & -\lambda_{2,2} & \cdots & 0 & 0 \\ -\lambda_{2,2} & \lambda_{1,2} + \lambda_{2,2} + \lambda_{2,3} & & 0 & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & & \lambda_{1,N-1} + \lambda_{2,N-1} + \lambda_{2,N} & -\lambda_{2,N} \\ 0 & 0 & \cdots & -\lambda_{2,N} & \lambda_{1,N} + \lambda_{2,N} \end{pmatrix} \quad (4.45)$$

and

$$\omega(t) = - \begin{pmatrix} \lambda_{1,1}u^0(t) + \lambda_{2,1}u(t-1) \\ \lambda_{1,2}u^0(t+1) \\ \lambda_{1,3}u^0(t+2) \\ \vdots \\ \lambda_{1,N}u^0(t+N-1) \end{pmatrix}. \quad (4.46)$$

Particularly if $\lambda_{1,j} = 0$ and $\lambda_{2,j} = \lambda_j$ ($j = 1, \dots, N$), the controls are weighted as in the ordinary GPC cost function (4.28) (Clarke *et al.* (1987A)).

It is also possible to choose $\Lambda(t)$ and $\omega(t)$ so that filtered controls, $u^f(t)$, are weighted in (4.44). Let

$$u^f(t) = F_t(q^{-1})u(t),$$

where $F_t(q^{-1})$ is a filter given by

$$F_t(q^{-1}) = 1 + f_{1,t}q^{-1} + f_{2,t}q^{-2} + \dots + f_{p,t}q^{-p}.$$

If, for instance, $u^f(t)$ should be weighted as

$$\sum_{j=1}^N \lambda_j [u^f(t+j-1)]^2,$$

the corresponding structures of $\Lambda(t)$ and $\omega(t)$ would turn out to be

$$\Lambda(t) = \begin{pmatrix} \overbrace{x \ \cdots \ x}^{p+1 \text{ columns}} & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ x & \cdots & x & \cdots & x & 0 \\ & \ddots & \vdots & \vdots & \ddots & \\ 0 & & x & \cdots & x & \cdots & x \\ \vdots & \ddots & & \ddots & \vdots & & \vdots \\ 0 & \cdots & 0 & & x & \cdots & x \end{pmatrix}$$

$$\omega(t) = \left(\overbrace{x \cdots x}^{p \text{ elements}} \ 0 \cdots 0 \right)^T,$$

where the x 's represent non-zero elements.

From the above discussion one can conclude that the cost function in (4.44) offers a very flexible criterion for defining the optimal control.

The cost horizon, N , which implicitly is included in the control criterion, should be chosen so that all future output values significantly affected by the current control are included in the cost function. The impulse response function of the system provides a natural tool for finding an appropriate value of N (cf. the advice for the choice of N_2 in Section 4.3.2).

4.4.3 The Resulting Controller

To minimize the cost function in (4.44) one can use the same technique as used in Section 4.3.3: The expectation of the cost function is evaluated and the resulting function is minimized by equating its derivative with respect to future controls with zero. If a vector of j -step-ahead prediction errors ($j = 1, \dots, N$),

$$\varepsilon(t) = (\varepsilon(t+1|t), \varepsilon(t+2|t), \dots, \varepsilon(t+N|t))^T,$$

is introduced, the expectation of the first term of the cost function can be written as

$$\begin{aligned} & E_t [(\hat{\mathbf{y}}(t) + \varepsilon(t) - \mathbf{y}^0(t))^T \mathbf{\Gamma}(t) (\hat{\mathbf{y}}(t) + \varepsilon(t) - \mathbf{y}^0(t))] \\ &= (\hat{\mathbf{y}}(t) - \mathbf{y}^0(t))^T \mathbf{\Gamma}(t) (\hat{\mathbf{y}}(t) - \mathbf{y}^0(t)) + E_t [\varepsilon(t)^T \mathbf{\Gamma}(t) \varepsilon(t)], \end{aligned}$$

where the fact that the predictions and the prediction errors are independent has been used. If we substitute for $\hat{\mathbf{y}}(t)$ from (4.43) the

cost function becomes:

$$\begin{aligned} J(\mathbf{\Gamma}(t), \mathbf{\Lambda}(t), \boldsymbol{\omega}(t); t, \mathbf{u}(t)) \\ = (\mathbf{H}(t)\mathbf{u}(t) + \mathbf{v}(t) - \mathbf{y}^0(t))^T \mathbf{\Gamma}(t) (\mathbf{H}(t)\mathbf{u}(t) + \mathbf{v}(t) - \mathbf{y}^0(t)) \\ + \mathbf{u}(t)^T \mathbf{\Lambda}(t) \mathbf{u}(t) + 2\boldsymbol{\omega}(t)^T \mathbf{u}(t) + E_t [\boldsymbol{\varepsilon}(t)^T \mathbf{\Gamma}(t) \boldsymbol{\varepsilon}(t)] . \end{aligned}$$

The last term in this expression does not depend on $\mathbf{u}(t)$ so the optimal control can be found by minimizing:

$$\begin{aligned} \tilde{J}(t, \mathbf{u}(t)) = (\mathbf{H}(t)\mathbf{u}(t) + \boldsymbol{\beta}(t))^T \mathbf{\Gamma}(t) (\mathbf{H}(t)\mathbf{u}(t) + \boldsymbol{\beta}(t)) \\ + \mathbf{u}(t)^T \mathbf{\Lambda}(t) \mathbf{u}(t) + 2\boldsymbol{\omega}(t)^T \mathbf{u}(t) , \end{aligned} \quad (4.47)$$

where

$$\boldsymbol{\beta}(t) = \mathbf{v}(t) - \mathbf{y}^0(t) . \quad (4.48)$$

Computing the derivative of this function with respect to $\mathbf{u}(t)$ gives (omitting the time arguments):

$$\frac{\partial \tilde{J}(\mathbf{u})}{\partial \mathbf{u}} = 2\mathbf{H}^T \mathbf{\Gamma} \mathbf{H} \mathbf{u} + 2\mathbf{H}^T \mathbf{\Gamma}^T \boldsymbol{\beta} + 2\mathbf{\Lambda} \mathbf{u} + 2\boldsymbol{\omega} .$$

The optimal control is obtained by equating this expression with zero and solving the resulting equation with respect to \mathbf{u} :

$$\mathbf{u}(t) = - [\mathbf{H}(t)^T \mathbf{\Gamma}(t) \mathbf{H}(t) + \mathbf{\Lambda}(t)]^{-1} [\mathbf{H}(t)^T \mathbf{\Gamma}(t)^T \boldsymbol{\beta}(t) + \boldsymbol{\omega}(t)] . \quad (4.49)$$

Here it is assumed that $\mathbf{\Gamma}(t)$ and $\mathbf{\Lambda}(t)$ have been chosen so that the matrix $\mathbf{H}(t)^T \mathbf{\Gamma}(t) \mathbf{H}(t) + \mathbf{\Lambda}(t)$ becomes non-singular. As only the first element of the control vector, $\mathbf{u}(t)$, is implemented, the control law may be written

$$\mathbf{u}(t) = -\mathbf{l}(t) [\mathbf{H}(t)^T \mathbf{\Gamma}(t)^T \boldsymbol{\beta}(t) + \boldsymbol{\omega}(t)] , \quad (4.50)$$

where $\mathbf{l}(t)$ is the first row of $[\mathbf{H}(t)^T \mathbf{\Gamma}(t) \mathbf{H}(t) + \mathbf{\Lambda}(t)]^{-1}$.

From (4.43) it appears that $\mathbf{v}(t)$ is a vector of input-free predictions, i.e. the predictions if $\mathbf{u}(t) = \mathbf{0}$, and the elements in the vector $\boldsymbol{\beta}(t)$

defined in (4.48) are predicted values of the input-free control errors. Thus the optimal control in (4.50) is obtained as a linear feedback from the predicted values of the input-free control errors plus a constant term coming from $\omega(t)$.

4.4.4 Linear Equality Constraints

For some practical systems it might be desirable to make the sequence of future controls (or some of them) follow a more or less fixed time function. In the GPC developed by Clarke *et al.* (1987A) (see Section 4.3), for instance, only the first N_u controls are free while the subsequent controls are considered to be constant, equal to the last free control. This restriction makes up a set of simple linear constraints: $u(t + j - 1) = u(t + j - 2)$ for $j > N_u$. The advantage of introducing such constraints is obvious: the dimension of the optimization problem associated with the computation of the optimal control is reduced to N_u . Furthermore, it does not reduce the performance of the controller significantly when reducing the degrees of freedom of the control in this way as long as the value of N_u is not too low compared to the complexity of the system.

In this section it is shown how to add general type linear constraints to the minimization criterion in (4.44) (or the equivalent in (4.47)). Consider the following set of constraints:

$$S(t)\mathbf{u}(t) = \mathbf{d}(t), \quad (4.51)$$

where the matrix $S(t)$ is $M \times N$ ($M \leq N$) and has full rank ($= M$). The vector $\mathbf{d}(t)$ is a column with M elements. To solve a minimization problem being subject to (4.51), the M equations can simply be solved with respect to a subset of the elements in $\mathbf{u}(t)$. This subset of the future control values can then be eliminated in the cost function. Since the cost function is a quadratic function of

the unknown controls, and the constraints are linear, this elimination leads to an unconstrained quadratic cost function. The procedure is shown below.

Let us split up the \mathbf{u} -vector⁴ into two parts, \mathbf{u}_1 and \mathbf{u}_2 :

$$\mathbf{u} = [\mathbf{u}_1^T \ \mathbf{u}_2^T]^T,$$

where \mathbf{u}_1 is a vector of the dimension $N - M$ and \mathbf{u}_2 is an M -vector. S is partitioned similarly,

$$S = [S_1 \ S_2],$$

where it is assumed that the columns of S_2 are linearly independent (i.e. S_2 is a non-singular matrix). (If, however, the last M columns of S are linearly dependent it is necessary to introduce suitable column permutations to ensure that S_2 becomes non-singular. The elements of \mathbf{u} must, of course, be permuted correspondingly.)

Now (4.51) can be written as

$$S_1 \mathbf{u}_1 + S_2 \mathbf{u}_2 = \mathbf{d}.$$

Solving this equation with respect to \mathbf{u}_2 yields

$$\mathbf{u}_2 = S_2^{-1}(\mathbf{d} - S_1 \mathbf{u}_1). \quad (4.52)$$

Matrices and vectors in (4.47) are partitioned so that they match the partitioning of \mathbf{u} :

$$\begin{aligned} H &= [H_1 \ H_2] \\ \Lambda &= \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix} \\ \omega &= [\omega_1^T \ \omega_2^T]^T. \end{aligned}$$

⁴The t -arguments are omitted in the derivation of the constrained cost function.

Using this (4.47) can be rewritten as

$$\begin{aligned} \tilde{J}([\mathbf{u}_1^T \ \mathbf{u}_2^T]^T) &= (\mathbf{H}_1 \mathbf{u}_1 + \mathbf{H}_2 \mathbf{u}_2 + \boldsymbol{\beta})^T \Gamma (\mathbf{H}_1 \mathbf{u}_1 + \mathbf{H}_2 \mathbf{u}_2 + \boldsymbol{\beta}) \\ &\quad + \mathbf{u}_1^T \boldsymbol{\Lambda}_{11} \mathbf{u}_1 + 2\mathbf{u}_1^T \boldsymbol{\Lambda}_{12} \mathbf{u}_2 + \mathbf{u}_2^T \boldsymbol{\Lambda}_{22} \mathbf{u}_2 + 2\boldsymbol{\omega}_1^T \mathbf{u}_1 + 2\boldsymbol{\omega}_2^T \mathbf{u}_2, \end{aligned}$$

where the fact that $\boldsymbol{\Lambda}$ is symmetric has been utilized (i.e. $\boldsymbol{\Lambda}_{12} = \boldsymbol{\Lambda}_{21}^T$). By substituting the expression in (4.52) for \mathbf{u}_2 , it is found that

$$\begin{aligned} \tilde{J}([\mathbf{u}_1^T \ \mathbf{u}_2^T]^T) \Big|_{\mathbf{u}_2 = \mathbf{S}_2^{-1}(\mathbf{d} - \mathbf{S}_1 \mathbf{u}_1)} &= \\ (\mathbf{H}^* \mathbf{u}_1 + \boldsymbol{\beta}^*)^T \Gamma (\mathbf{H}^* \mathbf{u}_1 + \boldsymbol{\beta}^*) &+ \mathbf{u}_1^T \boldsymbol{\Lambda}^* \mathbf{u}_1 + 2\boldsymbol{\omega}^{*T} \mathbf{u}_1 + \text{const.}, \end{aligned}$$

where “const.” means “independent of \mathbf{u} ” and

$$\begin{aligned} \mathbf{H}^* &= \mathbf{H}_1 - \mathbf{H}_2 \mathbf{S}_2^{-1} \mathbf{S}_1 \\ \boldsymbol{\beta}^* &= \mathbf{H}_2 \mathbf{S}_2^{-1} \mathbf{d} + \boldsymbol{\beta} \\ \boldsymbol{\Lambda}^* &= \boldsymbol{\Lambda}_{11} + (\mathbf{S}_1^T (\mathbf{S}_2^{-1})^T \boldsymbol{\Lambda}_{22} - 2\boldsymbol{\Lambda}_{12}) \mathbf{S}_2^{-1} \mathbf{S}_1 \\ \boldsymbol{\omega}^* &= (\boldsymbol{\Lambda}_{12} - \mathbf{S}_1^T (\mathbf{S}_2^{-1})^T \boldsymbol{\Lambda}_{22}) \mathbf{S}_2^{-1} \mathbf{d} + \boldsymbol{\omega}_1 - \mathbf{S}_1^T (\mathbf{S}_2^{-1})^T \boldsymbol{\omega}_2. \end{aligned} \tag{4.53}$$

Thus the optimal control is found by minimizing

$$\tilde{J}^*(\mathbf{u}_1) = (\mathbf{H}^* \mathbf{u}_1 + \boldsymbol{\beta}^*)^T \Gamma (\mathbf{H}^* \mathbf{u}_1 + \boldsymbol{\beta}^*) + \mathbf{u}_1^T \boldsymbol{\Lambda}^* \mathbf{u}_1 + 2\boldsymbol{\omega}^{*T} \mathbf{u}_1,$$

with respect to \mathbf{u}_1 . The structure of this cost function is similar to the structure of the cost function in (4.47). Therefore, the solution to the constrained optimization problem is obtained by replacing \mathbf{H} , $\boldsymbol{\beta}$, $\boldsymbol{\Lambda}$ and $\boldsymbol{\omega}$ in (4.49) by \mathbf{H}^* , $\boldsymbol{\beta}^*$, $\boldsymbol{\Lambda}^*$ and $\boldsymbol{\omega}^*$, respectively.

The Constraints used by Clarke *et al.*

Below it is shown how to choose \mathbf{H}^* , $\boldsymbol{\beta}^*$, $\boldsymbol{\Lambda}^*$ and $\boldsymbol{\omega}^*$ in order to obtain the GPC method presented by Clarke *et al.* (1987A). According to (4.29), this controller is subject to

$$u(t+j-1) = u(t+j-2), \quad j > N_u > 0.$$

These constraints can easily be expressed by means of (4.51). If, for instance, $N = 5$ and $N_u = 2$, then $M = N - N_u = 3$ and

$$S(t) = \underbrace{\begin{pmatrix} 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & -1 & 1 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{pmatrix}}_{S_1(t)}, \quad d(t) = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}.$$

If $N_1 = 1$ and $N_2 = N = 5$ in the cost function in (4.28) it is found that (cf. (4.45) and (4.46)):

$$\Lambda = \left(\begin{array}{cc|ccc} \lambda_1 + \lambda_2 & -\lambda_2 & 0 & 0 & 0 \\ -\lambda_2 & \lambda_2 + \lambda_3 & -\lambda_3 & 0 & 0 \\ \hline 0 & -\lambda_3 & \lambda_3 + \lambda_4 & -\lambda_4 & 0 \\ 0 & 0 & -\lambda_4 & \lambda_4 + \lambda_5 & -\lambda_5 \\ 0 & 0 & 0 & -\lambda_5 & \lambda_5 \end{array} \right)$$

$$\omega(t) = \underbrace{(-\lambda_1 u(t-1) \ 0 \ 0 \ 0 \ 0)^T}_{\omega_1(t)^T} \quad \underbrace{0 \ 0 \ 0 \ 0 \ 0^T}_{\omega_2(t)^T}$$

$$\Gamma = I \quad (5 \times 5 \text{ identity matrix}).$$

For the time-invariant ARIMAX model used by Clarke *et al.* (1987A) the “impulse response matrix”, H , introduced in (4.43) becomes time-invariant

$$H = \left(\begin{array}{cc|ccc} h_1 & 0 & 0 & 0 & 0 \\ h_2 & h_1 & 0 & 0 & 0 \\ h_3 & h_2 & h_1 & 0 & 0 \\ h_4 & h_3 & h_2 & h_1 & 0 \\ h_5 & h_4 & h_3 & h_2 & h_1 \end{array} \right).$$

Observing that

$$S_2^{-1} S_1 = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 0 & -1 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & -1 \\ 0 & -1 \\ 0 & -1 \end{pmatrix},$$

it is readily found from (4.53) that

$$\mathbf{H}^* = \mathbf{H}_1 + \mathbf{H}_2 \begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} h_1 & 0 \\ h_2 & h_1 \\ h_3 & h_1 + h_2 \\ h_4 & h_1 + h_2 + h_3 \\ h_5 & h_1 + h_2 + h_3 + h_4 \end{pmatrix}$$

$$\boldsymbol{\beta}^*(t) = \boldsymbol{\beta}(t)$$

$$\boldsymbol{\Lambda}^* = \begin{pmatrix} \lambda_1 + \lambda_2 & -\lambda_2 \\ -\lambda_2 & \lambda_2 \end{pmatrix}$$

$$\boldsymbol{\omega}^*(t) = \boldsymbol{\omega}_1(t) = (-\lambda_1 u(t-1) \ 0)^T,$$

where the vector $\boldsymbol{\beta}(t)$ contain the predicted values of the input-free control errors (defined in (4.48)).

Note that the first column of \mathbf{H}^* is equal to the first column of \mathbf{H} , and the second column is the sum of columns 2 to 5 of \mathbf{H} . As a general rule it can be established that the first $N_u - 1$ columns of \mathbf{H}^* are copied from \mathbf{H} while the N_u th column is the sum of the remaining columns of \mathbf{H} . This result was also heuristically proposed by Bjerre (1992). Note also that $\boldsymbol{\Lambda}^*$ and $\boldsymbol{\omega}^*$ depend on λ_1 and λ_2 , but not on λ_3 , λ_4 and λ_5 . This is quite natural since the latter λ 's have no influence on the optimal solution as the corresponding control increments have been fixed by means of the equality constraints. Moreover, this is in conformity with the results of Clarke *et al.* (1987A).

4.5 Simulation Experiments with XGPC

In this section, simulation experiments are used to compare the performance of the XGPC method with the suboptimal minimal variance control method. XGPC controllers with different choice of the tuning parameters are also compared.

In Section 4.1 it was pointed out that a multi-step predictive control is very relevant for control of the supply temperature in case that mass flow in the district heating network is near its maximum value. Section 4.2 describes how to use heat load models as a basis for control of the supply temperature, but it also suggests that models of the mass flow should be more adequate for control purpose.

A Simple Model of the Mass Flow

In this section the following model is considered

$$y(t) = -ay(t-1) + b_3u(t-3) + b_4u(t-4) + b_5u(t-5) + f(t) + e(t), \quad (4.54)$$

where

$$\begin{aligned} a &= -0.8, & b_3 &= -1, & b_4 &= -3.1, \\ b_5 &= -1.5, & \sigma_e &= 0.2, & f(t) &= 2 \sin\left(\frac{\pi}{12}t\right). \end{aligned}$$

This is a simple model of the relationship between the supply temperature, $u(t)$, and the mass flow, $y(t)$ (scaled variation around the mean levels). The impulse response function (Figure 4.8) shows that an impulse in the supply temperature results in a temporary reduction of the mass flow after some delay, here after 3 steps (cf. the discussion in Section 4.1).

The sampling interval is assumed to be one hour. Thus the term $2 \sin(\frac{\pi}{12}t)$ in (4.54) simulate a simple diurnal variation of the heat demand.

The roots of the polynomial

$$q^3B(q^{-1}) = -1 - 3.1q^{-1} - 1.5q^{-2}$$

are -0.6 and -2.5, i.e. one root outside the unit circle. Hence the model represents a non-minimum phase system. Consequently an ordinary

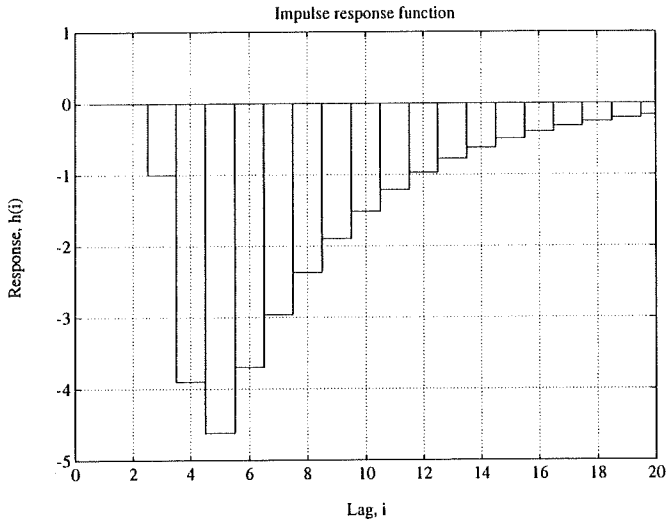


Figure 4.8: *Impulse response function from $u(t)$ to $y(t)$ in model (4.54).*

minimal variance controller corresponding to the model in (4.54) will become unstable. Therefore, in the simulation studies a suboptimal minimal variance control (SMVC) method which is stable is considered. The method implies that the unstable root is replaced by its reciprocal, see Åström (1970) and Sjøgaard (1988). The stabilized polynomial becomes

$$q^3 B'(q^{-1}) = -2.5 - 2.5q^{-1} - 0.6q^{-2} . \quad (4.55)$$

Reference Values

Both the SMVC and the XGPC require that output reference values are specified. These reference values are determined by probability considerations similar to those in Section 4.2. Suppose that the maximum value of the flow is $y_{\max} = 0$ (for this study, the actual choice of the maximum value is not important since it does not affect the dynamic and stochastic properties of the control signals). Then, at time t , the reference value, $y^0(t + j|t)$, is determined by

$$P\{y(t + j) \leq 0\} = \pi , \quad j > 0 , \quad (4.56)$$

where $\pi = 0.99$ is used. Another possibility would be to restrict the joint probability of having future output values less than zero. Note that a future reference value is a function of the present time. Thus, in general, $y^0(t + j|t) \neq y^0(t + j|t + 1)$. Since $y(t + j)$ is a sum of a prediction and a prediction error, i.e.

$$y(t + j) = \hat{y}(t + j|t) + \varepsilon(t + j|t) ,$$

$y^0(t + j|t)$ is found by inserting this into (4.56) with $\hat{y}(t + j|t) = y^0(t + j|t)$:

$$P\{y^0(t + j|t) + \varepsilon(t + j|t) \leq 0\} = \pi .$$

Assuming that $\{e(t)\}$ is normally distributed white noise with mean zero and that the predictions are computed as conditional expectations as described in Section 4.4.1, it is easily found that

$$y^0(t+j|t) = -u_{0.99}\sigma_j ,$$

where $u_{0.99} = 2.3263$ is the 99 % quantile of the standardized normal distribution and σ_j^2 is the j -step-ahead prediction error variance. In a stochastic sense (4.54) is an AR(1) model, and it can easily be verified that (see, e.g., Madsen (1989))

$$\sigma_j^2 = \sigma_e^2 \sum_{i=0}^{j-1} a^{2i} . \quad (4.57)$$

In the present case $a = -0.8$ and $\sigma_e^2 = 0.04$. Hence the first 8 reference values become

j	1	2	3	4	5	6	7	8
$y^0(t+j t)$	-0.47	-0.60	-0.67	-0.71	-0.73	-0.75	-0.76	-0.76

Simulation of Suboptimal Minimum Variance Control

The four-step-ahead prediction of $y(t+4)$ can be expressed as

$$\begin{aligned} \hat{y}(t+4|t) &= 0.8\hat{y}(t+3|t) + q^4 B(q^{-1})u(t) + f(t+4) \\ &= 0.8\hat{y}(t+3|t) - u(t+1) - 3.1u(t) - 1.5u(t-1) \\ &\quad + f(t+4) , \end{aligned}$$

where $\hat{y}(t+3|t)$ is computed as usual (by successive one-step-ahead predictions). The minimal variance controller is found by solving the equation $\hat{y}(t+4|t) = y^0(t+4|t)$ with respect to $u(t+1)$, but in

this case this would lead to an unstable closed-loop system⁵. Instead the B -polynomial is replaced by $B'(q^{-1})$ from (4.55), and with this replacement the equation $\hat{y}(t+4|t) = y^0(t+4|t)$, i.e.

$$0.8\hat{y}(t+3|t) - 2.5u(t+1) - 2.5u(t) - 0.6u(t-1) + f(t+4) = y^0(t+4|t),$$

is solved with respect to $u(t+1)$. The solution leads to the suboptimal minimal variance controller:

$$u(t+1) = 0.4v'_4(t) + 0.284, \quad (4.58)$$

where

$$v'_4(t) = 0.8\hat{y}(t+3|t) - 2.5u(t) - 0.6u(t-1) + f(t+4).$$

The suboptimal minimal variance control of a system governed by the model in (4.54) has been simulated over 2000 sampling intervals. Figure 4.9 shows parts of the input and output sequences. The output is kept below zero most of the time (97.5% of the time). This percentage should be compared with the original aim, namely 99.0%. This lack of performance should be attributed to the fact that the controller is only suboptimal.

Due to the linearity of the controller and the sinusoidal function in the model, the input becomes a sine wave superposed by noise. The amplitude and the irregularities of the input signal indicate the activity of the controller.

⁵As in Section 4.2 the observations of input and output are assumed to be hourly averages. Therefore, given observations of input and output up to time t , the controller calculates the input, $u(t+1)$, which should be implemented between time t and time $t+1$.

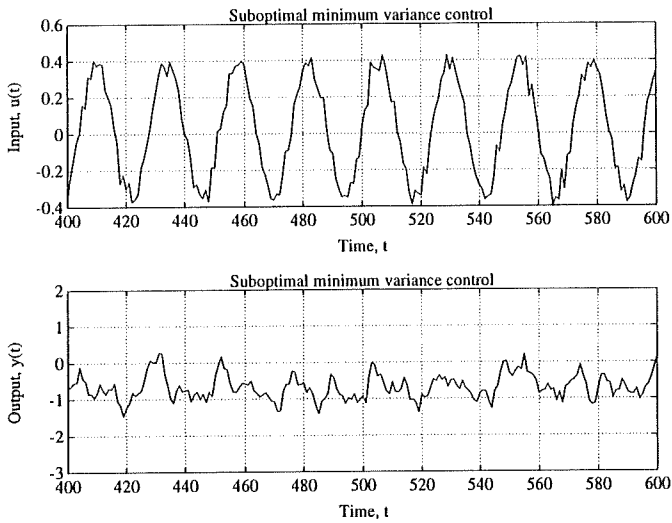


Figure 4.9: *Subsequences of the input and output signals obtained from the use of the suboptimal minimal variance control method.*

Simulation of XGPC

To carry out experiments with XGPC method, the controller in (4.49) is used. The following control parameters have been chosen:

$$\begin{aligned}
 N &= 8 \\
 \Gamma &= \text{diag}(\sigma_1^2, \sigma_2^2, \sigma_3^2, \sigma_4^2, \sigma_5^2, \sigma_6^2, \sigma_7^2, \sigma_8^2)^{-1} \\
 &= \text{diag}(25.00, 15.24, 12.20, 10.81, 10.08, 9.66, 9.41, 9.26) \\
 \Lambda &= \lambda I \\
 \omega &= \mathbf{0},
 \end{aligned}$$

where I is the 8×8 identity matrix and λ is a variable cost parameter. From the impulse response function, the H -matrix is obtained:

$$H = \begin{pmatrix}
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 -1.000 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 -3.900 & -1.000 & 0 & 0 & 0 & 0 & 0 & 0 \\
 -4.620 & -3.900 & -1.000 & 0 & 0 & 0 & 0 & 0 \\
 -3.696 & -4.620 & -3.900 & -1.000 & 0 & 0 & 0 & 0 \\
 -2.957 & -3.696 & -4.620 & -3.900 & -1.000 & 0 & 0 & 0
 \end{pmatrix}.$$

According to (4.49), the next control, $u(t+1)$, is found as the first element of

$$\mathbf{u}(t) = -[H^T \Gamma H + \lambda I]^{-1} H^T \Gamma [\mathbf{v}(t) - \mathbf{y}^0(t)],$$

where $\mathbf{y}^0(t) = (y^0(t+1|t), y^0(t+2|t), \dots, y^0(t+8|t))^T$. Observing that the j -step-ahead predictor of the model in (4.54) is

$$\begin{aligned}
 \hat{y}(t+j|t) &= 0.8y^*(t+j-1) - u(t+j-3) - 3.1u(t+j-4) \\
 &\quad - 1.5u(t+j-5) + f(t+j), \quad j = 1, 2, 3, \dots,
 \end{aligned}$$

Table 4.1: XGPC controllers for three different values of the costing parameter λ .

$\lambda =$	$u(t + 1) =$
1	$(0,0,0,0.1525,0.3109,-0.0997,0.0250,-0.0036)\mathbf{v}(t) + 0.2772$
100	$(0,0,0,0.0370,0.1072,0.0491,-0.0045,-0.0023)\mathbf{v}(t) + 0.1362$
1000	$(0,0,0,0.0077,0.0262,0.0235,0.0121,0.0069)\mathbf{v}(t) + 0.0567$

where

$$y^*(t + j - 1) = \begin{cases} y(t) & \text{if } j = 1 \\ \hat{y}(t + j - 1|t) & \text{if } j > 1 \end{cases},$$

the elements of $\mathbf{v}(t) = (v_1(t), v_2(t), \dots, v_8(t))^T$ are simply computed as

$$v_j(t) = \hat{y}(t + j|t)|_{u(t+i)=0, i=1, \dots, 8}, \quad j = 1, \dots, 8.$$

Table 4.1 shows the resulting controllers for three different values of λ . Figures 4.10, 4.11 and 4.12 show input and output signals for these three cases.

By comparing the input signals of Figures 4.9 and 4.10 it is seen that the XGPC controller with $\lambda = 1$ leads to more active control than the suboptimal minimal variance controller. However, the amplitudes of the input signals are almost the same. Concerning the output signals, it is obvious that the output of the SMVC controller is larger than zero more frequently than the output of the XGPC controller. Actually, the XGPC manages to keep the output below zero 99.05% of the time.

For $\lambda = 100$ it is seen (Figure 4.11) that the input sequence required by the XGPC controller has become more smooth. The output, on the other hand, is not significantly affected by the increase in the cost

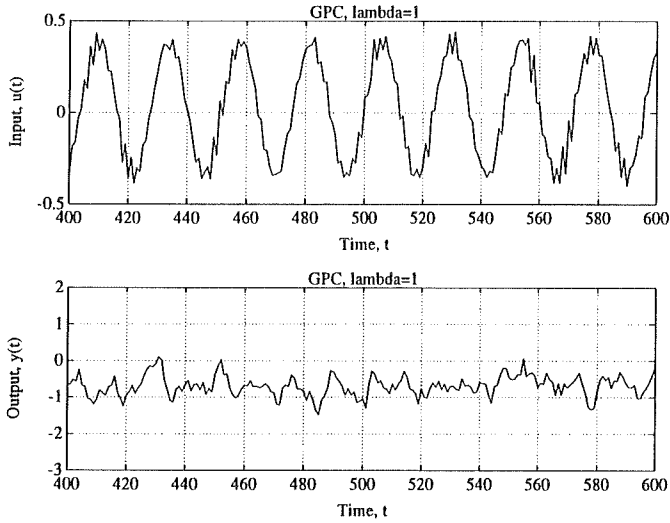


Figure 4.10: *Subsequences of the input and output signals obtained from the use of the XGPC with $\lambda = 1$.*

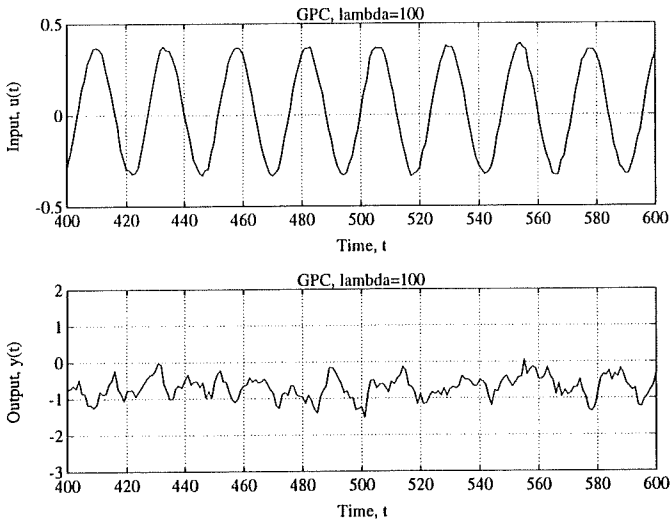


Figure 4.11: Subsequences of the input and output signals obtained from the use of the XGPC with $\lambda = 100$.

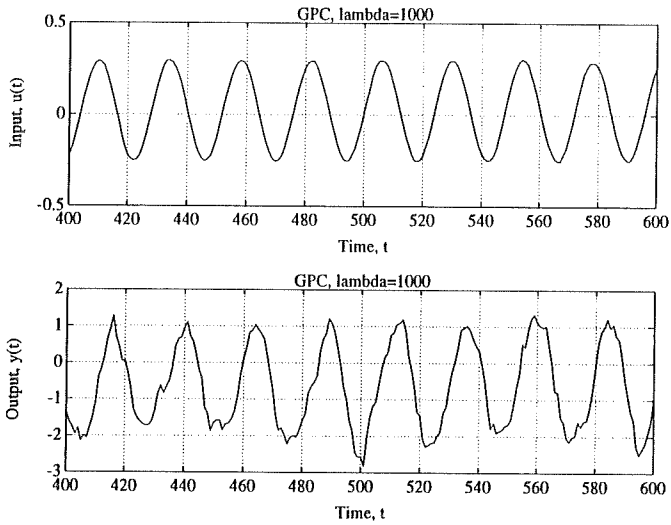


Figure 4.12: Subsequences of the input and output signals obtained from the use of the XGPC with $\lambda = 1000$.

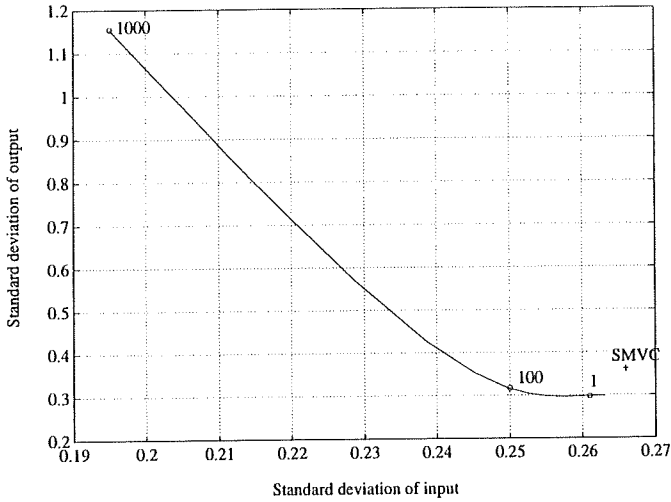


Figure 4.13: *Standard deviation of the output versus standard deviation of the input for the XGPC (solid curve) and the SMVC ('+'). Figures along the curve are values of λ for the XGPC.*

parameter from 1 to 100. 98.79% of the time the output is less than zero.

For $\lambda = 1000$ (Figure 4.12), the controller is no longer capable of controlling the output. The control cost has simply become too large.

The graph in Figure 4.13 shows the standard deviation of the output versus the standard deviation of the input for λ ranging from 0.1 to 1000. The curve seems to suggest $\lambda = 100$ as a reasonable choice of the cost parameter since the variation of input is reduced considerably, while the variation of the output is still close to its minimum. Compared to the SMVC, the superiority of the XGPC is obvious. In the case of $\lambda = 100$, the reductions of the input and output standard deviations are 6% and 12%, respectively.

4.6 Conclusion

Generalized predictive control (GPC) for both time-invariant and time-varying systems have been studied in this chapter. The GPC strategy is relevant to many different types of dynamic input-output systems. In this presentation, the application of the GPC method for the control of district heating systems was in focus, and the theoretical formulation of the controllers is, to some extent, seen from a district heating point of view.

Viewed as a production-consumption system, a district heating system constitutes a very complicated dynamic system involving maybe hundreds of different time-varying delays between the production plant and the various consumers. Thus, to meet the heat demand at any time, it is important to take the delays into account when controlling the supply temperature at the heat production plant. In order to do this, predictions of the heat demand are required. Therefore the heart of the problem is a predictive control problem.

At the combined heat and power plant, Vestkraft, in Esbjerg, a weighted predictive control (WPC) method is at present used to control the supply temperature. The WPC method can be used in case the dynamic relationship between input and output is not sufficiently described to permit the use of a true GPC method. Here insufficiency means that the model does not describe how the next control action affects the future output values. The purpose of the controller at Vestkraft is to keep the total flow of water close to, but below a critical maximum value. The resulting control input is computed by weighting predictions of the heat load 4, 5 and 6 hours ahead. Results from Esbjerg show that this control strategy leads to a significant lower supply temperature than previously used strategies.

In this chapter, the ordinary GPC (OGPC) method proposed by

Clarke *et al.* (1987A) was reviewed and a further generalized method (XGPC) allowing for embedded model structures was proposed. In brief, the GPC methods prescribe how to choose future input values in order to minimize future control errors and the control effort. The control error is the difference between the actual output and a reference value. At each sampling instant, the controller computes the optimal sequence of future input values (controls), but only the first value of this sequence is implemented.

Below it is briefly described how the XGPC method is generalized as compared to the OGPC method.

Model:

OGPC: A special time-invariant ARIMAX model is assumed: An ARMAX model with differenced input and output.

XGPC: A general time-varying ARMAX model which includes the OGPC ARIMAX model as a special case is assumed.

Cost function:

OGPC: Future control errors are weighted uniformly, and scalar weights are used to weight future control increments.

XGPC: Matrix weights for future control errors and the control values are introduced. By choosing special values of the matrix weights the OGPC cost function is obtained.

Control constraints:

OGPC: One of the design parameters of the OGPC method is a control horizon beyond which the control increments are zero.

XGPC: Future control values can be subject to any multi-dimensional linear equality constraint. This includes the zero-increment constraint as a special case.

Control of time-varying models:

OGPC: Time-varying parameters are not permitted since the multi-

step predictions are computed (recursively) by means of Diophantine equations.

XGPC: Multi-step predictions are achieved as successive one-step-ahead predictions. This approach allows for embedded models of time-varying parameters.

Notice that the great flexibility of the cost function and the control constraints of the XGPC method implies that the control signal beyond a certain control horizon can be constrained to follow any function of time which is a solution to a linear difference equation (a straight line, a parabola, a sine wave, an exponential etc.). Since the optimal supply temperature of a district heating system is known to be periodic with a period of 24 hours, the XGPC method could, e.g. be used to find the optimal sinusoidal variation of the future supply temperature.

The multi-step predictors of the ARMAX model are linear functions of the future control values. Furthermore, the cost function is a quadratic function of the future control errors and control values. Combining those two facts result in a cost function which is a quadratic function of the future control values, and the minimization problem is consequently a matter of solving a system of linear equations. This means that a matrix has to be inverted. However, if the parameters of the cost function and the model are time-invariant, the matrix inversion should be performed only once. If the time-variation of the parameters is periodic with period s , s different matrices must be inverted, and in the general time-varying case a matrix must be inverted at each sampling instant (e.g. in case of adaptive estimation of the parameters of embedded models).

Some simulation experiments with the XGPC method and suboptimal minimal variance control have been carried out in this chapter. The model and the controllers used for these experiments illustrate a simple model of the mass flow in a district heating system and as-

sociated controllers. Since the model is non-minimum phase, a true minimal variance controller would result in an unstable closed-loop system. A suboptimal minimal variance controller is used instead. The XGPC method can be applied directly without modifications. The results of the experiments show that the XGPC method is superior to the suboptimal minimal variance controller, and that parameters of the XGPC method can be chosen so that the variance of the input is reduced considerably without any significant increase in the variance of the output – i.e. the variance of the output is close to its minimum.

Chapter 5

Conclusions

IN both statistical modelling and control of dynamic systems it is often important and sometimes even necessary to consider embedded parameter variations. First of all the embedded physical parameters may be the key to the understanding of the underlying dynamic and stochastic mechanisms of the system. Secondly the consideration of the embedded parameters frequently leads to more natural parametrization of the models and makes the interpretability of the models easier. If the embedded parameters exhibit significant variation in time, this induces non-stationarity which must be taken into account when building models for prediction or control purposes.

All the models and methods studied in the thesis have some relation to identification, forecasting and control aspects of district heating systems. The distribution network of a district heating system together with the heat consumers represent a very complex dynamic and stochastic system. Due to the influence of the weather conditions and the diurnal rhythm of heating and hot tap water consumption, the total heat consumption and the transport times from the heating

plant to the consumers show very large variations during the day and the year. The variations of the heat consumption and the embedded time-delays represent very important factors when models for load forecasting and optimal control of the supply temperature are to be identified.

In **Chapter 2** methods for tracking variations of the time-delay and the dynamics of stochastic input-output systems are studied. The best results are obtained by using a recursive least squares technique for parallel estimation of a collection of models each representing one of the possible values of the real time-delay. This method is robust and easy to implement. Another method which is also based on the recursive least squares technique considers an approximative description of the embedded continuous time-delay by decomposing it into a discrete part and a continuous remainder. This method does only perform well for slowly varying time-delay. In case of fast variations, divergence tendencies are observed. This lack of robustness is most likely due to the actual approximation of the continuous time-delay, and a better approximation would probably improve the robustness.

Explicit modelling of the embedded parameter variations is an alternative to the methods based on recursive least squares estimation. A deterministic model of the embedded variation of the parameters in an ARX model can be obtained by letting the parameters depend on known functions of time. The diurnal variation of the time-delay in a district heating system can, for instance, be modelled by describing the coefficients of the B -polynomial by the first harmonics of a Fourier series. This leads to a linear regression model. The variation of the parameters and the time-delay of an ARX model can also be described by stochastic models. This results in a non-linear state-space model with the time-delay as one of its state variables. By means of an extended Kalman filter the state variables can be estimated recursively. A case study of district heating data shows that even with a simple model structure this method gives good results.

Chapter 3, which deals with dynamic models of the variations of ambient air temperature, falls into two parts. The first part describes exponential smoothing procedures which are empirical approaches for multi-step prediction (1 to 24 hours ahead). In the second part a physical approach based on linear stochastic models in continuous time is described.

As regards exponential smoothing procedures both linear and non-linear smoothing procedures are considered. The linear procedures correspond to embedded ARIMA model structures, and the resulting predictors are ARIMA predictors, although the parametrization is rather unusual. All the procedures include a description of the variation of the level and the seasonal (diurnal) profile in the air temperature data. Moreover, the non-linear procedures include a scaling factor for the amplitude of the seasonal profile, and experimental results show that this is a beneficial extension compared to the traditional linear procedures. The reason is that the scaling factor improves the adaptability to sudden changes in the diurnal amplitude of the air temperature which may occur when the weather changes from a cloudy to an unclouded period or vice versa. Generally the ability to detect such changes in the weather conditions seems to be crucial when forecasting ambient air temperature. This is one of the reasons why a new method which combines simple and seasonal exponential smoothing through an adaptively estimated regression model provides even better results than the above mentioned smoothing procedures.

The models in continuous time describe the relationship between the variations of the ambient air temperature and other climatic variables, primarily the net radiation. The parameters of the models are embedded quantities as thermal resistances and capacities which can be interpreted in physical terms. The models are formulated as linear state-space models and the parameters are identified directly using the maximum likelihood method. First a second order model which

alone utilizes the net radiation as an explanatory input is estimated, and next two third order models which utilize observations of the net radiation, the soil heat flux, the vapour pressure and the saturated vapour pressure are estimated. It turns out that the third order models which include embedded models of the heat balance at the surface of the earth provide the best description measured in terms of the one-step-ahead (one-hour-ahead) prediction error variance. The results indicate that the models give a good description of the thermodynamics of the upper layers of the soil (about 0.5 m) and the air up to about 75 m above the surface of the earth. The dynamics of the surrounding sea which is known to have influence on the air temperature over land is not reflected in the estimated time constants. The reason is that the time constant of the surrounding sea is almost two months, and such a large time constant is difficult to estimate with a sampling interval of only one hour when other time constants of the system is of same order of magnitude as the sampling time.

The topic of **Chapter 4** is multi-step predictive control of systems with embedded parameter variations. The common feature of multi-step predictive controllers is that the optimal control is obtained by a linear feedback from predicted values of the input-free control error. The optimal gain vector of this feedback can be determined by optimization of a suitable criterion as in generalized predictive control or it can be regarded as the actual design parameters as in weighted predictive control.

In district heating systems the presence of time-variation of both the time-delays and the heat consumption motivate heat load models with embedded parameter variations and multi-step predictive controllers which are able to utilize such models for optimal control of the plant supply temperature. Experiments with weighted predictive control of the plant supply temperature in a real district heating system have resulted in considerable lowering of the temperature and consequent savings. Furthermore, the frequency of critical situations

in connection with peak load in the mornings has been reduced significantly.

The ordinary generalized predictive control strategy is based upon a special ARIMAX model structure and a cost function that attributes costs to the differenced controls and the control errors. The multi-step predictions are obtained by considering the Diophantine equation, and therefore embedded variations of the model parameters are not allowed. However, an extended generalized predictive control strategy which allows embedded parameter variation is formulated. It is utilized that a general ARMAX model structure with time-varying parameters leads to multi-step predictions which can be formulated as a linear function of the next control actions. In the cost function more flexibility is obtained by assuming general cost matrices instead of scalar cost coefficients, and as an extension of the "cost horizon" idea in ordinary generalized predictive control, a general matrix equality constraint for the control sequence is introduced. Actually the ordinary generalized predictive control strategy is a special case of the new extended strategy with respect to both the model structure and the cost function.

Through simulation experiments excellent results are obtained when using extended generalized predictive control for control of a non-minimum phase system. Compared with suboptimal minimal variance control the extended control strategy leads to significant reduction of the variance of both the control signal and the output signal, i.e. less control effort and smaller control errors are obtained at the same time.

Appendix A

The Kalman Filter in Discrete Time

The Kalman filter in discrete time is mentioned in Chapters 2 and 3. Therefore, the following description is included as an appendix for reference purposes.

Consider the time-invariant linear, stochastic state-space model in discrete time:

$$\mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \mathbf{e}_{1,t} \quad (\text{A.1})$$

$$\mathbf{y}_t = \mathbf{C}\mathbf{x}_t + \mathbf{e}_{2,t}, \quad (\text{A.2})$$

where

\mathbf{x}_t is the state vector, which is not measurable (dimension $n \times 1$),

\mathbf{y}_t is the observation vector, which is measurable (dimension $m \times 1$),

A, C are the system matrices (dimensions $n \times n$ and $m \times n$), and $\{e_{1,t}\}, \{e_{2,t}\}$ are n and m dimensional white noise processes with

$$E[e_{1,t}] = \mathbf{0} \quad (\text{A.3})$$

$$E[e_{2,t}] = \mathbf{0} \quad (\text{A.4})$$

$$\text{Var} \left[\begin{pmatrix} e_{1,t} \\ e_{2,t} \end{pmatrix} \right] = \begin{pmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \Sigma_2 \end{pmatrix}. \quad (\text{A.5})$$

The Equations (A.1) and (A.2) are called the system equation and the observation equation, respectively.

Given the model in (A.1) and (A.2), the optimal linear prediction, $\hat{\mathbf{x}}_{t|t-1}$, and reconstruction, $\hat{\mathbf{x}}_{t|t}$, of \mathbf{x}_t can be obtained by means of the ordinary Kalman filter in discrete time (Abraham and Ledolter (1983), Madsen (1989)),

$$\hat{\mathbf{x}}_{t|t-1} = A\hat{\mathbf{x}}_{t-1|t-1} \quad (\text{A.6})$$

$$P_{t|t-1} = AP_{t-1|t-1}A^T + \Sigma_1 \quad (\text{A.7})$$

$$\hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + K_t(\mathbf{y}_t - C\hat{\mathbf{x}}_{t|t-1}) \quad (\text{A.8})$$

$$P_{t|t} = P_{t|t-1} - K_tCP_{t|t-1} \quad (\text{A.9})$$

$$K_t = P_{t|t-1}C^T(CP_{t|t-1}C^T + \Sigma_2)^{-1}. \quad (\text{A.10})$$

By $\hat{\mathbf{x}}_{t|t-1}$, e.g., we denote the conditional expectation of \mathbf{x}_t , conditioned on the observations $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{t-1}$.

Equations (A.6) and (A.7) are prediction equations that specify the one-step-ahead prediction of the state vector and its covariance matrix. In (A.8) and (A.9) (the reconstruction equations) the estimate and the covariance matrix of the state vector are updated after a new observation, y_t , has become available. The vector K_t is a weight (the Kalman gain) determining how much the most recent one-step-ahead prediction error should influence the updating steps. Before

the first iteration, the Kalman filter has to be supplied with the prior distribution of \mathbf{x}_0 ; i.e. $\hat{\mathbf{x}}_{0|0}$ and $\mathbf{P}_{0|0}$.

The results can easily be generalized to allow time-varying systems.

Appendix B

The Computer Program PRESS

Some of the models and methods described in the thesis are implemented in the computer program PRESS. Therefore, a brief description of PRESS and practical experience from using PRESS at the combined heat and power plant Vestkraft in Esbjerg are given here (see also Madsen *et al.* (1992)).

B.1 PRESS

The name PRESS is an abbreviation for “Prognose- og Energisty-ringssystem” (Prediction and Energy Control System). The software is a tool for automatic data collection, supervision, forecasting and control in district heating systems and was developed by IMSOR in co-operation with the power plant company, I/S Vestkraft, in Es-

bjerg and the municipal distribution services in Esbjerg and Varde (Forsyningsvirksomhederne i Esbjerg Kommune og Varde Kommunale Værker). The development costs in connection with the first version of PRESS were financed by I/S Vestkraft, the municipal distribution services in Esbjerg and Varde. Until now PRESS has been installed at I/S Vestkraft in Esbjerg and in a modified version at the district heating plant in Sønderborg (Sønderborg Fjernvarme a.m.b.a.). PRESS has a menu-driven user interface and runs on IBM compatible personal computers. A new version of PRESS based upon a modern graphic user interface environment (Motif under X11 and UNIX) is being developed.

The idea behind PRESS is to minimize the costs in connection with the production and distribution of heat. For the district heating system in Esbjerg/Varde (and several other systems in Denmark as well) lower production costs and less heat loss in the distribution network can be obtained by minimization of the supply temperature. It has been estimated that the production costs are reduced by about 1% and the heat loss by about 0.5% per °C the plant supply temperature is decreased. Therefore, in consideration of certain restrictions, PRESS minimizes the supply temperature of the water leaving the combined heat and power plant. Two restrictions are imposed when the currently acceptable minimum of the supply temperature is computed:

- 1) It must be possible to meet the total heat demand at any time (taking a physical maximum limit of the flow in the distribution network into account).
- 2) The supply temperature in the service pipes at the consumers must be above a certain minimum which depends on the ambient air temperature at the time of delivering (see Figure B.1).

In order to comply with the first restriction, PRESS computes forecasts of the heat demand and uses these forecasts in a weighted pre-

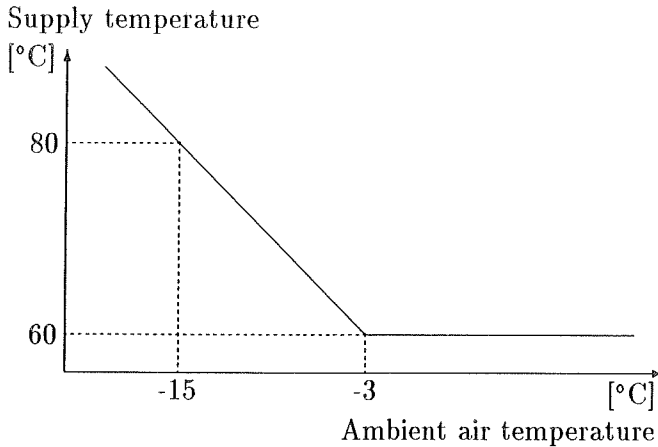


Figure B.1: *Temperature demand at the consumers.*

dictive control strategy (see the discussion in Section 4.1 and the prediction models and control strategy in Section 4.2). This gives one lower limit of the plant supply temperature. The second restriction is fulfilled by using minimal variance controllers for each of five representative points in the distribution network (see the discussion in Section 2.1). Thereby another lower limit of the plant supply temperature is obtained. Both of the two lower limits of the plant supply temperature must be observed and therefore the highest of them is used. The process of calculating forecasts and the optimal supply temperature is repeated once each hour.

Figure 1.2 in Chapter 1 (page 6) shows a simplified diagram of the data flow and the data processing taking place in PRESS. Measurements of supply temperatures, return temperature, heat production, ambient air temperature and wind speed are collected automatically, and hourly averages are computed (based on 12 instantaneous measurements per hour). Every hour forecasts of the heat demand up to 12 hours ahead (in Esbjerg, 18 hours ahead in Sønderborg) and forecasts of the ambient air temperature up to 20 hours ahead are

computed. These forecasts together with other measurements are used to compute the optimal supply temperature for the following hour as briefly described above.

B.2 Results from Using PRESS in Esbjerg

Experiences from use of PRESS in Esbjerg are reported by Madsen *et al.* (1992). In this report hourly measurements of the plant supply temperature and the ambient air temperature from the period January - March 1987 are compared with corresponding data from the period February - April 1991. The former period was before and the latter period after PRESS was brought into operation at Vestkraft. For each of the periods an estimation of a regression line for the relationship between the ambient air temperature and the supply temperature have been carried out. The result of this regression analysis is shown in Figure B.2. In order to estimate the average difference between the supply temperature in the two periods, the regression lines have been restricted to have the same slope. The regression analysis is based on observations for which the ambient air temperature is within the interval from $-1.5\text{ }^{\circ}\text{C}$ and $7\text{ }^{\circ}\text{C}$. This is the overlapping interval of the two periods with respect to the ambient air temperature. It appears from the figure that the supply temperature was about $9\text{ }^{\circ}\text{C}$ lower in 1991 than in 1987. It also appears that on average the supply temperature in 1991 was even lower than the minimum control curve used at Vestkraft in 1987 (this curve is shown in full in Figure 1.1 on page 4). The conclusion from this is that the use of PRESS imply a significant reduction of the supply temperature in typical winter periods, and this in turn results in reduction of the heat production costs and the heat loss from the district heating pipes. Although the control strategy leads to lower average supply temperature, the results from Esbjerg also show that

during few hours in the night-time PRESS requires higher supply temperature than the previously used control strategy. The reason is that due to the transport time the temperature is increased in due time before the heat demand culminates between 7 a.m. and 8 a.m. (at workdays). In this way PRESS eliminates most of the critical situations which occurred in winter mornings under the previously used control strategy.

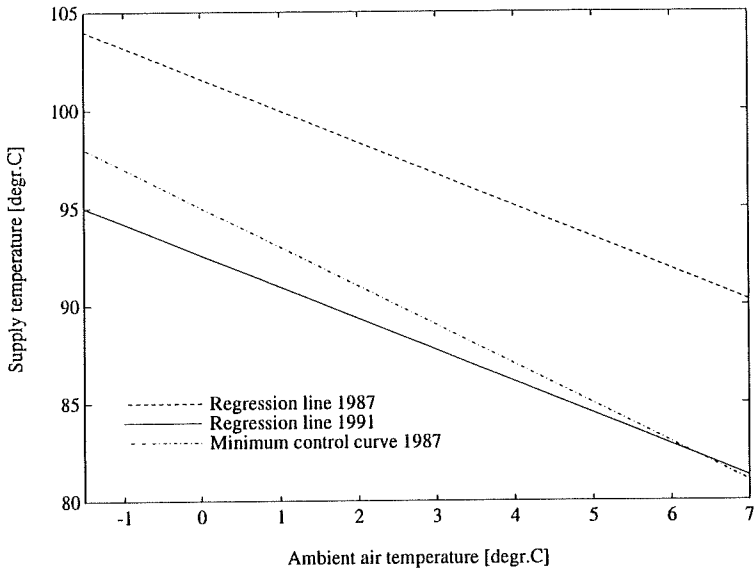


Figure B.2: *Estimated regression lines for the relationship between the ambient air temperature and the supply temperature based on data from 1987 (before PRESS) and 1991 (with PRESS).*

References

Abraham, B. and J. Ledolter, 1983: *Statistical Methods for Forecasting*. Wiley & Sons.

Åström, K. J., 1970: *Introduction to Stochastic Control Theory*. Academic Press.

Åström, K. J. and B. Wittenmark, 1989: *Adaptive Control*. Addison-Wesley.

Bányász, Cs. and L. Keviczky, 1988: A New Recursive Time Delay Estimation Method for ARMAX Models. *8th IFAC/IFORS Symposium on Identification and System Parameter Estimation*, Pergamon Press, 1452-1457.

Benonysson, A., 1991: *Dynamic Modelling and Operational Optimization of District Heating Systems*. Ph.D. Thesis, The Laboratory of Heating and Air Conditioning, The Technical University of Denmark, Lyngby, Denmark.

Berkowicz, R. and L. P. Prahm, 1982: Sensible Heat Flux Estimated from Routine Meteorological Data by the Resistance Method. *J. Appl. Meteor.*, **21**, 1845-1864.

Bitmead, R. R., M. Gevers and V. Wertz, 1990: *Adaptive Optimal Control, The Thinking Man's GPC*. Prentice-Hall.

- Bjerre, T.**, 1992: *Generalized Predictive Control of Energy Systems*. (In Danish: *Generel prædiktiv kontrol af energisystemer*). Masters Thesis No. 6/92, IMSOR, The Technical University of Denmark, Lyngby, Denmark.
- Box, G. E. P.** and **G. M. Jenkins**, 1976: *Time Series Analysis, Forecasting and Control*. San Francisco: Holden-Day.
- Brown, R. G.**, 1963: *Smoothing, Forecasting and Prediction of Discrete Time Series*. Prentice-Hall, Englewood Cliffs, NJ.
- Chen, H.-F.** and **J.-F. Zhang**, 1990: Identification and Adaptive Control for Systems with Unknown Orders, Delay and Coefficients. *IEEE Transactions on Automatic Control*, **35**, No. 8, 866-877.
- Clarke, D. W.**, **C. Mohtadi** and **P. S. Tuffs**, 1987: Generalized Predictive Control – Part I. The Basic Algorithm. *Automatica*, **23**, 137-148.
- Clarke, D. W.**, **C. Mohtadi** and **P. S. Tuffs**, 1987: Generalized Predictive Control – Part II. Extensions and Interpretations. *Automatica*, **23**, 149-160.
- Davis, M. H. A.** and **R. B. Vinter**, 1985: *Stochastic Modelling and Control; Monographs on Statistics and Applied Probability*. Chapman and Hall, London.
- Edlund, P. O.**, 1984: Identification of the Multi-input Box-Jenkins Transfer Function Model. *Journal of Forecasting*, **3**, 297-308.
- Edlund, P. O.** and **H. T. Søgaaard**, 1993: Fixed versus Time-varying Transfer Functions for Modelling Business Cycles. *Journal of Forecasting*, **12**, 345-364.
- Elnaggar, A.**, **G. A. Dumont** and **A.-L. Elshafei**, 1989: Recursive Estimation for System of Unknown delay. *Proceedings of the 28th Conference on Decision and Control*, 1809-1810.

Hansen, S., S. E. Jensen and H. C. Aslyng, 1981: *Jordbrugsmeteorologiske Observationer, Statistik Analyse og Vurdering, 1955 - 1979*. Hydrological Laboratory, Royal Danish Veterinary and Agricultural University, Denmark.

Harvey, A. C., 1989: *Forecasting, Structural Time Series Models and the Kalman Filter*. Cambridge University Press.

Holst, J., 1977: *Adaptive Prediction and Recursive Estimation*. Ph.D. Thesis TFRT-1013, Department of Automatic Control, Lund Institute of Technology, Lund, Sweden.

Holst, U., G. Lindgren and J. Holst, 1992: Recursive Estimation in Switching Autoregressions with Markov Regime. Report TFMS-3084, Department of Mathematical Statistics, Lund Institute of Technology, Lund, Sweden.

Holt, C. C., 1957: *Forecasting Trends and Seasonals by Exponentially Weighted Moving Average*. O. N. R. Memorandum, No. 52, Carnegie Institute of Technology.

Houghton, J. T., 1977: *The Physics of Atmospheres*. Cambridge University Press.

Kabaila, P. V., 1981: Estimation Based on One Step Ahead Prediction Versus Estimation Based on Multi-Step Ahead Prediction. *Stochastics*, **6**, 43-55.

Kurz, H. and W. Goedecke, 1981: Digital Parameter-Adaptive Control of Processes with Unknown Dead Time. *Automatica*, **17**, 245-252.

Ljung, L., 1987: *System Identification: Theory for the User*. Prentice-Hall, Englewood Cliffs, NJ.

Ljung, L. and Söderström, 1983: *Theory and Practice of Recursive Identification*. MIT Press, Cambridge, Massachusetts.

Madsen, H., 1985: *Statistically Determined Dynamical Models for Climate Processes*. Ph.D. Thesis No. 45, IMSOR, The Technical University of Denmark, Lyngby, Denmark.

Madsen, H., 1989: *Time Series Analysis*. (In Danish: *Tidsrækkeanalyse*). IMSOR, The Technical University of Denmark, Lyngby, Denmark.

Madsen, H. and H. Melgaard, 1991: The Mathematical and Numerical Methods Used in CTLSM – a program for ML-estimation in stochastic, continuous time dynamical models. Research Report No. 7/1991, IMSOR, The Technical University of Denmark, Lyngby, Denmark.

Madsen, H., P. Thyregod and J. Holst, 1987: A Continuous Time Model for the Variations of Air Temperature. *Preprints of the Tenth Conference on Probability and Statistics in Atmospheric Sciences*, 52-58.

Madsen, H., O. P. Palsson, K. Sejling and H. Søgaaard, 1990: *Models and Methods for Optimization of District Heating Systems, Part I: Models and Identification Methods*. EFP 1323/89-14, The Danish Energy Research Program.

Madsen, H., O. P. Palsson, K. Sejling and H. Søgaaard, 1992: *Models and Methods for Optimization of District Heating Systems, Part II: Models and Control Methods*. EFP 1323/89-14, The Danish Energy Research Program.

Makridakis, S., A. Andersen, R. Carbone, R. Fildes, M. Hibon, R. Lewandowski, J. Newton, E. Parzen and R. Winkler, 1984: *The Forecasting Accuracy of Major Time Series Methods*. Wiley & Sons.

Martin, R. D. and V. J. Yohai, 1985: Robustness in Time Series and Estimating ARMA Models. *Handbook of Statistics*, 5, 119-155.

Melgaard, H. and H. Madsen, 1992: Methods for Parameter Estimation in Stochastic Differential Equations. *Proceedings of the 1st Workshop on Stochastic Numerics*, 53-64.

Melgaard, H. and H. Madsen, 1993: CTLSM, Continuous Time Linear Stochastic Modeling, Version 2.6. Technical Report No. 1/1993, IMSOR, The Technical University of Denmark, Lyngby, Denmark.

Newbold, P. and C. W. J. Granger, 1974: Experience with Forecasting Univariate Time Series and the Combination of Forecasts. *J. of the Royal Stat. Soc., Series A*, **137**, 131-165.

Ng, C. N. and P. C. Young, 1990: Recursive Estimation and Forecasting of Non-stationary Time Series. *Journal of Forecasting*, **9**, 173-204.

OECD, 1987: *OECD Leading Indicators and Business Cycles in Member Countries 1960-1985*. OECD Department of Economics and Statistics, Main Economic Indicators, Source and Methods, No. 39, January 1987.

Parkum, J. E., 1992: *Recursive Identification of Time-Varying Systems*. Ph.D. Thesis No. 57, IMSOR, The Technical University of Denmark, Lyngby, Denmark.

Poulsen, N. K., 1985: *Robust Self-tuning Controllers*. Ph.D. Thesis No. 44, IMSOR, The Technical University of Denmark, Lyngby, Denmark.

Poulsen, N. K. and J. Holst, 1988: Simultaneous Estimation of Innovation Variance and States in a Dynamic System. *8th IFAC/IFORS Symposium on Identification and System Parameter Estimation*, Pergamon Press, 1263-1269.

Sejling, K., 1993: *Modelling and Prediction of Load in District Heating Systems*. Ph.D. Thesis No. 65, IMSOR, The Technical University of Denmark, Lyngby, Denmark.

Shumway, R. H. and **D. S. Stoffer**, 1982: An Approach to Time Series Smoothing and Forecasting Using the EM Algorithm. *Journal of Time Series Analysis*, **3**, 253-264.

Shumway, R. H., 1988: *Applied Statistical Time Series Analysis*. Prentice-Hall.

Stoica, P. and **A. Nehorai**, 1989: On Multistep Prediction Error Methods for Time Series Models. *J. of Forecasting*, **8**, 1989.

Söderström, T. and **P. Stoica**, 1989: *System Identification*. Prentice-Hall.

Søgaard, H. T., 1988: *Identification and Adaptive Control of District Heating Systems*. (In Danish: *Identifikation og adaptiv regulering af fjernvarmesystemer*). Masters Thesis No. 10/88, IMSOR, The Technical University of Denmark, Lyngby, Denmark.

Søgaard, H. T. and **Madsen, H.**, 1989: Methods for Tracking Time Varying Delays in Dynamic Systems. Preprint No. 3/89, IMSOR, The Technical University of Denmark, Lyngby, Denmark.

Winters, P. R., 1960: Forecasting Sales by Exponentially Weighted Moving Averages. *Manage. Sci.*, **6**, 324-342.

Young, P., 1984: *Recursive Estimation and Time-Series Analysis. An Introduction*. Springer-Verlag, Berlin, Heidelberg.

IMSOR Ph.D. Theses

1. **Sigvaldason, Helgi.** (1963). *Beslutningsproblemer ved et hydrotermisk elforsyningssystem.* 92 pp.
2. **Nygaard, Jørgen.** (1966). *Behandling af et dimensioneringsproblem i telefonien* 157 pp.
3. **Krarpup, Jakob.** (1967). *Fixed-cost and other network flow problems as related to plant location and to the design of transportation and computer systems.* 159 pp.
4. **Hansen, Niels Herman.** (1967). *Problemer ved forudsigelse af lyd hastighed i danske farvande. Analyse af et stokastisk system. Del 1: Tekst. Del 2: Figurer og tabeller.* 104 pp. + 95 pp.
5. **Larsen, Mogens E.** (1968). *Statistisk analyse af elementære kybernetiske systemer.* 210 pp.
6. **Punhani, Amrit Lal.** (1968). *Decision problems in connection with atomic power plants.* 133 pp.
7. **Clausen, Svend.** (1969). *Kybernetik, systemer og modeller.* 205 pp.

8. **Vidal, R.V. Valqui.** (1970.) *Operations research in production planning. Interconnections between production and demand. Volume 1-2.* 321 pp.
9. **Bilde, Ole.** (1970). *Nonlinear and discrete programming in transportation, location and road design. Volumes 1-2.* 291 pp.
10. **Rasmusen, Hans Jørgen.** (1972). *En decentraliseret planlægningsmodel.* 185 pp.
11. **Dyrberg, Christian.** (1973). *Tilbudsgivning i en entreprenør virksomhed.* 158 pp.
12. **Madsen, Oli B.G.** (1973). *Dekomposition og matematisk programmering.* 271 pp.
13. **Dahlgard, Peter.** (1973). *Statistical aspects of tide prediction. Volume 1. Volume 2: Figures and tables.* 202 pp. + 170 pp.
14. **Spliid, Henrik.** (1973). *En statistisk model for stormflodsvarsling.* 205 pp.
15. **Pinochet, Mario.** (1973). *Operations research in strategic transportation planning. The decision process in a multiharbour system.* 374 pp.
16. **Christensen, Torben.** (1973). *Om semi-markov processer. Udvidelser og anvendelser inden for den sociale sektor.* 239 pp.
17. **Jacobsen, Søren Kruse.** (1973). *Om lokaliseringsproblemer, modeller og løsninger.* 355 pp.
18. **Marqvardsen, Hans.** (1973). *Skemalægning ved numerisk simulation.* 222 pp.
19. **Mortensen, Jens Hald.** (1974). *Interregionale godstransporter. Teoridannelser og modeller.* 223 pp.

20. **Severin, Juan Melo.** (1974). *Introduction to operations research in systems synthesis. A chemical process design synthesis application.* 249 pp.
21. **Spliid, Iben & Uffe Bundgaard-Jørgensen.** (1974). *Skitse*
22. *til en procedure for kommunalplanlægning.* 544 pp.
23. **Mosgaard, Christian.** (1975). *International planning in disaster situations.* 187 pp.
24. **Holm, Jan.** (1975). *En optimeringsmodel for kollektiv trafik.* 246 pp.
25. **Jesson, Pall.** (1975). *Stokastisk programmering. Del 1: Modeller. Del 2: Metodologiske overvejelser og anvendelser.* 333 pp.
26. **Iversen, Villy Bæk.** (1976). *On the accuracy in measurements of time intervals and traffic intensities with application to teletraffic and simulation.* 202 pp.
27. **Drud Arne.** (1976). *Methods for control of complex dynamic systems. Illustrated by econometric models.* 209 pp.
28. **Togsverd, Tom.** (1976). *Koordinering af kommunernes ressourceforbrug.* 295 pp.
29. **Jensen, Olav Holst.** (1976). *Om planlægning af kollektiv trafik. Operationsanalytiske modeller og løsningsmetoder.* 321 pp.
30. **Beyer, Jan E.** (1976). *Ecosystems. An operational research approach.* 315 pp.
31. **Bille, Thomas Bastholm.** (1977). *Vurdering af Egnsudviklingsprojekter. Samspil mellem benefit-cost analyse og den politiske vurdering i en tid under forandring.* 260 pp.

32. **Holst, Erik.** (1979). *En statistisk undersøgelse af tabletsierier.* 316 pp.
33. **Aagaard-Svendsen, Rolf.** (1979). *Econometric methods and Kalman filtering.* 300 pp.
34. **Hansen, Steen.** (1979). *Project control by quantitative methods.* 230 pp.
35. **Scheufens, Ernst Edvard.** (1980). *Statistisk analyse og kontrol af tidsafhængige vandkvalitetsdata.* 152 pp.
36. **Lyngvig, Jytte.** (1981). *Samfundsøkonomisk planlægning.* 252 pp.
37. **Troelsgård, Birgitte.** (1981). *Statistisk bestemmelse af modeller for rumlufttemperatur.* 213 pp.
38. **Raft, Ole.** (1981). *Delivery planning by modular algorithms.* 220 pp.
39. **Jensen, Sigrid M.** (1981). *Analyse af interregionale togrejser. + Figurer og appendices.* 212 pp. + 174 pp.
40. **Ravn, Hans.** (1982). *Technology and underdevelopment. The case of Mexico.* 376 pp.
41. **Hansen, Sten.** (1983). *Phase-type distributions in queueing theory.* 209 pp.
42. **Ferreira, Jose A.S.** (1984). *Optimal control of discrete-time systems with applications.* 252 pp.
43. **Behrens, Jens Christian.** (1985). *Mathematical modelling of aquatic ecosystems applied to biological waste water treatment + Appendix 1-2.* 32 pp. + 389 pp. + 180 pp.

-
44. **Poulsen, Niels Kjølstad.** (1985). *Robust self tuning controllers.* 240 pp.
 45. **Madsen, Henrik.** (1985). *Statistically determined dynamic models for climate processes. Part 1-2.* 428 pp.
 46. **Sørensen, Bo.** (1986). *Interactive distribution planning.* 253 pp.
 47. **Lethan, Helge B.** (1986). *Løsning af store kombinatoriske problemer.* 173 pp
 48. **Boelskifte, Søren.** (1988). *Dispersion and current measurements. An investigation based on time series analysis and turbulence models.* Risø-M-2566. 154 pp.
 49. **Nielsen, Bo Friis.** (1988). *Modelling of multiple access systems with phase type distributions.* 253 pp.
 50. **Christensen, John M.** (1988). *Project planning and analysis. Methods for assessment of rural energy projects in deveoping countries.* Risø-M-2706. 158 pp.
 51. **Olsen, Klaus Juel.** (1988). *Texture analysis of ultrasound images of livers.* 162 pp.
 52. **Holst, Helle.** (1988). *Statistisk behandling af nærinfrarøde reflektionsmålinger.* 309 pp. + app.
 53. **Knudsen, Torben.** (1989). *Start/stop strategier for vind-diesel systemer.* 275 pp.
 54. **Ersbøll, Bjarne Kjær.** (1989). *Transformations and classifications of remotely sensed data. Theory and geological cases.* 297 pp.
 55. **Kragh, Anders Laage.** (1990). *Kø-netværksmodeller til analyse af FMS anlæg.* 205 pp.

56. Hansen, Christian Kornerup. (1991). *Statistical methods in the analysis of repairable systems reliability*. 56 pp.
57. Parkum, Jens Ejnar. (1992). *Recursive identification of time-varying systems*. 206 pp.
58. Bilbo, Carl M. (1992). *Statistical analysis of multivariate degradation models*. 167 pp.
59. Carstensen, Jens Michael. (1992). *Description and simulation of visual texture*. 234 pp.
60. Halse, Karsten. (1992). *Modeling and solving complex vehicle routing problems*. 372 pp.
61. Hendricks, Elbert. (1992). *Identification and estimation of nonlinear systems using physical modelling*. 273 pp.
62. Windfeld, Kristian. (1992). *Application of computer intensive data analysis methods to the analysis of digital images and spatial data*. 190 pp.
63. Iwersen, Jørgen. (1992). *Statistical control charts : Performance of Shewhart and Cusum charts*. 326 pp.
64. Olsson, Carsten Kruse. (1993). *Image processing methods in materials science*. 274 pp.
65. Sejling, Ken. (1993). *Modelling and prediction of load in heating systems*. 283 pp.
66. Søggaard, Henning T. (1993). *Stochastic systems with embedded parameter variations - applications to district heating*. 280 pp.