

Sparse estimation: applications in atrial fibrillation

Citation for published version (APA):

Zeemering, S. (2015). *Sparse estimation: applications in atrial fibrillation*. [Doctoral Thesis, Maastricht University]. Datawyse / Universitaire Pers Maastricht. <https://doi.org/10.26481/dis.20151126sz>

Document status and date:

Published: 01/01/2015

DOI:

[10.26481/dis.20151126sz](https://doi.org/10.26481/dis.20151126sz)

Document Version:

Publisher's PDF, also known as Version of record

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.umlib.nl/taverne-license

Take down policy

If you believe that this document breaches copyright please contact us at:

repository@maastrichtuniversity.nl

providing details and we will investigate your claim.

SPARSE ESTIMATION
APPLICATIONS IN ATRIAL FIBRILLATION

Copyright ©Stef Zeemering, Maastricht, The Netherlands, 2015

All rights reserved. No part of this book may be reproduced, stored in a database or retrieval system, or transmitted in any form or by any means, without the written permission of the author.

Printed by Datawyse, Universitaire Pers Maastricht

SPARSE ESTIMATION
APPLICATIONS IN ATRIAL FIBRILLATION

PROEFSCHRIFT

ter verkrijging van de graad van doctor aan de Universiteit
Maastricht, op gezag van de Rector Magnificus, Prof. dr. L.L.G. Soete
volgens het besluit van het College van Decanen, in het openbaar te
verdedigen op donderdag 26 november 2015 om 10.00 uur

door

Stef Zeemering
geboren op 7 maart 1979
te Oldenzaal, Nederland

Promotores

Prof. dr. ir. R.L.M. Peeters

Prof. dr. U. Schotten

Copromotor

Dr. R.L. Westra

Beoordelingscommissie

Prof. dr. T. Delhaas (voorzitter)

Dr. J.M.H. Karel

Prof. dr. P. Martens

Prof. dr. P.G. Platonov (Lund University, Zweden)

Prof. dr. S. Weiland (Technische Universiteit Eindhoven)

The research presented in this thesis was supported in part by a grant from the European Union (FP7 Collaborative project EUTRAF, 261057) and by a grant from the Center of Translational Molecular Medicine (CTMM – COHFAR).

Contents

1	General introduction	1
1.1	Background	1
1.2	Thesis outline	2
	Sparse estimation	3
2	A framework for sparse estimation	5
2.1	Introduction	5
2.2	Maximizing sparsity	6
2.2.1	Minimizing the ℓ_1 -norm of the parameter vector θ	9
2.2.2	Relationship between primal and dual problem	13
2.2.3	Practical issues	16
2.3	Mixed optimization	17
2.3.1	Technical description	18
2.3.2	Practical issues	19
3	Application: Linear regression models	21
3.1	Introduction	21
3.2	The model class	22
3.3	Experiments	25
3.3.1	Experimental setup	25
3.3.2	Underdetermined problems	26
3.3.3	Performance of the least absolute shrinkage and selection operator (LASSO)	43
3.3.4	Overdetermined problems and subset selection problems	44
3.4	Conclusions	46
4	Application: State-space models	47
4.1	Introduction	47
4.2	The model class	48
4.3	Prediction error identification	49
4.4	Parameterization and identifiability issues	50

4.5	Sparse state-space estimation	51
4.5.1	Jacobian and Hessian	52
4.5.2	Relation to DDLC	56
4.5.3	Example: small non-sparse models and full parameterization	57
4.6	Discrete-time network models	60
4.6.1	Experiments	61
4.6.2	Results	64
4.7	Discretized continuous-time network models	65
4.7.1	Experiments	69
4.7.2	Results	70
4.8	Continuous-time network models	72
4.8.1	Experiments	75
4.8.2	Results	76
4.9	Conclusions	80
Applications in atrial fibrillation		83
5	Introduction to atrial fibrillation	85
5.1	Definition and treatment of atrial fibrillation	85
5.2	Quantitive description of atrial fibrillation	86
5.2.1	Invasive analysis	88
5.2.2	Noninvasive analysis	88
6	Electrogram analysis and wave reconstruction	91
6.1	Introduction	91
6.2	Methods	92
6.2.1	Data acquisition	92
6.2.2	Electrogram pre-processing	92
6.2.3	Intrinsic deflection detection	92
6.2.4	Fibrillation wave construction	94
6.2.5	Validation	96
6.3	Results and discussion	96
6.4	Conclusion	98
7	Identification of recurring wavefront patterns	99
7.1	Introduction	99
7.2	Methods	100
7.2.1	Sparse multivariate autoregression model	100
7.2.2	Dominant pathway identification	101
7.2.3	High-density contact mapping	103
7.3	Results	103
7.4	Discussion and Conclusions	105

8	Noninvasive prediction of cardioversion outcome	107
8.1	Introduction	107
8.2	Methods	108
8.3	Results	112
8.4	Risk of progression to persistent AF	117
8.5	Discussion	121
8.A	Parameter selection via elastic net logistic regression	123
8.B	Effect of time interval between echo and CV	126
9	General discussion	129
9.1	Sparse estimation	129
9.1.1	Sparse linear regression	129
9.1.2	Sparse state-space identification	130
9.1.3	Sparse network identification	130
9.2	Applications in atrial fibrillation	132
9.2.1	Recurrent propagation pattern identification	132
9.2.2	Feature selection to predict pharmacological cardioversion . . .	132
	Bibliography	135
	Summary	143
	Samenvatting (Summary in Dutch)	145
	Valorization	149
	Dankwoord / Acknowledgements	151
	List of publications	155
	Curriculum vitae	157
	Acronyms	159
	Notation index	161
	Index	163

Chapter 1

General introduction

1.1 Background

For many practical purposes it is useful to develop mathematical models which describe the essence of the system, the process, or the phenomenon under consideration, in a simplified, yet sufficiently accurate way. Within the area of system identification, the methodology is focused on the estimation of a model from a selected model class followed by a model validation step to establish the feasibility of the result. Restricting the discussion to the class of linear time-invariant models of finite order, when it comes to measuring the *complexity* of a given model, the model order and the number of parameters are the most widely used characteristics. Here, the model order is defined by e.g. the number of regressors in a regression model, the dimension of the state-space in a state-space model, the maximum lag in an autoregressive-moving-average model (ARMA) model, or the McMillan degree of a rational transfer function. Model simplification is often carried out by model order reduction techniques, such as Krylov subspace methods and balanced truncation (see [9] for a recent detailed review of model reduction approaches) and well known information theoretic criteria such as AIC, BIC, FPE, etc. [93], are based on a trade off between the model order, the number of parameters, and the quality of the achieved fit between the model and the data. Another well known measure for the complexity of a model is *the number of nonzero parameters* used in the model. Such a measure is standard in statistics, e.g., when using multivariate linear regression models and ANOVA tables. Many statistical tests have been developed to decide whether the value of an estimated parameter differs significantly from zero, and to determine the contribution of a regressor in explaining the observed data. The outcomes of such tests determine in part the iterative process by which regressors are included in or discarded from a regression model.

The number of nonzero parameters in a model determines the *sparsity* of that model, which may be defined as the proportion of zero parameters among the total number of model parameters. The question of how to obtain an accurate sparse model is relevant for many applications, for instance in engineering with respect to the design of filters implemented in hardware on analog or digital chips (where space limitations

and the complexity of component interconnection is an issue) [7, 50], in systems biology to determine the dominant interactions between a large number of genes and proteins [38, 14], in biomedicine and epidemiology to determine the main risk factors for certain pathologies [56, 101], in data sensing and compression to describe a signal using only the most relevant components [27, 18], and so on. It becomes especially relevant to take sparsity into account at an early stage of the system identification procedure in situations where only a limited amount of input-output data is available, possibly of relatively low quality (due to high noise levels, limited opportunity to carry out experiments, high costs involved, etc.). Also, the use of a sparse model, containing only a small number of nonzero parameters, may contribute to resolving identifiability problems.

1.2 Thesis outline

In this thesis an approach to sparse identification is advocated which employs an ℓ_2 -norm to optimize the fit between a model and the data (using a conventional least squares criterion with respect to the vector of prediction errors) and an ℓ_1 -norm to minimize the size of the parameter vector to achieve model sparsity. In Part I of this thesis, a general framework for sparse estimation is presented as a widely applicable method to improve model sparsity (Chapter 2). It is applied to the class of regression models in Chapter 3 to study the settings in which sparse estimation is likely to be successful. In Chapter 4 sparse estimation is applied to the class of state-space models in innovations form, using either a full parameterization, data-driven local coordinates (DDLC) or a structured parameterization. The focus is then shifted to network models, where the sparse estimation algorithm is employed to find dominant network interactions in discrete- and continuous-time networks.

In Part II of this thesis a number of applications is presented of sparse estimation in the field of atrial fibrillation research. Chapters 5 and 6 provide some background on the characteristics of atrial fibrillation, the analysis of invasive atrial measurements and complexity quantification by means of fibrillation wave construction. In Chapter 7 one of the techniques presented in Chapter 4 is applied to find dominant interactions between sites in invasive high-density recordings of atrial fibrillation. Finding dominant predictors of pharmacological cardioversion outcome is the subject of Chapter 8, where from a large number of candidate predictors of successful cardioversion, a subset is selected using elastic net regression.



“In der Beschränkung zeigt sich erst der Meister.”

— J.W. von Goethe, *Was wir bringen*

Chapter 2

A framework for sparse estimation

2.1 Introduction

The algorithm presented in this chapter is designed for the sparse estimation of a model from any model class \mathcal{M} , defined in the following way: a parameterization P is defined as a mapping $P : \Theta \rightarrow \mathcal{M}$ with $\Theta \subseteq \mathbb{R}^N$, so that every $\theta \in \Theta$ is mapped to a model from the model class \mathcal{M} : $\theta \mapsto P(\theta)$. The model class \mathcal{M} is implicitly defined by the set of all possible model representations $\{P(\theta) \mid \theta \in \Theta\}$. It is assumed that the fit between a model $P(\theta)$ and observed “data” can be expressed by an error vector $e(\theta) = f(P(\theta), \text{data})$. The criterion to assess the goodness of fit between the model and the data is taken to be a *least squares* criterion $V(\theta) = \|e(\theta)\|_2^2$, possibly normalized by the length of the error vector. The goal of the algorithm is 1) to maximize the quality of the fit of the model and 2) to maximize the sparsity of the model, while retaining the quality of the fit for certain (optimal) parameters θ^* . The space to search for better model sparsity is the space of *equivalent* models. There are a number of sources of equivalence between models from the model class:

- I Depending on the context there may be a natural equivalence relation in the sense that the observed model behavior is identical, i.e. two models parameterizations $\theta, \tilde{\theta} \in \Theta$ are equivalent if $P(\theta) = P(\tilde{\theta})$.
- II In the context of *model fitting* (maybe from the perspective of achieving good prediction properties) two models are considered equivalent if $e(\theta) = e(\tilde{\theta})$.
- III While optimizing the least squares error criterion $V(\theta)$ two models are equivalent if $V(\theta) = V(\tilde{\theta})$. This constitutes larger equivalence spaces, that contain the previous case, since models with equal error vectors e have the same criterion value.

Every source of equivalence can be perceived as an equivalence relation on Θ , that can be used for improving sparsity.

Table 2.1: Model equivalence sources

Context	Equivalence relation
Natural equivalence	$P(\theta) = P(\hat{\theta})$
Model fit equivalence	$e(\theta) = e(\hat{\theta})$
Criterion equivalence	$V(\theta) = V(\hat{\theta})$

The sparse system identification procedure can be stated in terms of the steps in the general identification process as described in for instance [58]. A schematic overview of these steps is given in Figure 2.1.

Prior knowledge Known facts about the identification problem at hand can be incorporated into the process: a specific model class and the associated least squares criterion, known parameter values, the subset of the parameter vector θ that is subject to sparse maximization.

Observed data Data can be anything that can be used to express the fit between a model and a desired outcome, ranging from measurements (e.g. input/output data or frequency domain data) to a reference model that has to be approximated.

Model set The model class \mathcal{M} and its parameterization $P(\theta)$ have to be determined.

Calculate model In this phase, a two-step iterative algorithm is executed to fit the model to the data and to maximize the model sparsity. In Section 2.2 the procedure to maximize sparsity is explained. The iterative algorithm is described in Section 2.3

Validate model The model validation step is very much dependent on the nature of the observed data and the goal of the estimation procedure as a whole. Therefore, a detailed discussion on model validation is postponed to the chapters containing practical examples.

The following sections describe the algorithm in general terms, since it is applicable to a large range of model classes. The basic components of the algorithm can be explained using only the assumptions that the model can be parametrized by a vector θ and that the fit of the model can be expressed as a least squares criterion of an error vector that is only dependent on θ : $e(\theta)$. Practical applications of the algorithm are given in the next chapters.

2.2 Maximizing sparsity

Maximizing the sparsity of a model is possible if there exists an equivalence class for the model at hand as defined in the previous section. By definition, maximizing sparsity means maximizing the number of zero entries in the model parameter vector θ or,

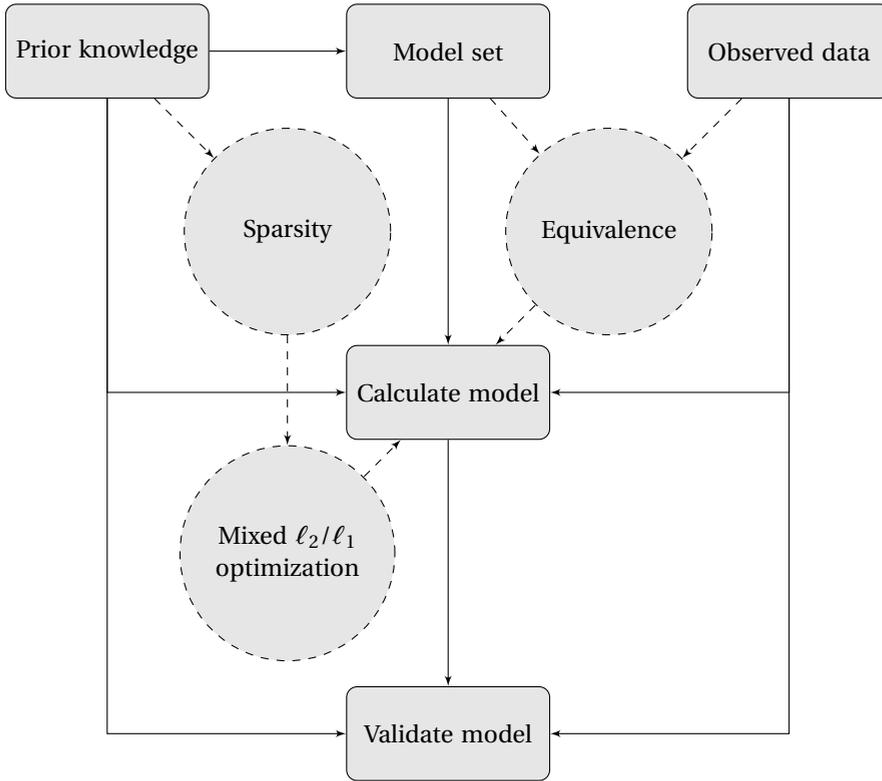


Figure 2.1: Sparse identification in the system identification loop

equivalently, minimizing the number of non-zero entries $\|\theta\|_0$, where $\|\cdot\|_0$ denotes the pseudo-norm ℓ_0 of a vector:

$$\|\theta\|_0 := \#\{\theta_i \mid \theta_i \neq 0\}. \quad (2.1)$$

Unfortunately this results in a difficult problem, often NP hard, that has to be solved using a combinatorial approach which can be time-consuming. A heuristic approach which is often much more efficient is to try and minimize the ℓ_1 -norm of the parameter vector θ , $\|\theta\|_1$. This heuristic is applied here since it has been shown that in a linear least squares setting conditions can be formulated (see [10, 36]) under which minimizing the ℓ_1 -norm can produce a parameter vector with maximum sparsity. The proposed method to maximize sparsity also features the possibility to target only a subset of the parameter vector, which allows one to incorporate prior knowledge about the value of certain parameters into the estimation procedure.

To see how the space of equivalent models can be defined starting from a specific parameterization θ_0 , consider the following situation: at a given model representation $P(\theta_0)$ and M data records, the error vector $e(\theta_0) = (e_1(\theta_0), e_2(\theta_0), \dots, e_M(\theta_0))^T$ and the least squares error criterion $V(\theta_0) = \|e(\theta_0)\|_2^2$ are computed. As seen in the previous section they offer three possible equivalence spaces to search for sparsity: 1) a natural equivalence space where the parameterization of a different vector $\theta \neq \theta_0$ yields a

model with identical behavior, 2) the space in which the error vector $e(\theta_0)$ does not change and 3) the space in which the criterion $V(\theta_0)$ does not change. Note that the second space is the space that contains models that result in the same errors with regard to the data, while the third space contains models that share the same goodness of fit, but not necessarily the same errors with respect to the data. Formally, equivalence is defined as:

(I) $\theta \cong \theta_0 : \Leftrightarrow P(\theta) = P(\theta_0)$, in the case of natural equivalence,

(II) $\theta \cong \theta_0 : \Leftrightarrow e(\theta) = e(\theta_0)$, when the individual errors need to remain unchanged, and

(III) $\theta \cong \theta_0 : \Leftrightarrow V(\theta) = V(\theta_0)$, when only the goodness of fit matters.

The definitions II and III can be used to determine a *local direction* s to search for sparsity, starting at the given model $P(\theta_0)$: $\theta = \theta_0 + s$. First, the error vector $e(\theta_0)$ is not allowed to change in the direction s . The Taylor series expansion of $e(\theta)$ about θ_0 is:

$$\begin{aligned} e(\theta) &= e(\theta_0) + \frac{\partial e}{\partial \theta}(\theta_0)(\theta - \theta_0) + \mathcal{O}(\|\theta - \theta_0\|^2) \\ &= e(\theta_0) + J(\theta_0)s + \mathcal{O}(\|s\|^2), \end{aligned} \quad (2.2)$$

where $J(\theta_0)$ is the $M \times N$ Jacobian matrix of partial derivatives of e to the parameter vector θ_0 . It follows that in order to stay within the first equivalence space in first order approximation, the search direction has to satisfy $J(\theta_0)s = 0$. This means that the space in which the entire error vector $e(\theta_0)$ does not change in first order approximation is described by the kernel of $J(\theta_0)$: $s \in \text{Ker}(J(\theta_0))$. The third equivalence space consists of those directions that do not change the least squares error criterion $V(\theta_0)$. The Taylor series expansion for $V(\theta)$ about θ_0 is:

$$\begin{aligned} V(\theta) &= V(\theta_0) + \frac{\partial V}{\partial \theta}(\theta_0)(\theta - \theta_0) + \frac{1}{2}(\theta - \theta_0)^T \frac{\partial^2 V}{\partial \theta^2}(\theta_0)(\theta - \theta_0) + \mathcal{O}(\|\theta - \theta_0\|^3) \\ &= e^T(\theta_0)e(\theta_0) + 2e^T(\theta_0)\frac{\partial e}{\partial \theta}(\theta_0)s + \frac{1}{2}s^T \frac{\partial^2 V}{\partial \theta^2}(\theta_0)s + \mathcal{O}(\|s\|^3) \\ &= e^T(\theta_0)e(\theta_0) + 2e^T(\theta_0)J(\theta_0)s + \frac{1}{2}s^T H(\theta_0)s + \mathcal{O}(\|s\|^3) \end{aligned} \quad (2.3)$$

where $H(\theta_0)$ denotes the Hessian matrix of second order partial derivatives of $V(\theta_0)$ to the parameter vector θ_0 . Clearly, if $J(\theta_0)s = 0$, then the value of $V(\theta)$ does not change in first approximation either. Suppose θ_0 is a local optimum of the criterion V , then the gradient of $V(\theta_0)$, $\frac{\partial V}{\partial \theta}(\theta_0)$, must be zero. But then of course any search direction does not change the value of V in first order approximation. Limiting the search direction to $s \in \text{Ker}(J(\theta_0))$ is still a valid option, but the space in which the criterion does not change in second order approximation is a more meaningful source to search for promising directions when θ_0 is a locally optimal parameterization:

$$\frac{1}{2}s^T H(\theta_0)s = 0 \Leftrightarrow H(\theta_0)s = 0 \Leftrightarrow s \in \text{Ker}(H(\theta_0)). \quad (2.4)$$

Under certain conditions the two spaces coincide. To see this, write the Hessian matrix in terms of the error vector $e(\theta_0)$:

$$\begin{aligned}
\frac{\partial V}{\partial \theta}(\theta_0) &= 2e^T(\theta_0) \frac{\partial e}{\partial \theta}(\theta_0) = 2 \sum_{k=1}^M e_k(\theta_0) \frac{\partial e_k}{\partial \theta}(\theta_0) \\
H(\theta_0) &= \frac{\partial^2 V}{\partial \theta^2}(\theta_0) \\
&= 2 \sum_{k=1}^M \left(\frac{\partial e_k^T}{\partial \theta}(\theta_0) \frac{\partial e_k}{\partial \theta}(\theta_0) + e_k(\theta_0) \frac{\partial^2 e_k}{\partial \theta^2}(\theta_0) \right) \\
&= 2J^T(\theta_0)J(\theta_0) + 2 \sum_{k=1}^M e_k(\theta_0) \frac{\partial^2 e_k}{\partial \theta^2}(\theta_0) \\
&= 2J^T(\theta_0)J(\theta_0) + 2S(\theta_0), \tag{2.5}
\end{aligned}$$

where

$$S(\theta_0) = \sum_{k=1}^M e_k(\theta_0) \frac{\partial^2 e_k}{\partial \theta^2}(\theta_0). \tag{2.6}$$

When $S(\theta_0)$ is sufficiently small, then (2.5) can be approximated by:

$$H(\theta_0) \approx 2J^T(\theta_0)J(\theta_0), \tag{2.7}$$

an approximation which is also used by Gauss-Newton type methods for least squares optimization problems. Substituting this into (2.4) gives:

$$s^T (J^T(\theta_0)J(\theta_0)) s = 0 \Leftrightarrow \|J(\theta_0)s\| = 0 \Leftrightarrow J(\theta_0)s = 0 \Leftrightarrow s \in \text{Ker}(J(\theta_0)), \tag{2.8}$$

which shows that in this case both approaches lead to a local equivalence space formed by the kernel of the Jacobian matrix of the error vector.

These observations clarify the available approximations to the local equivalence spaces at a given parameteriation θ_0 . The next step is to maximize the number of zero entries in the parameter vector, starting at θ_0 , limiting the search space to the selected equivalence space.

2.2.1 Minimizing the ℓ_1 -norm of the parameter vector θ

Sparsity is sought after by minimizing the ℓ_1 -norm of the parameter vector θ of the system at hand. Minimizing the ℓ_1 -norm of a parameter vector from an affine space can be expressed in terms of a linear programming (LP) problem. The search space is constrained, either by the kernel of the Jacobian matrix J of the error vector $e(\theta)$ at a certain model estimate θ_0 or by the kernel of the Hessian matrix H of the least squared error criterion V . This means that the search is limited to those models that a) produce *exactly* the same error vector e in first order approximation when J is used, or b) models that produce the same value of the least squares error criterion V in second order approximation when H is used. The minimization problem can be formulated as:

$$\min_s \|\theta_0 + s\|_1 \quad \text{s.t.} \quad Ks = 0 \tag{2.9}$$

where K represents either J or H . This LP problem is not in standard form. A linear programming problem in standard form only features a linear objective function, linear (in)equality constraints and positive decision variables (see for instance [22]):

$$\begin{aligned} \min_x c^T x \quad \text{s.t.} \quad & Ax \geq b \\ & x \geq 0, \end{aligned} \quad (2.10)$$

with the N -vector of decision variables x , A an $M \times N$ matrix, b an $M \times 1$ vector, and c an $N \times 1$ vector. Its corresponding *dual* problem is given by:

$$\begin{aligned} \max_y b^T y \quad \text{s.t.} \quad & A^T y \leq c \\ & y \geq 0, \end{aligned} \quad (2.11)$$

with y an $M \times 1$ vector of dual decision variables. There are several options to bring (2.9) into standard form. One way is by applying two substitutions. First, let θ be the result of taking a step s starting from θ_0 : $\theta = \theta_0 + s$. Then the formulation becomes:

$$\begin{aligned} \min_{\theta} \|\theta\|_1 \\ \text{s.t.} \quad & K(\theta - \theta_0) = 0 \\ & \Leftrightarrow \\ \min_{\theta} \|\theta\|_1 \\ \text{s.t.} \quad & K\theta = K\theta_0 \end{aligned} \quad (2.12)$$

To linearize the non-linear objective function, a second substitution can be applied:

$$\theta_i = \theta_i^+ - \theta_i^-, \quad \theta_i^+, \theta_i^- \geq 0, \quad i = 1, 2, \dots, N \quad (2.13)$$

which means that each variable θ_i (positive or negative) is split into two positive variables θ_i^+ and θ_i^- . The ℓ_1 -norm of the vector θ in the objective function can now be replaced by a linear expression:

$$\begin{aligned} \min_{\theta^+, \theta^-} (|\theta_1^+ - \theta_1^-| + |\theta_2^+ - \theta_2^-| + \dots + |\theta_N^+ - \theta_N^-|) \\ \text{s.t.} \quad & K(\theta^+ - \theta^-) = K\theta_0 \\ & \theta_i^+, \theta_i^- \geq 0 \quad i = 1, 2, \dots, N \end{aligned} \quad (2.14)$$

$$\begin{aligned} \Leftrightarrow \\ \min_{\theta^+, \theta^-} (\theta_1^+ + \theta_1^- + \theta_2^+ + \theta_2^- + \dots + \theta_N^+ + \theta_N^-) \\ \text{s.t.} \quad & K(\theta^+ - \theta^-) = K\theta_0 \\ & \theta_i^+, \theta_i^- \geq 0 \quad i = 1, 2, \dots, N, \end{aligned} \quad (2.15)$$

because an optimal solution to (2.15) will consist of pairs $\langle \theta_i^+, \theta_i^- \rangle$ with at least one zero entry, ($\langle \theta_i^+, 0 \rangle$ when $\theta_i > 0$ and $\langle 0, \theta_i^- \rangle$ when $\theta_i < 0$), since a mixed pair that represents the same value of θ_i yields a higher sum. The expression in (2.15) can be rewritten

to include the inequality constraint in (2.10), completing the conversion to standard form, with the following substitutions for the notation used in (2.10):

$$\begin{aligned} c &= \left(1 \quad 1 \quad 1 \quad \dots \quad 1 \right)^T \\ x &= \begin{bmatrix} \theta^+ & \theta^- \end{bmatrix} \\ A &= \begin{bmatrix} K & -K \end{bmatrix} \\ b &= K\theta_0. \end{aligned} \tag{2.16}$$

This leads to the following theorem:

Theorem 2.2.1. *The minimization problem in (2.9) can be restated as a LP problem in standard form. Solving the linear program*

$$\min_s \|\theta_0 + s\|_1 \quad \text{s.t.} \quad Ks = 0$$

is equivalent to solving:

$$\begin{aligned} \min_{\theta^+, \theta^-} & (\theta_1^+ + \theta_1^- + \theta_2^+ + \theta_2^- + \dots + \theta_N^+ + \theta_N^-) \\ \text{s.t.} & K(\theta^+ - \theta^-) = K\theta_0 \\ & \theta_i^+, \theta_i^- \geq 0 \quad i = 1, 2, \dots, N \end{aligned}$$

where $\theta = \theta_0 + s$ and $\theta = \theta^+ - \theta^-$.

The LP problem in standard form can be solved by the simplex algorithm and numerous other solvers such as the active set methods and interior point methods. Some solvers require the computation of the dual problem as well. The dual problem is easily derived from the primal problem in standard form, using the (free) dual variable vector δ :

$$\begin{aligned} \max_{\delta} & (K\theta_0)^T \delta \\ \text{s.t.} & K^T \delta \leq 1 \quad (\text{restricted variables } \theta^+) \\ & -K^T \delta \leq 1 \quad (\text{restricted variables } \theta^-) \end{aligned} \tag{2.17}$$

The number of constraints in the primal problem is M when $K = J(\theta_0)$ and N when $K = H(\theta_0)$. This number can possibly be reduced by applying singular value decomposition (SVD) to the matrix K . Suppose that K is an $(M \times N)$ matrix. Singular value decomposition factors the matrix into a product of three matrices $K = UDV^T$ where U ($M \times M$) and V ($N \times N$) are orthonormal matrices and D ($M \times N$) is a diagonal matrix whose elements are the singular values of K . An orthonormal matrix is a matrix whose columns are perpendicular to each other (orthogonal) and have unit length. A convenient property of an orthonormal matrix is that multiplying it with its transpose gives the identity matrix: $U^T U = I$. The matrix D can be partitioned into a strictly positive square diagonal matrix D_1 and zero blocks. If one or more of the singular values of K

is zero, then this results in a smaller set of constraints for the minimization problem.

$$\begin{aligned}
 K &= UDV^T \\
 &= \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} D_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_1 & V_2 \end{bmatrix}^T \\
 &= U_1 D_1 V_1^T.
 \end{aligned} \tag{2.18}$$

Substituting $K = U_1 D_1 V_1^T$ into (2.15) allows to reformulate the constraints for the LP problem, taking advantage of the orthonormality of U_1 and the invertibility of D_1 :

$$\begin{aligned}
 K(\theta^+ - \theta^-) &= K\theta_0 \\
 U_1 D_1 V_1^T (\theta^+ - \theta^-) &= U_1 D_1 V_1^T \theta_0 \\
 U_1^T U_1 D_1 V_1^T (\theta^+ - \theta^-) &= U_1^T U_1 D_1 V_1^T \theta_0 \\
 D_1 V_1^T (\theta^+ - \theta^-) &= D_1 V_1^T \theta_0 \\
 D_1^{-1} D_1 V_1^T (\theta^+ - \theta^-) &= D_1^{-1} D_1 V_1^T \theta_0 \\
 V_1^T (\theta^+ - \theta^-) &= V_1^T \theta_0,
 \end{aligned} \tag{2.19}$$

which implies that the matrix K can be substituted by V_1^T everywhere.

A selection matrix C can be introduced to target only a subset $y = C\theta$ of the model parameters in the search for sparsity:

$$\begin{aligned}
 \min_{y, \theta} & \|y\|_1 \\
 \text{s.t.} & K\theta = K\theta_0 \\
 & y = C\theta
 \end{aligned} \tag{2.20}$$

with y a $(n \times 1)$ vector and $n \leq N$. The problem (2.20) can again be reformulated by a substitution for y , similar to the one for θ in (2.13):

$$y_i = y_i^+ - y_i^-, \quad y_i^+, y_i^- \geq 0 \quad i = 1, 2, \dots, n. \tag{2.21}$$

This gives the partial sparsity maximization problem with restricted variables y^+ and y^- :

$$\begin{aligned}
 \min_{y^+, y^-, \theta} & (y_1^+ + y_1^- + y_2^+ + y_2^- + \dots + y_n^+ + y_n^-) \\
 \text{s.t.} & K\theta = K\theta_0 \\
 & y^+ - y^- - C\theta = 0 \\
 & y_i^+, y_i^- \geq 0 \quad i = 1, 2, \dots, n
 \end{aligned} \tag{2.22}$$

The dual problem can be formulated by introducing dual variable vectors δ for the constraints $K\theta = K\theta_0$ and ϵ for the constraints $y^+ - y^- - C\theta = 0$ (see also [77]):

$$\begin{aligned}
 \max_{\delta, \epsilon} & (K\theta_0)^T \delta \\
 \text{s.t.} & K^T \delta - C^T \epsilon = 0 \quad (\text{free variables } \theta) \\
 & \epsilon \leq 1 \quad (\text{restricted variables } y^+) \\
 & -\epsilon \leq 1 \quad (\text{restricted variables } y^-)
 \end{aligned} \tag{2.23}$$

2.2.2 Relationship between primal and dual problem

Formulating both the primal *and* the dual problem of the minimization problem in (2.9) has the advantage that the dual problem can be solved first to reduce the dimensions of the primal problem, which may speed up computations in certain cases. In the remainder of this section, the matrix K is assumed to be an $(M \times N)$ matrix to make the notation unambiguous, although M is equal to N when the Hessian matrix is used.

Theorem 2.2.2. *An optimal solution δ^* to the dual problem in (2.17) is related to a corresponding optimal solution $\langle \theta^{+*}, \theta^{-*} \rangle$ of the primal problem in (2.15) by the following rules:*

$$\theta_j^{+*} = \begin{cases} 0 & \text{if } \sum_{i=1}^M K_{ij} \delta_i^* < 1 \\ \geq 0 & \text{if } \sum_{i=1}^M K_{ij} \delta_i^* = 1 \end{cases} \quad j = 1, 2, \dots, N \quad (2.24)$$

$$\theta_j^{-*} = \begin{cases} 0 & \text{if } \sum_{i=1}^M K_{ij} \delta_i^* > -1 \\ \geq 0 & \text{if } \sum_{i=1}^M K_{ij} \delta_i^* = -1 \end{cases} \quad j = 1, 2, \dots, N. \quad (2.25)$$

Proof. According to the duality theorem (see e.g. [22]) an optimal solution y^* of the dual problem corresponds to an optimal solution of the primal problem x^* such that

$$\sum_{j=1}^n c_j x_j^* = \sum_{i=1}^m b_i y_i^*. \quad (2.26)$$

This property is sometimes referred to as “strong” duality, opposed to “weak” duality which states that a feasible solution to the primal problem has a criterion value that is equal or larger than any feasible solution of the dual problem and vice versa. To be able to write down Equation (2.26) in terms of θ and δ , the $(2N \times 1)$ vector $\tilde{\theta}$ is introduced to hold both θ^+ and θ^- : $\tilde{\theta} = \begin{bmatrix} \theta^+ \\ \theta^- \end{bmatrix}$, and the matrix $\tilde{K} = [K \ -K]$. The primal problem in (2.15) can now be written as

$$\begin{aligned} & \min_{\tilde{\theta}} (\tilde{\theta}_1 + \tilde{\theta}_2 + \dots + \tilde{\theta}_{2N}) \\ & \text{s.t. } \tilde{K} \tilde{\theta} = K \theta_0 \\ & \quad \tilde{\theta}_i \geq 0 \quad i = 1, 2, \dots, 2N, \end{aligned} \quad (2.27)$$

and the dual (2.17) as

$$\begin{aligned} & \max_{\delta} (K \theta_0)^T \delta \\ & \text{s.t. } \tilde{K}^T \delta \leq 1. \end{aligned} \quad (2.28)$$

Equation (2.26) now translates in this specific case to

$$\sum_{j=1}^{2N} \tilde{\theta}_j^* = \sum_{i=1}^M (K \theta_0)_i \delta_i^*. \quad (2.29)$$

where δ^* is an optimal solution to the dual problem and $\tilde{\theta}^*$ its corresponding optimal solution to the primal problem. The inequality constraints in (2.28) imply that

$$\tilde{\theta}_j^* \geq \left(\sum_{i=1}^M \tilde{K}_{ij} \delta_i^* \right) \tilde{\theta}_j^*, \quad j = 1, 2, \dots, 2N, \quad (2.30)$$

and the equality constraints in (2.27) imply that

$$\left(\sum_{j=1}^{2N} \tilde{K}_{ij} \tilde{\theta}_j^* \right) \delta_i^* = (K\theta_0)_i \delta_i^*, \quad i = 1, 2, \dots, M. \quad (2.31)$$

These two observations can be combined to relate the objective functions of the primal and dual problem to each other:

$$\sum_{j=1}^{2N} \tilde{\theta}_j^* \geq \sum_{j=1}^{2N} \left(\sum_{i=1}^M \tilde{K}_{ij} \delta_i^* \right) \tilde{\theta}_j^* = \sum_{i=1}^M \left(\sum_{j=1}^{2N} \tilde{K}_{ij} \tilde{\theta}_j^* \right) \delta_i^* = \sum_{i=1}^M (K\theta_0)_i \delta_i^*. \quad (2.32)$$

Now the equality in (2.29) holds if and only if equalities hold in (2.30). This means that for each j either $\tilde{\theta}_j^* = 0$ or $\sum_{i=1}^M \tilde{K}_{ij} \delta_i^* = 1$. The following rules can be derived to obtain information about the primal solution from the dual solution:

$$\tilde{\theta}_j^* = \begin{cases} 0 & \text{if } \sum_{i=1}^M \tilde{K}_{ij} \delta_i^* < 1 \\ \geq 0 & \text{if } \sum_{i=1}^M \tilde{K}_{ij} \delta_i^* = 1 \end{cases} \quad j = 1, 2, \dots, 2N. \quad (2.33)$$

Splitting the vector $\tilde{\theta}^*$ in θ^{+*} and θ^{-*} gives

$$\theta_j^{+*} = \begin{cases} 0 & \text{if } \sum_{i=1}^M K_{ij} \delta_i^* < 1 \\ \geq 0 & \text{if } \sum_{i=1}^M K_{ij} \delta_i^* = 1 \end{cases} \quad j = 1, 2, \dots, N \quad (2.34)$$

$$\theta_j^{-*} = \begin{cases} 0 & \text{if } \sum_{i=1}^M K_{ij} \delta_i^* > -1 \\ \geq 0 & \text{if } \sum_{i=1}^M K_{ij} \delta_i^* = -1 \end{cases} \quad j = 1, 2, \dots, N. \quad (2.35)$$

which concludes the proof. \square

Using the substitution in (2.13), the dual solution can be linked to the original LP formulation not in standard form in (2.9):

$$\theta_j^* = \begin{cases} 0 & \text{if } -1 < \sum_{i=1}^M K_{ij} \delta_i^* < 1 \\ \geq 0 & \text{if } \sum_{i=1}^M K_{ij} \delta_i^* = 1 \\ \leq 0 & \text{if } \sum_{i=1}^M K_{ij} \delta_i^* = -1 \end{cases} \quad j = 1, 2, \dots, N \quad (2.36)$$

which means that an optimal solution to the dual problem contains information on the signs of the primal variables in the corresponding optimal solution to the primal problem.

Corollary 2.2.3. *An optimal solution $\langle \delta^*, \epsilon^* \rangle$ to the partial sparsity dual problem in (2.23) is related to a corresponding optimal solution $\langle y^{+*}, y^{-*}, \theta^* \rangle$ of the primal problem in (2.22) by the following rules:*

$$y_j^{+*} = \begin{cases} 0 & \text{if } \epsilon_j^* < 1 \\ \geq 0 & \text{if } \epsilon_j^* = 1 \end{cases} \quad j = 1, 2, \dots, n. \quad (2.37)$$

$$y_j^{-*} = \begin{cases} 0 & \text{if } \epsilon_j^* > -1 \\ \geq 0 & \text{if } \epsilon_j^* = -1 \end{cases} \quad j = 1, 2, \dots, n. \quad (2.38)$$

Proof. The relation between an optimal solution to the *partial* sparsity dual problem in (2.23) and its corresponding optimal solution of the primal problem in (2.22) can be established in a similar manner as in the proof of Theorem 2.2.2, by introducing $\tilde{y} = \begin{bmatrix} y^+ \\ y^- \end{bmatrix}$ and $\tilde{I} = [I_n - I_n]$. The primal problem becomes:

$$\begin{aligned} & \min_{\tilde{y}, \theta} (\tilde{y}_1 + \tilde{y}_2 + \dots + \tilde{y}_{2n}) \\ & \text{s.t. } K\theta = K\theta_0 \\ & \quad \tilde{I}\tilde{y} - C\theta = 0 \\ & \quad \tilde{y}_i \geq 0 \quad i = 1, 2, \dots, 2n, \end{aligned} \quad (2.39)$$

and the dual:

$$\begin{aligned} & \max_{\delta, \epsilon} (K\theta_0)^T \delta \\ & \text{s.t. } K^T \delta - C^T \epsilon = 0 \\ & \quad \tilde{I}^T \epsilon \leq 1. \end{aligned} \quad (2.40)$$

Equation (2.26) now translates to

$$\sum_{j=1}^{2n} \tilde{y}_j^* = \sum_{i=1}^M (K\theta_0)_i \delta_i^*. \quad (2.41)$$

where δ^* is (part of) an optimal solution to the dual problem and \tilde{y}^* its corresponding optimal solution to the primal problem (2.39). Just as in the proof of Theorem 2.2.2, the two objective functions can be linked to each other, but in this case it requires a little bit more work to make the connection. The inequalities in (2.40) imply that

$$\tilde{y}_j \geq \left(\sum_{i=1}^n \tilde{I}_{ij} \epsilon_i \right) \tilde{y}_j \quad j = 1, 2, \dots, 2n \quad (2.42)$$

and the equalities in (2.39) imply that

$$\left(\sum_{j=1}^n K_{ij} \theta_j \right) \delta_i = (K\theta_0)_i \delta_i \quad i = 1, 2, \dots, M \quad (2.43)$$

Using the relations $\tilde{I}\tilde{y} = C\theta$ and $K^T\delta = C^T\epsilon$, equation (2.32) becomes:

$$\begin{aligned}
 \sum_{j=1}^{2n} \tilde{y}_j^* &\geq \sum_{j=1}^{2n} \left(\sum_{i=1}^n \tilde{I}_{ij} \epsilon_i^* \right) \tilde{y}_j^* = \sum_{i=1}^n \left(\sum_{j=1}^{2n} \tilde{I}_{ij} \tilde{y}_j^* \right) \epsilon_i^* \\
 &= \sum_{i=1}^n \left(\sum_{j=1}^N C_{ij} \theta_j^* \right) \epsilon_i^* = \sum_{j=1}^N \left(\sum_{i=1}^n C_{ij} \epsilon_i^* \right) \theta_j^* \\
 &= \sum_{j=1}^N \left(\sum_{i=1}^M K_{ij} \delta_i^* \right) \theta_j^* = \sum_{i=1}^M \left(\sum_{j=1}^N K_{ij} \theta_j^* \right) \delta_i^* \\
 &= \sum_{i=1}^M (K\theta)_i \delta_i^*. \tag{2.44}
 \end{aligned}$$

The equality in (2.41) now holds if and only if equalities hold in (2.42), which in turn means that either $\tilde{y}_j^* = 0$ or $\sum_{i=1}^n \tilde{I}_{ij} \epsilon_i^* = 1$. The following rules can be derived to obtain information about the primal solution from the dual solution:

$$\tilde{y}_j^* = \begin{cases} 0 & \text{if } \sum_{i=1}^n \tilde{I}_{ij} \epsilon_i^* < 1 \\ \geq 0 & \text{if } \sum_{i=1}^n \tilde{I}_{ij} \epsilon_i^* = 1 \end{cases} \quad j = 1, 2, \dots, 2n. \tag{2.45}$$

\Leftrightarrow

$$y_j^{+*} = \begin{cases} 0 & \text{if } \sum_{i=1}^n I_{ij} \epsilon_i^* < 1 \\ \geq 0 & \text{if } \sum_{i=1}^n I_{ij} \epsilon_i^* = 1 \end{cases} \quad j = 1, 2, \dots, n. \tag{2.46}$$

$$y_j^{-*} = \begin{cases} 0 & \text{if } \sum_{i=1}^n I_{ij} \epsilon_i^* > -1 \\ \geq 0 & \text{if } \sum_{i=1}^n I_{ij} \epsilon_i^* = -1 \end{cases} \quad j = 1, 2, \dots, n. \tag{2.47}$$

\Leftrightarrow

$$y_j^{+*} = \begin{cases} 0 & \text{if } \epsilon_j^* < 1 \\ \geq 0 & \text{if } \epsilon_j^* = 1 \end{cases} \quad j = 1, 2, \dots, n. \tag{2.48}$$

$$y_j^{-*} = \begin{cases} 0 & \text{if } \epsilon_j^* > -1 \\ \geq 0 & \text{if } \epsilon_j^* = -1 \end{cases} \quad j = 1, 2, \dots, n. \tag{2.49}$$

□

The value of the dual variable ϵ_j^* can be translated to the sign of y_j^* in (2.20):

$$y_j = \begin{cases} \geq 0 & \text{if } \epsilon_j = 1 \\ \leq 0 & \text{if } \epsilon_j = -1 \\ 0 & \text{if } -1 < \epsilon_j < 1 \end{cases} \quad j = 1, 2, \dots, n, \tag{2.50}$$

and the selection matrix C determines which variable signs in the original parameter vector θ can be determined from these rules.

2.2.3 Practical issues

In a practical setting, the algorithm to maximize sparsity has to deal with a number of numerical uncertainties: applying singular value decomposition to the matrix K (Ja-

cobian or Hessian) can possibly return some singular values close to zero. A threshold needs to be applied to determine which values are assumed to be zero and which not. This threshold has to be chosen carefully: a threshold that is too low will limit the search space and perhaps exclude valid directions to search for sparsity; a threshold that is too high can possibly result in search directions that change the model behavior. A similar choice has to be made in the case that the solution to the dual LP problem is computed and the solution to the primal LP problem is derived from this dual solution: when a value of $\sum_{i=1}^M \tilde{K}_{ij} \delta_i$ is close to 1 or -1 (or the value of ϵ_j in the case of the partial sparsity problem), a decision has to be made about whether to set the corresponding primal variable θ_i^+ or θ_i^- (or y_i^+ or y_i^-) to zero or not. Again, the threshold for this decision determines the quality of the maximization procedure. When deciding on appropriate threshold values one needs to take into account the model class (and data set) at hand and its robustness with respect to (small) aberrations in search direction accuracy during the optimization procedure.

2.3 Mixed optimization: minimizing the least squares error and maximizing sparsity

The procedure to maximize the sparsity of a certain model estimate, requires that there exists a procedure to arrive at a model that provides a good fit to the observed data. Only then there is an opportunity to search for a sparser solution that retains the quality of fit of the estimated model. In Section 2.2 the kernel of the Jacobian matrix J of the model error or the Hessian H of the least squares error criterion V is the subspace that contains the search directions to maximize sparsity. The orthogonal complement of the kernel of J , $\text{Ker}(J)^\perp$, determines the subspace of the parameter space in which the search directions generated by Gauss-Newton type methods are contained, see [102], to improve the fit of the model to the data. These observations lead to the following two key elements of the sparse identification procedure:

- (1) The orthogonal complement of the kernel of J is employed as the subspace in which to select a search direction for improving the (nonlinear) least squares error criterion V (such as, e.g., achieved by Gauss-Newton type optimization methods).
- (2) The kernel of J or the kernel of H if the parameter estimate θ in the first step is a (local) optimum of V , is employed as the subspace in which to select a search direction for improving the ℓ_1 -norm of (a part of) the parameter vector θ .

There are many ways to build an actual sparse identification algorithm from these two key elements. Algorithms may differ with respect to the amount of iterations of these two types, and the combinations and order in which they occur. For instance: one may perform steps in the two subspaces $\text{Ker}(J)$ and $\text{Ker}(J)^\perp$ simultaneously at each iteration; one may perform steps in the two subspaces alternatingly; one may first perform optimization of V by only taking steps in subspaces of the type $\text{Ker}(J)^\perp$ and then minimize the ℓ_1 -norm of θ afterwards by taking steps in subspaces of the type

$\text{Ker}(J)$. In the latter case, the value of V may deteriorate after a number of steps, so that it becomes necessary to incorporate intermediate optimization steps focused on the re-minimization of V . Other aspects in which algorithms may differ are: the computation of the search directions in the two subspaces (e.g., depending on the actual Gauss-Newton type method used); the computation of the corresponding step sizes in the computed directions.

2.3.1 Technical description

Suppose that a data record is given, which stems from a system S . At a given model $P(\theta)$ corresponding to a current parameter vector θ , one may compute the error vector $e(\theta)$ by running the associated error filter, and the Jacobian matrix $J(\theta)$ by running the associated sensitivity system of the system (see e.g. [40, 42, 76, 100]). Then the gradient of the criterion of fit $V(\theta) = \|e(\theta)\|^2$ is given up to a scalar factor by $J^T(\theta)e(\theta)$ and the Gauss-Newton (GN) method produces a search direction s given by

$$s = -(J^T(\theta)J(\theta))^{-1}J^T(\theta)e(\theta). \quad (2.51)$$

Variants of GN (including Levenberg-Marquardt and robust GN) differ in their choice for the inverse of the matrix $J^T(\theta)J(\theta)$ which appears in this formula: it may be replaced by $(\lambda I_N + J^T(\theta)J(\theta))^{-1}$ for some positive scalar λ , or by a suitably chosen pseudo-inverse. Note that all such methods yield a search direction s in the subspace $\text{Ker}(J(\theta))^\perp$, see also [102]. In contrast, the (damped) Newton method does not necessarily possess this property (see Table 2.2) and is therefore not applicable in this specific situation. Next, to improve the value of $V(\theta)$, the parameter vector θ is modified according to

$$\theta_{k+1} = \theta_k + \alpha s \quad (2.52)$$

for some step size parameter $\alpha > 0$. In the undamped GN method $\alpha = 1$, but to achieve good convergence behavior it is preferred to determine α by a suitable line minimization procedure (cf. [33]). $\text{Ker}(J)$ is the space in which the error vector $e(\theta)$ does not change locally around θ in first order approximation. The value of $\|\theta\|_1$ can be improved by computing a search direction in $\text{Ker}(J)$, as this will not affect the value of V in first order approximation. This leads to the following optimization problem:

$$\text{minimize } \|\theta + s\|_1 \quad \text{subject to: } Js = 0 \quad (2.53)$$

which can be rewritten as an LP problem in standard form, as described in Section 2.2. It clearly admits a finite feasible solution and can be solved with standard LP software. Suppose θ^* is a (local) minimum of V , then it is also possible to use the kernel of the Hessian H of the error criterion, which in this case is equivalent to the space in which $V(\theta^*)$ does not change locally around θ^* in *second* order approximation:

$$\text{minimize } \|\theta^* + s\|_1 \quad \text{subject to: } Hs = 0, \quad (2.54)$$

with θ^* a (local) minimum of V . Note that last approach can only be used in an algorithm that performs the sparsity maximization step *after* the minimization of the least squares error criterion.

Table 2.2: Overview of non-linear least squares minimization methods

Method	Search direction $s(\theta)$	Applicable
Steepest descent	$-\alpha J(\theta)^T e(\theta)$	Yes
Newton (damped)	$-\alpha H(\theta)^{-1} J(\theta)^T e(\theta)$	No
Gauss-Newton	$-\alpha (J(\theta)^T J(\theta))^{-1} J(\theta)^T e(\theta)$	Yes
Levenberg-Marquardt	$-(J(\theta)^T J(\theta) + \lambda I)^{-1} J(\theta)^T e(\theta)$	Yes

2.3.2 Practical issues

If s^* is an optimal solution giving an improved value in (2.53), then $\|\theta + \beta s^*\|_1$ gives an improved value too for all $0 < \beta < 1$. An actual choice of β should take into account that second and higher order changes in V do not significantly compromise the quality of the fit between the data and the model. It also must be tuned in such a way that convergence of the overall optimization algorithm can be guaranteed to a point for which $\|\theta\|_1$ is minimal among the set of points for which V is (locally) minimal. One practical heuristic way by which one may attempt to achieve this, is to restrict the maximal relative change in the value of V that is allowed to occur when a value for β is chosen. However, it is not easy to choose an appropriate bound which guarantees monotonic convergence: several experiments have been carried out which exhibit chaotic iteration behavior or cyclic behavior near a local optimum value. For a bound that is not restrictive enough, it has been witnessed that an *increase* of $\|\theta\|_1$ may happen instead of a decrease, resulting as the net effect of an optimization step with respect to $\|\theta\|_1$ followed by an optimization step with respect to V . This makes clear that the choice of β should be treated with care.

Chapter 3

Application: Linear regression models

3.1 Introduction

The class of linear regression models is treated here as a case study to investigate the properties and applicability of the proposed algorithm for maximizing sparsity. Linear regression models are applied in virtually every area of scientific research, rendering it one of the most influential model classes around. There are three elements to a linear regression model:

1. The *dependent variable(s)*, for instance observed outputs, measurements, classifiers.
2. The *regressors*, the entities that are candidates to explain the behavior of the dependent variables. In an estimation setting the regressors are associated with a *design matrix* which is a matrix with regressor instances. The regressors are functions of the *independent variables*.
3. The *parameter vector* of coefficients that characterize the dependent variables as a linear combination of the regressors.

The linear relationship between the regressors and the dependent variables makes the analysis and interpretation of the model relatively straightforward. The definition of linear regression models is also very flexible, enabling for instance the use of time-dependent regressors (or time itself) and the use of a larger or smaller amount of model parameters to model the observed behavior, thus increasing or decreasing the model order. This is where the sparse estimation algorithm comes into play: linear regression models are often used to determine which (combination of) entities determine the output or behavior of a certain phenomenon. A relevant problem is to select the most significant contributors from a set of candidate regressors. There are numerous methods that attempt to solve this problem. *Subset selection* is the general class

of methods that selects a subset of the regressors based on certain statistical criteria. *Forward selection*, for instance, starts with a zero order model and iteratively adds a new regressor from the set of candidates that best improves the model fit until k regressors are selected. *Backwards elimination* starts with a full regression of all available regressors and iteratively eliminates the ones that do not deteriorate the quality of the fit significantly. See for instance [66] and [67] for an overview of subset selection strategies. *Ridge regression* (or Tikhonov regularization) applies a regularization of the optimization criterion to favor solutions with certain properties, for instance solutions with smaller norms [45]. A related technique is the *lasso*, where the sum of the absolute value of the model parameters is constrained to be less than a constant value [95]. *Elastic net* [105] is a more recent variable selection and regularization method that combines the qualities of both ridge regression and the lasso technique. It is particularly useful in the case where the number of regressors is much larger than the number of observations, and when there are highly correlated regressors.

The sparse estimation algorithm described in the previous chapter can be applied to the problem of subset selection. In a situation where only a small portion of all the candidate regressors are assumed to contribute to the modeled phenomenon, the parameter vector is sparse. The *dominant* regressors can be identified by trying to minimize the number of non-zero parameters, while still retaining an optimal fit of the model to the observed data. The case under investigation will be the situation in which the amount or quality of the available data is insufficient to uniquely identify the contribution of each of the candidate regressors. In this situation the problem is *underdetermined* and there are multiple solutions that provide an optimal fit. The space of optimal solutions is the equivalence space to maximize the sparsity. The applicability of the sparse estimation algorithm in this linear regression settings is evaluated, specifically to get an insight into which conditions allow for reproducibility of the non-zero data generating parameters by sparse linear regression and which conditions do not.

3.2 The model class

The scalar linear regression problem can be stated as follows: the data available is of the form (ϕ_i, y_i) , $i = 1, 2, \dots, M$, where $\phi_i = [\phi_{i1}, \phi_{i2}, \dots, \phi_{iN}]^T$ is the regression vector and y_i is the dependent variable. The index i is used to distinguish between different states of the regression vector and the corresponding dependent variable. It can be time-related or just refer to different measurement records. The *deterministic* linear regression model class parametrized by θ , $\{P(\theta) | \theta \in \mathbb{R}^N\}$ is formulated in the following way:

$$y_i = \phi_i^T \theta, \quad i = 1, 2, \dots, M, \quad (3.1)$$

where θ is the $(N \times 1)$ parameter vector that describes the linear relationship between the regressors and the dependent variable. To be able to determine the N entries in the parameter vector θ uniquely, at least $M \geq N$ independent data records are required.

The model can be expressed in matrix form:

$$Y = \Phi\theta, \quad (3.2)$$

with $Y = [y_1, y_2, \dots, y_M]^T$ and $\Phi = [\phi_1, \phi_2, \dots, \phi_M]^T$. If $M = N$ and the matrix Φ is non-singular, there exists a unique solution for θ , given by $\theta = \Phi^{-1}Y$. In general, the number of measurements M will be different from the number of parameters, often much higher, which leads to an overdetermined system of equations that is not likely to have an exact solution, due to measurement disturbances or modeling errors. In this case the goal is to find θ that optimizes the fit between the model and the measurements. The error $e_i(\theta)$ between a model estimate and the observed output can be expressed as

$$e_i(\theta) = y_i - \phi_i^T \theta, \quad i = 1, 2, \dots, M, \quad (3.3)$$

and the least squares error criterion to express the quality of the model fit is

$$V(\theta) = \sum_{i=1}^M e_i^2(\theta). \quad (3.4)$$

or

$$V(\theta) = e^T(\theta)e(\theta) = \|e(\theta)\|^2, \quad (3.5)$$

where $e(\theta) = [e_1(\theta), e_2(\theta), \dots, e_M(\theta)]^T$. If $M \geq N$ and the corresponding matrix Φ has full column rank N , the unique least squares solution θ_{LS} is given by

$$\theta_{LS} = (\Phi^T \Phi)^{-1} \Phi^T Y. \quad (3.6)$$

If $M < N$ and/or the matrix Φ has a column rank $r < N$, there is no unique least squares solution, but an optimal least squares solution space of dimension $N - r$. The Moore-Penrose pseudo inverse $\Phi^+ = \Phi^T (\Phi \Phi^T)^{-1}$ is then commonly used to determine the minimum norm solution θ_{MN} with the smallest Euclidean norm of θ :

$$\theta_{MN} = \Phi^+ Y, \quad (3.7)$$

but it is also possible to compute an optimal solution that has at least $N - r$ zero entries, using QR decomposition and back-substitution, which intuitively provides an upper bound for the sparse maximization problem.

The two prerequisites for the sparse estimation algorithm are present in this model class: a parameterization by θ and a least squares criterion $V(\theta)$ to determine the goodness of fit of the model to the observed data. The iterative procedure described in the previous chapter can be brought down to just one iteration. Only in the underdetermined case, an equivalence space exists for the optimal least squares solution that can be used to maximize the sparsity of the parameter vector while retaining the optimal least squares fit. In the overdetermined case the equivalence space contains just the least squares solution θ_{LS} .

Theorem 3.2.1. *The equivalence space with respect to the least squares criterion V at a least squares solution θ_{LS} is determined by the kernel of the regressor matrix Φ*

$$V(\theta_{LS} + s) = V(\theta_{LS}) \Leftrightarrow s \in \text{Ker}(\Phi). \quad (3.8)$$

Proof. It is easily seen that a solution $\theta_{LS} + s$ with s in the kernel of the matrix Φ yields an error vector $e(\theta_{LS} + s)$ that is equal to $e(\theta_{LS})$ (i.e. model fit equivalence as defined in 2.1. An equal error vector produces an equal least squares criterion value, which proves the forward indication of the theorem. An equal criterion value for θ_{LS} and $\theta_{LS} + s$ also requires that s is in the kernel of Φ . Every optimal least squares estimate θ^* has the following property: $\Phi^T(Y - \Phi\theta^*) = 0$. Substituting $(\theta^* + s)$ for θ^* implies $-\Phi^T\Phi s = 0$, which is true if and only if $s \in \text{Ker}(\Phi)$. \square

In the linear least squares setting criterion equivalence implies model fit equivalence, which can be explained by observing that in this case $\Phi = J$ and $H = J^T J$.

Following the general sparse estimation scheme in the previous chapter, the algorithm can be stated as follows in the case of linear regression models. Given dependent variable observations Y and the design matrix Φ :

1. Compute a least squares solution θ_{LS}
2. Solve the linear program

$$\min_{\theta} \|\theta\|_1 \quad \text{subject to } \Phi\theta = \Phi\theta_{LS}. \quad (3.9)$$

The algorithm can also be stated as

$$\min_{\theta} \|\theta\|_1 \quad \text{subject to } \Phi\theta = Y^* \quad (3.10)$$

with $Y^* = \text{proj}_{\Phi} Y$ the orthogonal projection of the vector Y on the column space of Φ . There are a number of other methods that are closely related to this method. One of them is the more general class of convex problems that minimize the ℓ_1 -norm of the parameter vector under a non-linear constraint (see [37]):

$$\min_{\theta} \|\theta\|_1 \quad \text{subject to } \|\Phi\theta - Y\|_p \leq \rho_p \quad (3.11)$$

where $\|\cdot\|_p$ denotes the p -norm. For $\rho_p = 0$ the formulations (3.10) and (3.11) are equivalent formulations, that lead to the same solution(s). Ridge regression imposes an ℓ_2 -norm constraint on the parameter vector θ :

$$\min_{\theta} \|\Phi\theta - Y\|_2 \quad \text{subject to } \|\theta\|_2 \leq t, \quad (3.12)$$

Another method is the method known as the lasso (least absolute shrinkage and selection operator) [95] that is based on regression shrinkage, which minimizes a least squares criterion under a ℓ_1 -norm constraint for the parameter vector:

$$\min_{\theta} \|\Phi\theta - Y\|_2 \quad \text{subject to } \|\theta\|_1 \leq t, \quad (3.13)$$

where t is a tuning parameter. An optimal solution θ^* to (3.10) is an optimal solution to (3.13) if $t = \|\theta^*\|_1$. The elastic net method combines the penalties used in the ridge regression approach and the lasso:

$$\min_{\theta} \|\Phi\theta - Y\|_2 \quad \text{subject to } (1 - \alpha)\|\theta\|_1 + \alpha\|\theta\|_2 \leq t. \quad (3.14)$$

These methods promote a non-optimal least squares estimate (for $p = 2$ in (3.13) and (3.14)) in favor of a better ℓ_1 -norm of the parameter vector θ . If the sparse data generating model parameters are not an optimal solution to the least squares problem due to disturbances, this provides a opportunity to still retrieve that sparse solution. In this situation the sparse estimation algorithm will not produce a sparse solution, but an optimal least squares solution with minimum parameter ℓ_1 -norm. Under certain conditions the non-zero data generating parameters can still be recovered from this solution, which will be illustrated in the experiments.

3.3 Experiments

The sparse estimation algorithm is applied to the problem setting that has been previously defined in the case of linear regression models: underdetermined problems, where the regressor matrix Φ allows for a nontrivial space of solutions that minimize the least squares criterion V , giving way to an equivalence space to maximize the sparsity of the parameters. The case of over-determined problems corrupted by noise, where there exists a unique solution that minimizes V (and the equivalence space is restricted to that solution), will be briefly discussed, but is not the focus of this chapter.

3.3.1 Experimental setup

In the following experiments, the linear regression model is chosen to be to a linear multivariable model where the regressors are simply the independent variables themselves, providing an easy way to generate and analyze the simulations and results. A model is created by generating a random parameter vector θ of length N containing k non-zero entries, determining the sparsity $(1 - \frac{k}{N})$ of the model. The values of the non-zero parameters are drawn from a uniform distribution on the interval $[-2\alpha, -\alpha] \cup [\alpha, 2\alpha]$. The value of α is chosen in such a way that there is a clear distinction between zero and non-zero parameter values. A typical value is $\alpha = 1$. An $(M \times N)$ matrix Φ of regressor values is generated from a Gaussian distribution with zero mean and unit variance. The columns of Φ are normalized to unit length. The parameter vector θ and the regressor matrix Φ together produce the $(M \times 1)$ measurement vector Y . These measurements are optionally disturbed by Gaussian white noise with zero mean and variance σ^2 .

The performance of the algorithm is measured by the number of incorrectly estimated parameter values S_e , counting both false negatives (parameters that should be zero are assigned a non-zero value) and false positives (parameters that should be non-zero are assigned a zero value). In the noiseless case the parameter values are required to be equal up to machine precision, but in the presence of noise a threshold is

applied (see point 5 in Section 3.3.2). In the following experiments, the performance is visualized by plotting the mean performance over a number of trials and the mean minus and plus the standard deviation over these trials. A second way of assessing the algorithm performance is to compute the *probability* that it will succeed, meaning that the original non-zero and zero parameters are all estimated correctly. This performance parameter is denoted by P_0 . The performance of the algorithm is compared to the performance of the least squares solution that produces a parameter vector containing at most r non-zero entries, with r the rank of Φ and $r \leq N$. The probability that a least squares solution with at most r non-zero entries is the original data generating parameter vector is determined by the probability that the original k non-zero parameter positions are among the r selected non-zero least squares parameters. This probability can be expressed by the following formula

$$P(S_c = 0 | \theta_{LS}) = \binom{r}{k} \frac{k!(N-k)!}{N!}, \quad r \geq k. \quad (3.15)$$

All computations are performed using MATLAB. The least squares solution is computed using the backslash operator, giving an estimated parameter vector with at most r non-zero entries. The linear programming problem is solved by the `linprog` function that is available from the Optimization Toolbox.

3.3.2 Underdetermined problems

In situations where the column rank r of Φ is strictly smaller than the number of entries in the parameter vector θ , the problem is called underdetermined. More than one solution provides an optimal fit to the data. To see how this works, consider a small-scale example, where $N = 3$, $k = 2$ and $M = 2$.

$$\begin{bmatrix} 5 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 & 2 & -1 \\ 1 & -1 & 2 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}, \quad (3.16)$$

generated by $\theta^* = [1 \ 2 \ 0]^T$. The minimum norm solution for this problem is

$$\theta_{MN} = \left[1\frac{1}{3} \ 1\frac{2}{3} \ -\frac{1}{3} \right]^T. \quad (3.17)$$

The solution space for the system of equations in (3.16) is the line

$$\theta = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} + \theta_3 \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}. \quad (3.18)$$

The minimum norm solution is retrieved by choosing $\theta_3 = -\frac{1}{3}$. In this case the solution with maximum sparsity θ^* coincides with the solution with a minimal ℓ_1 -norm, which is confirmed by visual inspection of Figure 3.1. The sparse estimation algorithm is able to find this optimum starting from the minimum norm solution θ_{MN} . Note that there are actually more maximally sparse solutions with identical ℓ_0 -norm, but with higher ℓ_1 -norm for $\theta_3 = -2$ and $\theta_3 = 1$.

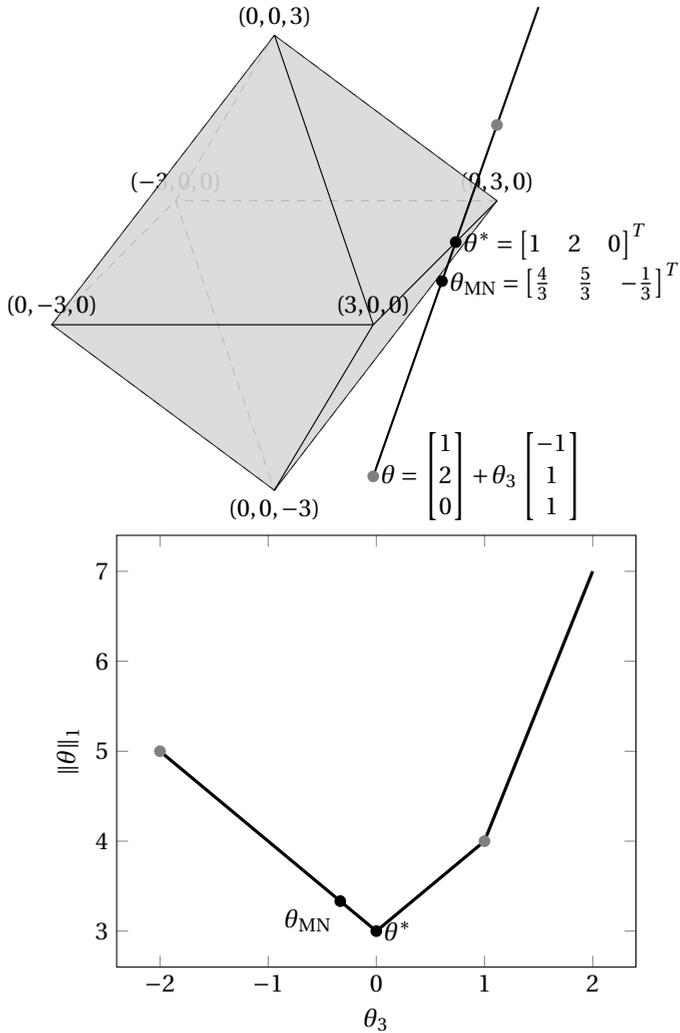


Figure 3.1: Visualization of the minimum norm solution θ_{MN} and the maximum sparsity solution θ^* for the problem in (3.16): (a) Visualization of the space $\|\theta\|_1 = 3$ and the least squares solution line. (b) $\|\theta\|_1$ along the least squares solution line, parameterized by θ_3 . Maximally sparse solutions with minimal ℓ_0 -norm but higher ℓ_1 -norm are marked in gray.

However, it is also possible to construct a problem that allows multiple maximally sparse solutions with equal minimal ℓ_1 -norms. Another problem of equal dimensions is

$$\begin{bmatrix} 1 \\ -2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & -2 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}, \quad (3.19)$$

generated by $\theta^* = [1 \ -2 \ 0]^T$. The minimum norm solution for this problem is

$$\theta_{\text{MN}} = [1\frac{1}{2} \ -1\frac{1}{2} \ \frac{1}{2}]^T. \quad (3.20)$$

The solution space for the system of equations in (3.19) is the line

$$\theta = \begin{bmatrix} 1 \\ -2 \\ 0 \end{bmatrix} + \theta_3 \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}. \quad (3.21)$$

In this case there are actually *two* maximally sparse solutions with minimal ℓ_1 -norm. Besides θ^* , also $\theta' = [2 \ 0 \ 1]^T$ contains as many zeros as possible and has the same ℓ_1 -norm. The solution space for the minimal ℓ_1 -norm of the parameter vector θ is the line segment between these two solutions, as can be seen in Figure 3.2.

Finally, it is also relatively straightforward to design a problem where the minimum ℓ_1 -norm solution does *not* correspond to a maximally sparse solution. Consider the problem

$$\begin{bmatrix} 3 \\ 6 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 1 \\ 2 & 2 & 4 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}, \quad (3.22)$$

generated by $\theta^* = [3 \ 0 \ 0]^T$. The solution space for the system of equations in (3.22) is the line

$$\theta = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} + \theta_3 \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}. \quad (3.23)$$

Here the solution with minimal ℓ_0 -norm (at $\theta_3 = 0$) has a higher ℓ_1 -norm than the minimal ℓ_1 -norm solution (at $\theta_3 = 1$), as is visualized in Figure 3.3.

To ensure that the sparse estimation algorithm selects one of the maximum sparsity solutions in the second example (Figure 3.2), the algorithm that computes the solution to the linear programming problem has to select an optimal vertex of the polytope that defines the feasible region. Only then a point is selected where one or more of the variables is zero. The simplex algorithm and the class of active set algorithms both possess this property, but for instance the interior point algorithm does not, which makes it not the best candidate to solve this particular LP problem. Only the simplex algorithm and the active set algorithm are used in the remainder of this chapter. The last example also shows that the data generating system may not be the only maximally

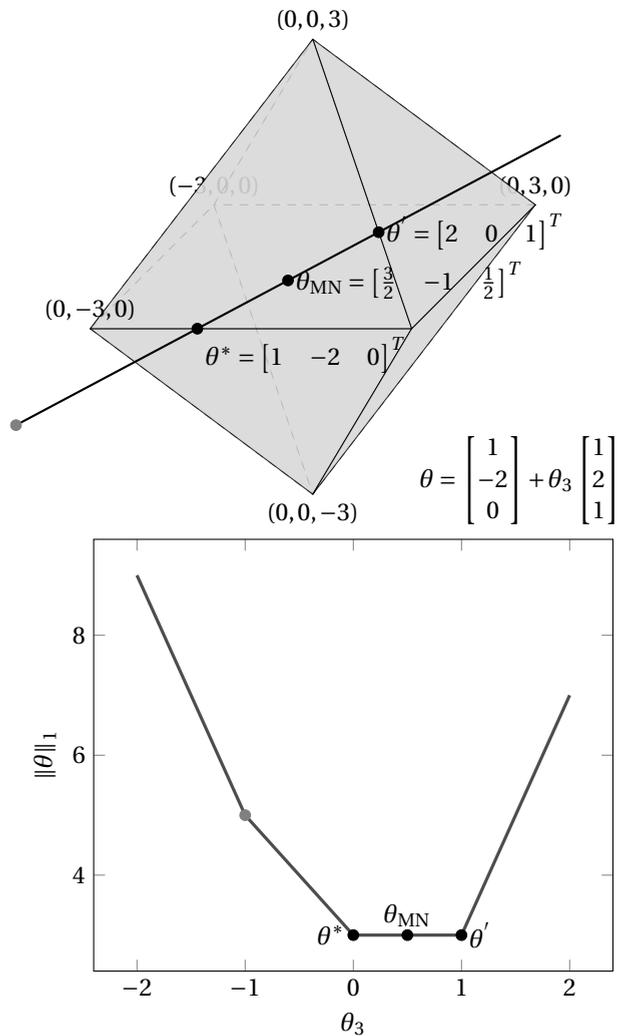


Figure 3.2: Visualization of the minimum norm solution θ_{MN} and the maximum sparsity solutions θ^* and θ' for the problem in (3.19): (a) Visualization of the space $\|\theta\|_1 = 3$ and the least squares solution line. (b) $\|\theta\|_1$ along the least squares solution line, parameterized by θ_3 .

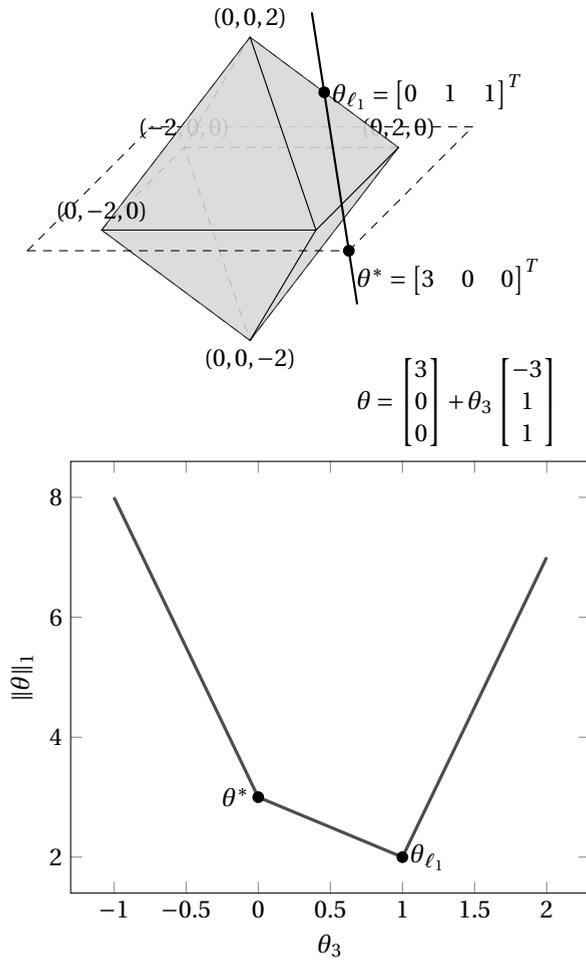


Figure 3.3: Visualization of the the maximum sparsity solutions θ^* and minimum ℓ_1 -norm solution θ_{ℓ_1} for the problem in (3.22): (a) Visualization of the space $\|\theta\|_1 = 2$ and the least squares solution line. (b) $\|\theta\|_1$ along the least squares solution line, parameterized by θ_3 .

sparse solution, or may even not be a maximally sparse solution. The algorithm can produce a different solution θ' in these cases, that has either the same or lower number of zero entries, or has a lower ℓ_1 -norm altogether. The question whether a original sparse vector can be retrieved from a set of measurements in a noiseless situation, is well studied. In [36] and [35] (and references mentioned there), conditions are given under which a parameter vector θ is the unique sparsest solution, namely:

$$\|\theta\|_0 \leq \frac{1}{C} \quad \text{with } C = \sup_{i \neq j} |\Phi_i^T \Phi_j|, \quad (3.24)$$

where Φ_i denotes the i th column of the matrix Φ , and conditions under which the sparsest solution is also the unique solution of the LP problem that minimizes the ℓ_1 -norm of the parameter vector:

$$\|\theta\|_0 \leq \frac{1}{2} \left(1 + \frac{1}{C} \right) \quad \text{with } C \text{ as above.} \quad (3.25)$$

These conditions are only valid if the columns of the regressor matrix Φ have unit length, a condition that was not met in the previous examples. If for instance the problem in (3.22) is scaled accordingly, the ℓ_1 -norm of the solution $\theta = [3 \ 0 \ 0]^T$ increases to $3\sqrt{5}$ (≈ 6.71), but the ℓ_1 -norm of the original minimal ℓ_1 -norm solution $\theta = [0 \ 1 \ 1]^T$ now exceeds this value by increasing to $\sqrt{8} + \sqrt{17}$ (≈ 6.95). The parameter C expresses the *mutual coherence* of the matrix Φ (usually denoted as M , but here a different notation is used to avoid confusion with the number of available measurements M). This gives a very strict bound, in many practical situations the conditions will not be met. It is therefore useful to investigate how the algorithm performs in a broader sense, i.e. how it performs *on average* in an number of different conditions. To investigate the probability that the sparse estimation algorithm does retrieve the original data generating model, and how many errors are produced when it does not, a number of experiments were carried out, where the algorithm performance was evaluated for varying problem dimensions N , k and M .

1. Dependence on the number of measurements M

The number of available measurements is an important factor in practical situations where the cost of acquiring measurements can be high. The algorithm is applied for a fixed parameter vector size N and a varying number of available measurements M , where $M < N$. The generation of the matrix Φ is carried out in such a way that $\text{rank}(\Phi) = M$. Results for $N = 10$ and $N = 100$ are shown in Figures 3.4 and 3.5. They indicate that in the situation where there is low sparsity together with a low number of measurements, the algorithm is not always able to find the desired solution. Furthermore, the number of measurements needs to be higher than the number of non-zero entries to obtain reasonable *average* performance in an underdetermined situation: $k < M < N$.

2. Dependence on the number of parameters N

The number of parameters represents the number of candidate regressors. Figure 3.6 shows the influence of the number of parameters on the number of mea-

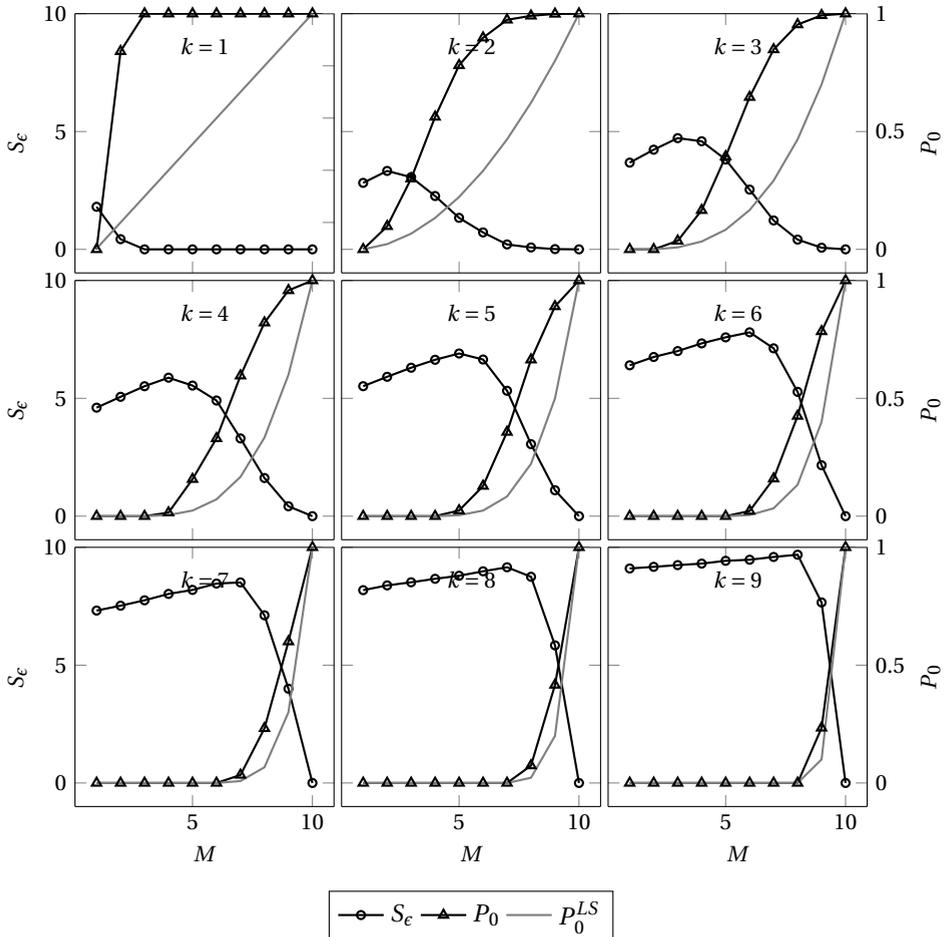


Figure 3.4: Average number of errors S_e and probability of a correct estimate P_0 depending on M ranging from 1 to 10, k ranging from 1 to 9 and $N = 10$. The data is not disturbed by noise and the number of trials for each combination of parameters is 1000.

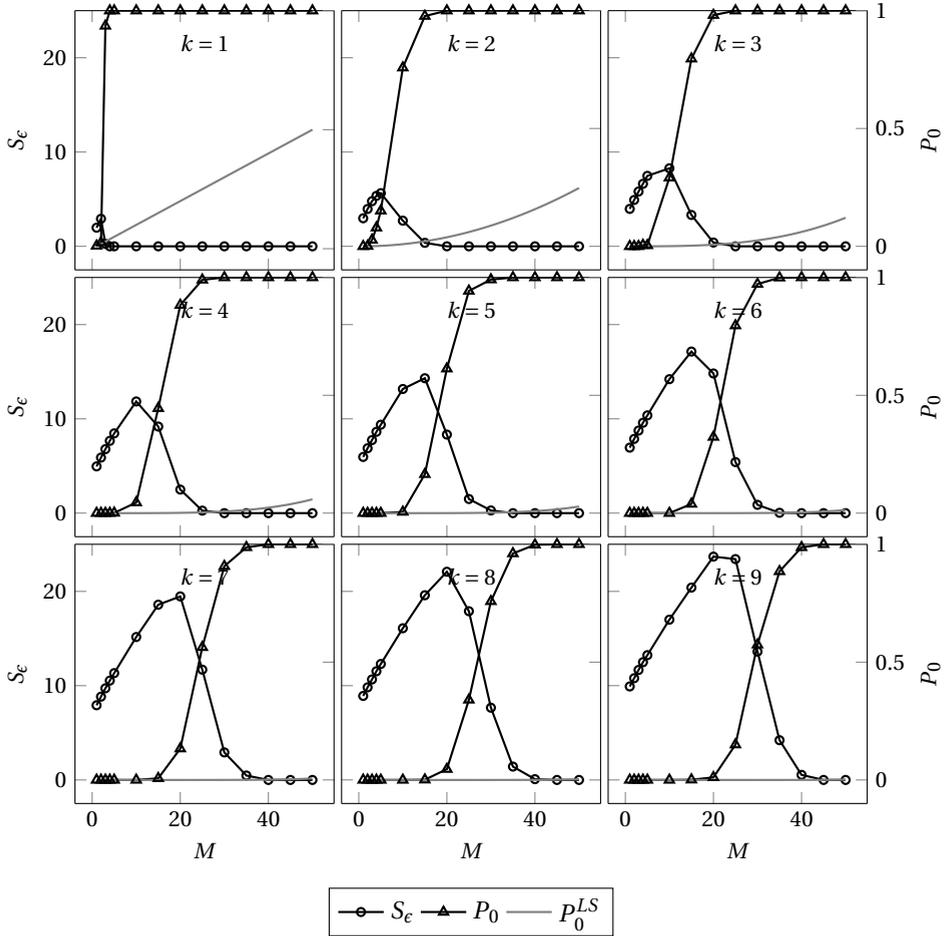


Figure 3.5: Average number of errors S_e and probability of a correct estimate P_0 depending on M ranging from 1 to 50, k ranging from 1 to 9 and $N = 100$. The data is not disturbed by noise and the number of trials for each combination of parameters is 1000.

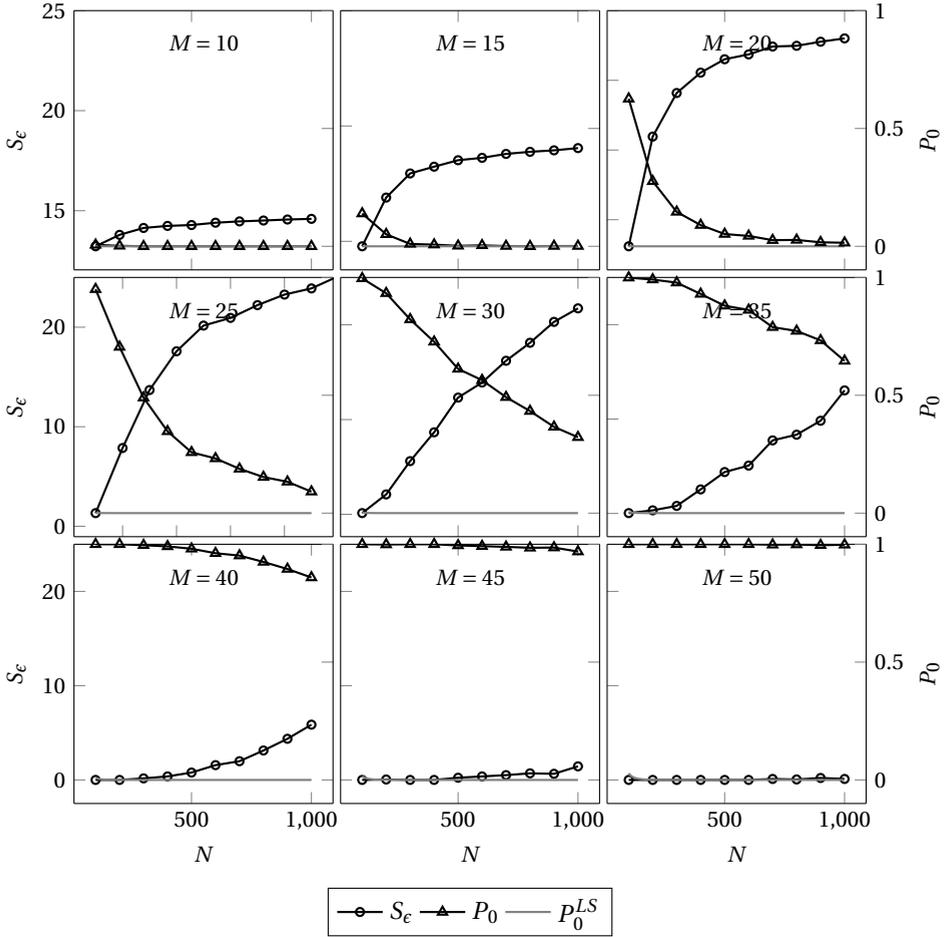


Figure 3.6: Average number of errors S_e and probability of a correct estimate P_0 depending on N , ranging from 100 to 1000, $k = 5$, M ranging from 10 to 50. The number of trials is 1000.

measurements needed to find the right solution for a fixed number of non-zero parameters. Clearly, the higher the total number of parameters, the higher the number of required measurements, but the number needed to estimate correctly is relatively low at high sparsity, for instance in the case where $N = 1000$, $k = 5$, only about 50 measurements are needed to succeed, compared to 1000 measurements needed to solve the same problem using only a conventional least squares criterion.

3. Dependence on the number of non-zero parameters k

The number of non-zero parameters k determines, together with the number of parameters N , the sparsity of the model. Figure 3.7 shows the effect of increasing the number of non-zero parameters on the algorithm performance. As k

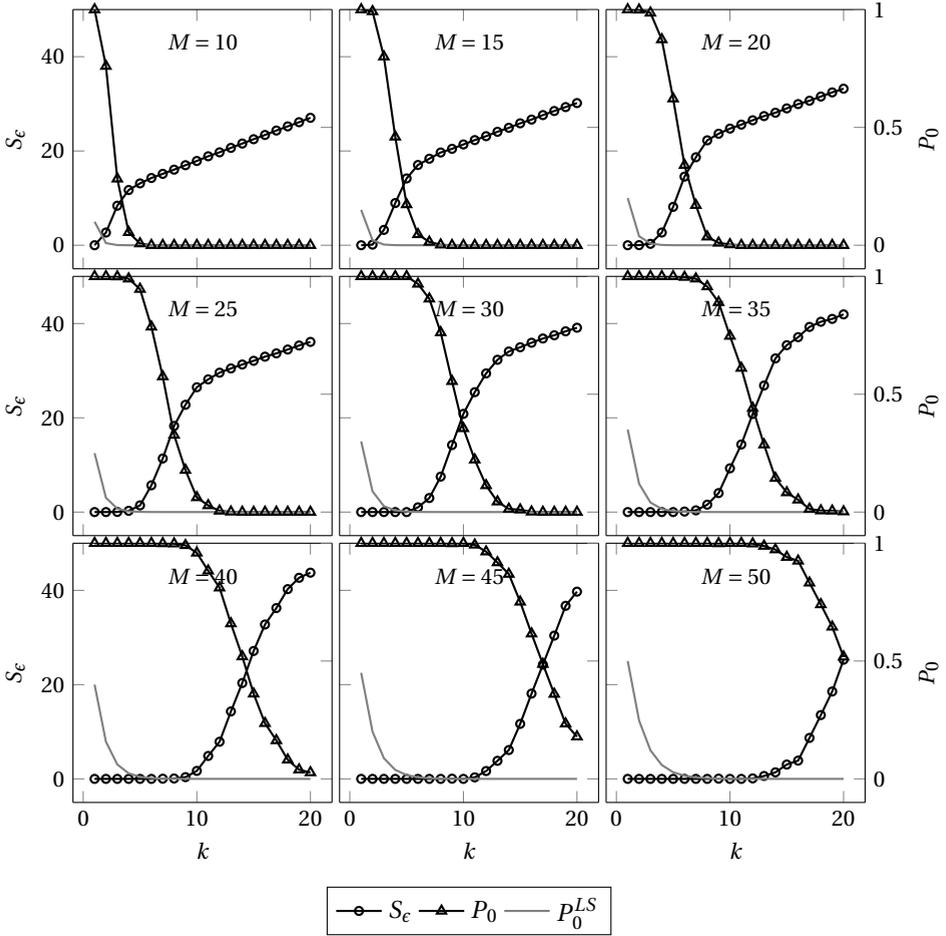


Figure 3.7: Average number of errors S_e and probability of a correct estimate P_0 depending on k , ranging from 1 to 20, $N = 1000$, M ranging from 10 to 50.

increases, the algorithm performance deteriorates and more measurements are needed for a correct estimate.

4. Relation between N , k and M

The previous analyses indicate that there exists a relation between the number of total available parameters, the number of non-zero parameters and the number of available measurements when it comes to correctly estimating the data generating model parameters θ . To investigate the nature of this relation further, an experiment is performed that records the results for ranges of N , k and M . The range of N is $\mathcal{N} = \{10, 20, \dots, 100\}$, and the range of k is $\{1, 2, \dots, n\} \forall n \in \mathcal{N}$. For each combination of N , k the minimum number of measurements $M_{\min}^{(N,k)}$ is determined that is needed to correctly estimate the data generating model parameter vector. The number of trials for each combination is 20.

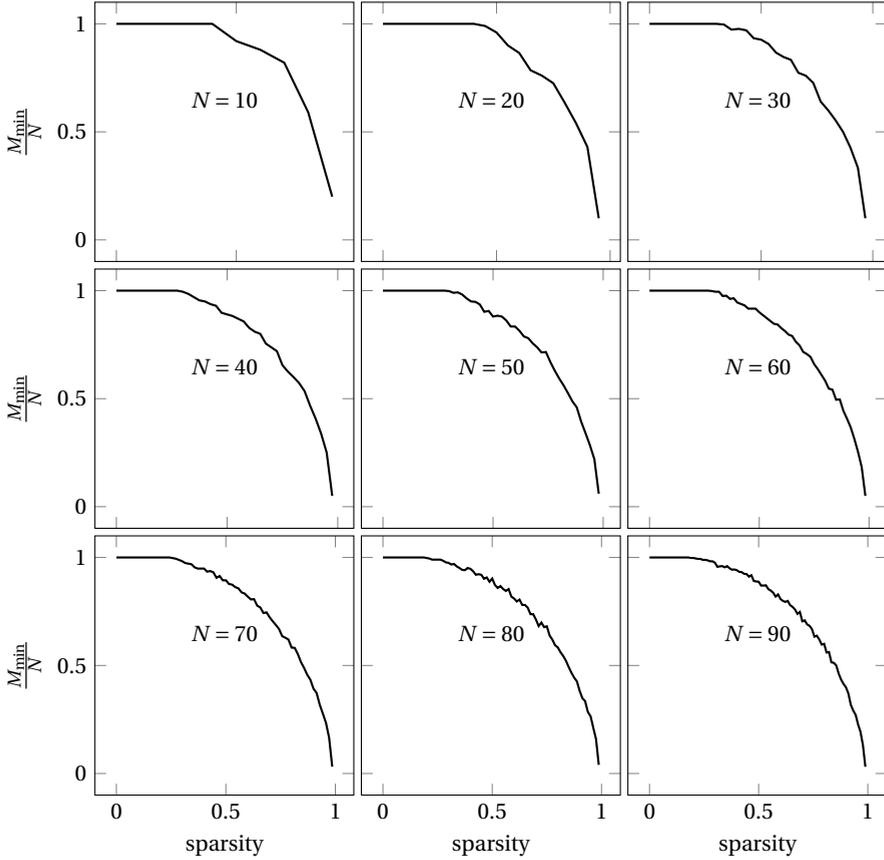


Figure 3.8: The minimum relative number of measurements $\frac{M_{\min}}{N}$ needed to estimate the data generating model parameters correctly, as a function of the sparsity of the parameter vector $(1 - \frac{k}{N})$ for $N \in \{10, 20, \dots, 100\}$ and $k = (1, 2, \dots, N)$. The number of trials for each combination is 20.

The main result of this experiment is depicted in Figure 3.8. It shows that the performance of the algorithms improves with higher sparsity of the parameter vector and that at higher sparsity, a lower amount of measurements is sufficient to compute an accurate estimate. This is in contrast to the strict bound that is given by the threshold based on matrix coherence C , stated in Equation 3.25, as can be seen in Figure 3.9. The theoretical upper bound on the number of non-zero parameters for example does not even reach 3 at $N = 500$, given a relatively large number of measurements $M < N$. A similar observation was made by Donoho [28], who showed that for most large underdetermined systems the parameter vector with minimal ℓ_1 -norm corresponds to the sparsest solution.

5. Dependence on the noise level

In many practical situations, the measurement data Y is corrupted by a $(M \times 1)$

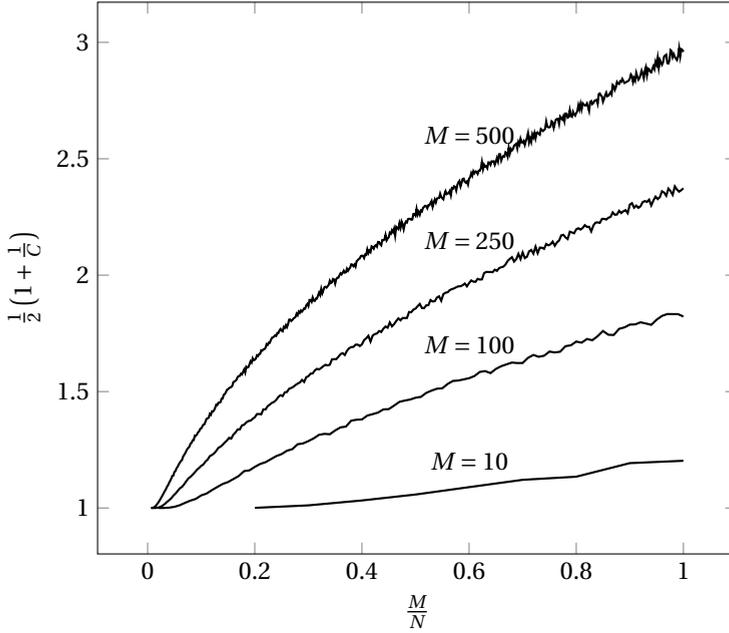


Figure 3.9: Relative number of measurements $\frac{M}{N}$ and the average threshold for $\|\theta\|_0$, based on the measurement matrix coherence C : $\|\theta\|_0 \leq \frac{1}{2} \left(1 + \frac{1}{C}\right)$.

noise vector E , producing the noisy data $\hat{Y} = Y + E$. A number of experiments have been performed to assess the sensitivity of the sparse estimation algorithm to a noise vector E , generated from a Gaussian distribution with zero mean and standard deviation σ . The signal to noise ratio (SNR) is used to express the noise level:

$$\text{SNR} = \frac{\text{Var}(Y)}{\text{Var}(E)} = \frac{\text{Var}(Y)}{\sigma^2}, \quad (3.26)$$

often converted to decibels:

$$\text{SNR} = 10 \log_{10} \left(\frac{\text{Var}(Y)}{\text{Var}(E)} \right) \text{dB}. \quad (3.27)$$

The variance $\text{Var}(Y)$ and $\text{Var}(E)$ can be interpreted here as the power of the “signals” Y and E , since they both have zero mean. Y has zero mean because $Y = \Phi\theta$ and the independent columns Φ_j of the matrix Φ and θ_j both have zero mean. The variance of Y also depends on the values of Φ and θ and can be expressed as

$$\text{Var}(Y) = \frac{7k\alpha^2}{3M}. \quad (3.28)$$

This can be derived as follows: since every column j in Φ has length 1, the energy in $\Phi_j\theta_j$ is θ_j^2 . Therefore, the total energy in Y is $\theta_1^2 + \theta_2^2 + \dots + \theta_N^2 = \|\theta\|^2$. The variance of Y is the expected power: $\text{Var}(Y) = \|\theta\|^2/M$. The expected value of θ_j ,

coming from a uniform distribution on $[-2\alpha, -\alpha] \cup [\alpha, 2\alpha]$, is $\frac{7\alpha^2}{3}$. If only k elements of θ are non-zero, $\text{Var}(Y)$ is $\frac{7k\alpha^2}{3M}$. The required noise variance σ^2 is computed in each experiment to match the desired noise level in decibels. The expected effect of adding noise to the measurements is that more than k elements of the estimated parameter vector $\hat{\theta}$ will be non-zero. In general M elements of $\hat{\theta}$ will be non-zero. When $M \geq k$ and at reasonable noise levels, the original k non-zero parameters will also be non-zero in the parameter estimate $\hat{\theta}$. Figures 3.10 and 3.11 explore the effects of adding noise on the *minimum* estimated value of a non-zero parameter and the *maximum* value of a zero parameter.

A threshold is needed to distinguish the small (and consequently probably zero) non-zero parameters from the dominant non-zero parameters in $\hat{\theta}$. To decide on the computation of the value of this threshold, one can look at how the expected values in Y depend on the values in θ , and the inverse relation, what the expected value of an entry in an estimated $\hat{\theta}$ is based on a data record \hat{Y} . From the computation of $\text{Var}(Y)$ it can be derived that the expected power of a non-zero element of $\hat{\theta}$ is $\frac{M}{N}\text{Var}(\hat{Y})$. At moderate noise levels, originally non-zero parameters are more likely to have a power in $\hat{\theta}$ higher than this value, since they represent the power in Y , while originally zero parameters that are non-zero in $\hat{\theta}$ are more likely to have a power value that is less than $\frac{M}{N}\text{Var}(\hat{Y})$, since they represent the power in E . The square root of $\frac{M}{N}\text{Var}(\hat{Y})$ is therefore a good threshold candidate to determine which non-zero entries in $\hat{\theta}$ are actually zero entries, moreover since it does not depend on any a priori knowledge about the noise level or k . Figures 3.12 and 3.13 show the effect of different noise levels on the performance of the algorithm and the probability of correctly classifying the non-zero and zero parameters at different noise levels. It can be concluded that the algorithm is relatively sensitive to noise, especially at low sparsity. At higher sparsity the negative effect of noise can be diminished by adding more measurements.

6. Dependence on LP solver

The time that is required to compute a solution to a given sparse linear regression problem is a relevant aspect of the algorithm. The two steps of the algorithm can be performed using a number of different solvers. This section concentrates on the choice of linear programming solver. The least squares minimization method is the same in each setup. The LP solvers under investigation are the simplex method, the activate set method and the interior point method as implemented in Matlab release 2010a. The performance of each solver is dependent on the dimensions of the problem (N, k, M) . Figure 3.14 shows that the active set method is the fastest option to solve the ℓ_1 -minimization problem as it outperforms both the simplex method and the interior point method.

Another property that influences the performance of the algorithm is the formulation of the linear programming problem. As explained in Chapter 2, Section 2.2, the problem of minimizing the ℓ_1 -norm of the parameter vector can be stated as a linear programming problem in two equivalent ways: the *primal*

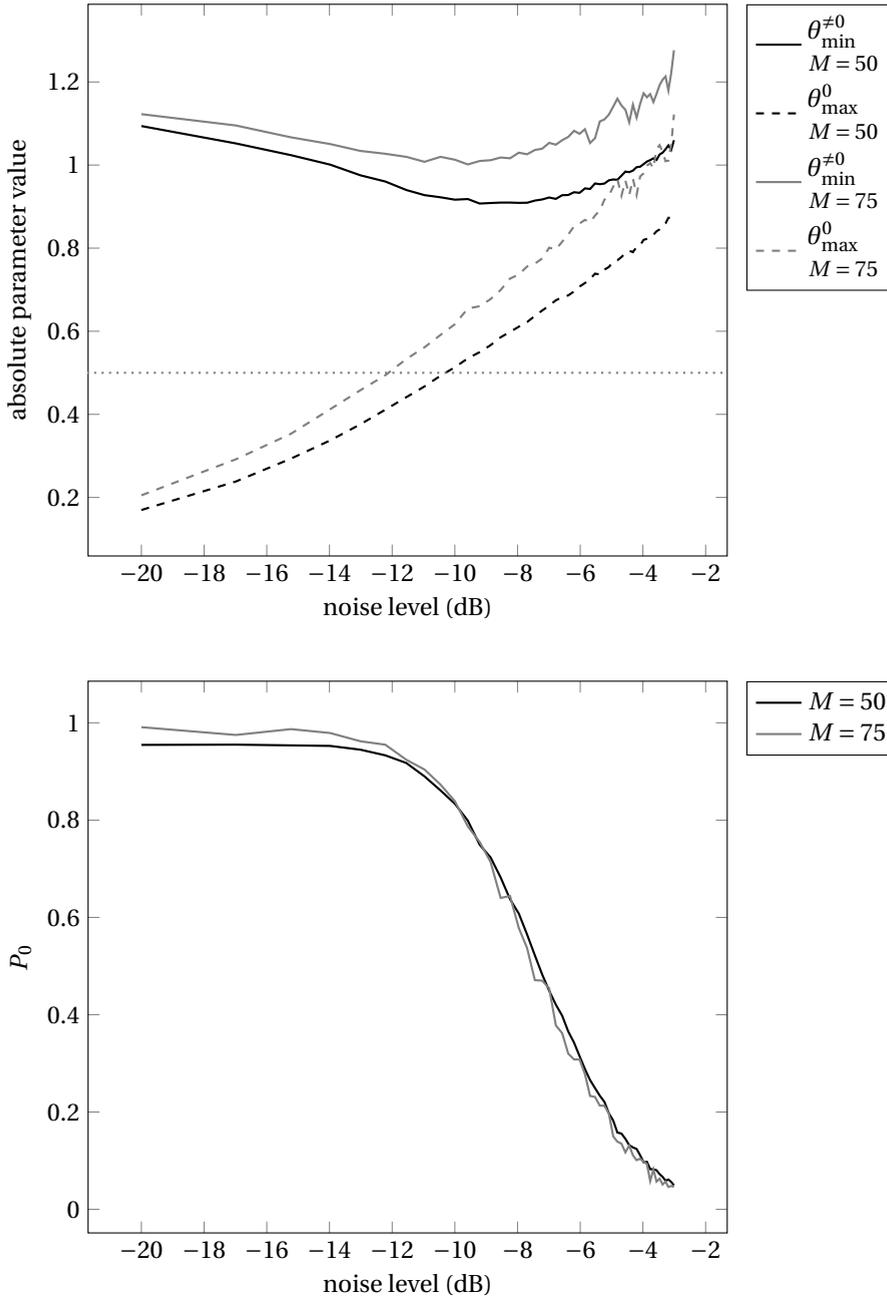


Figure 3.10: Effect of the noise level on the minimum estimated non-zero parameter value ($\theta_{\min}^{\neq 0}$) and the maximum estimated zero parameter value (θ_{\max}^0), and the probability P_0 that the non-zero and zero parameters can be correctly separated. $N = 100$, $k = 5$, $M = 50$ and 75 .

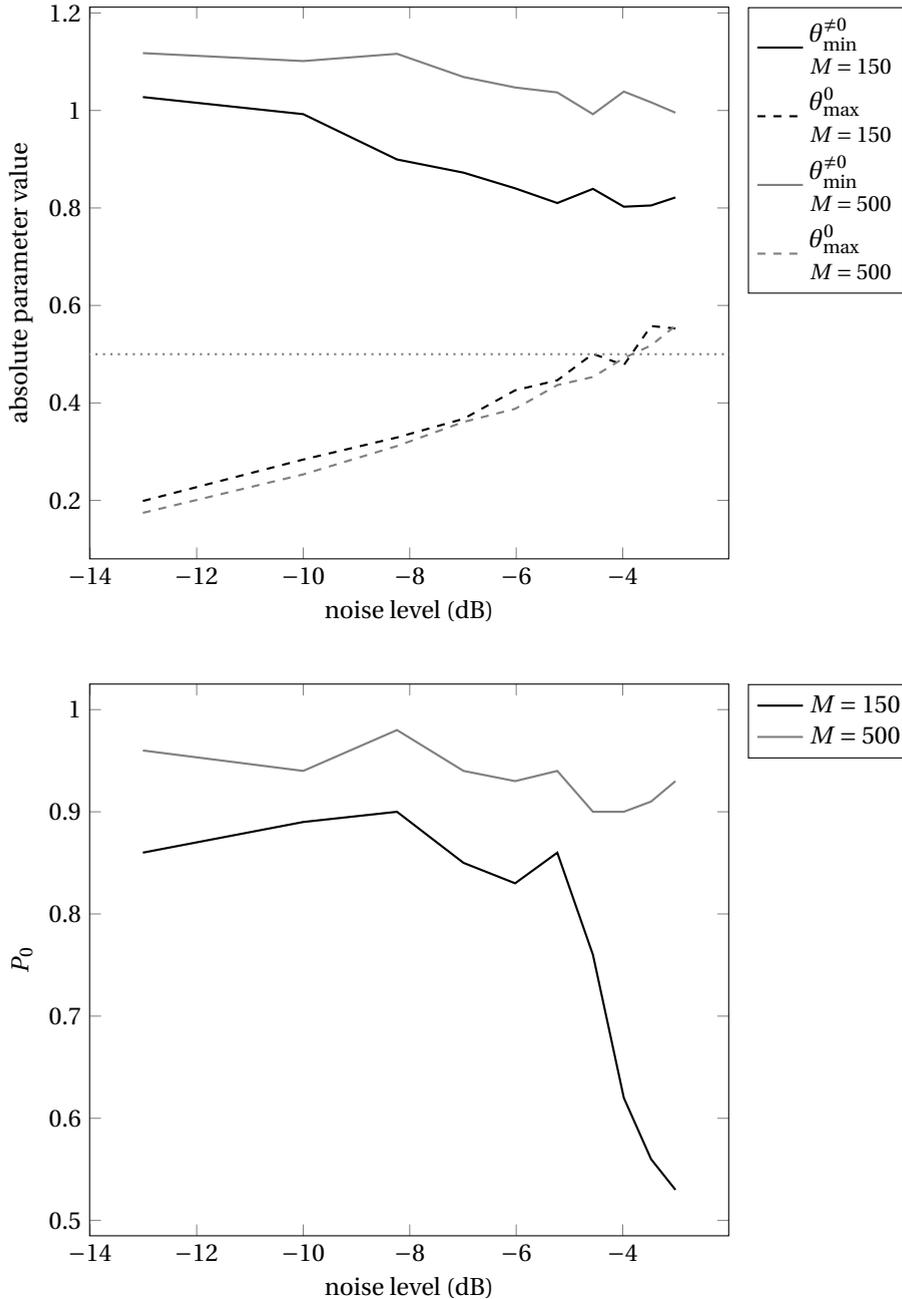


Figure 3.11: Effect of the noise level on the minimum estimated non-zero parameter value ($\theta_{\min}^{\neq 0}$) and the maximum estimated zero parameter value (θ_{\max}^0), and the probability P_0 that the non-zero and zero parameters can be correctly separated. $N = 1000$, $k = 5$, $M = 150$ and 500 .

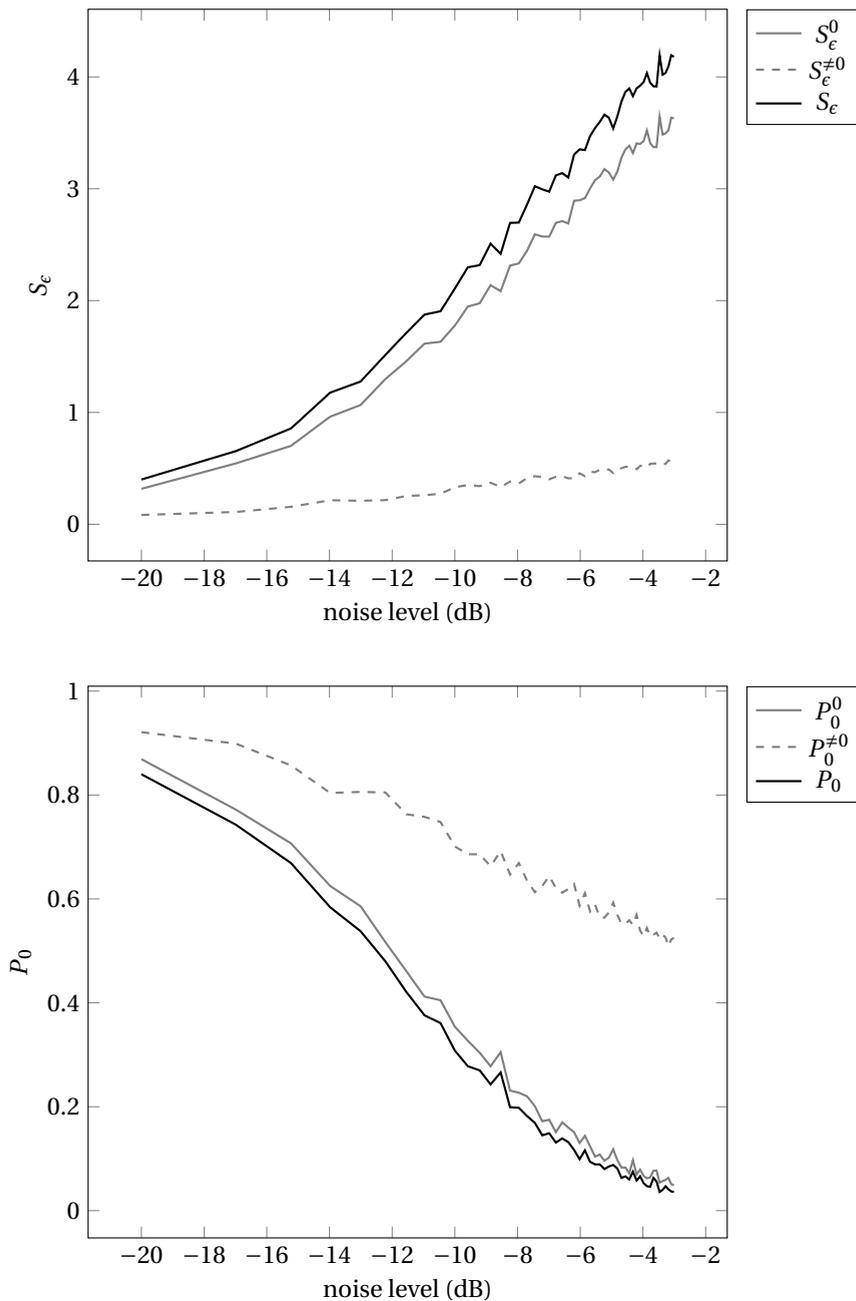


Figure 3.12: Effect of the noise level on the number of errors S_ϵ and the probability P_0 that no parameters are incorrectly classified, both split into zero (0) and non-zero ($\neq 0$) parameters. $N = 10$, $k = 2$ and $M = 9$.

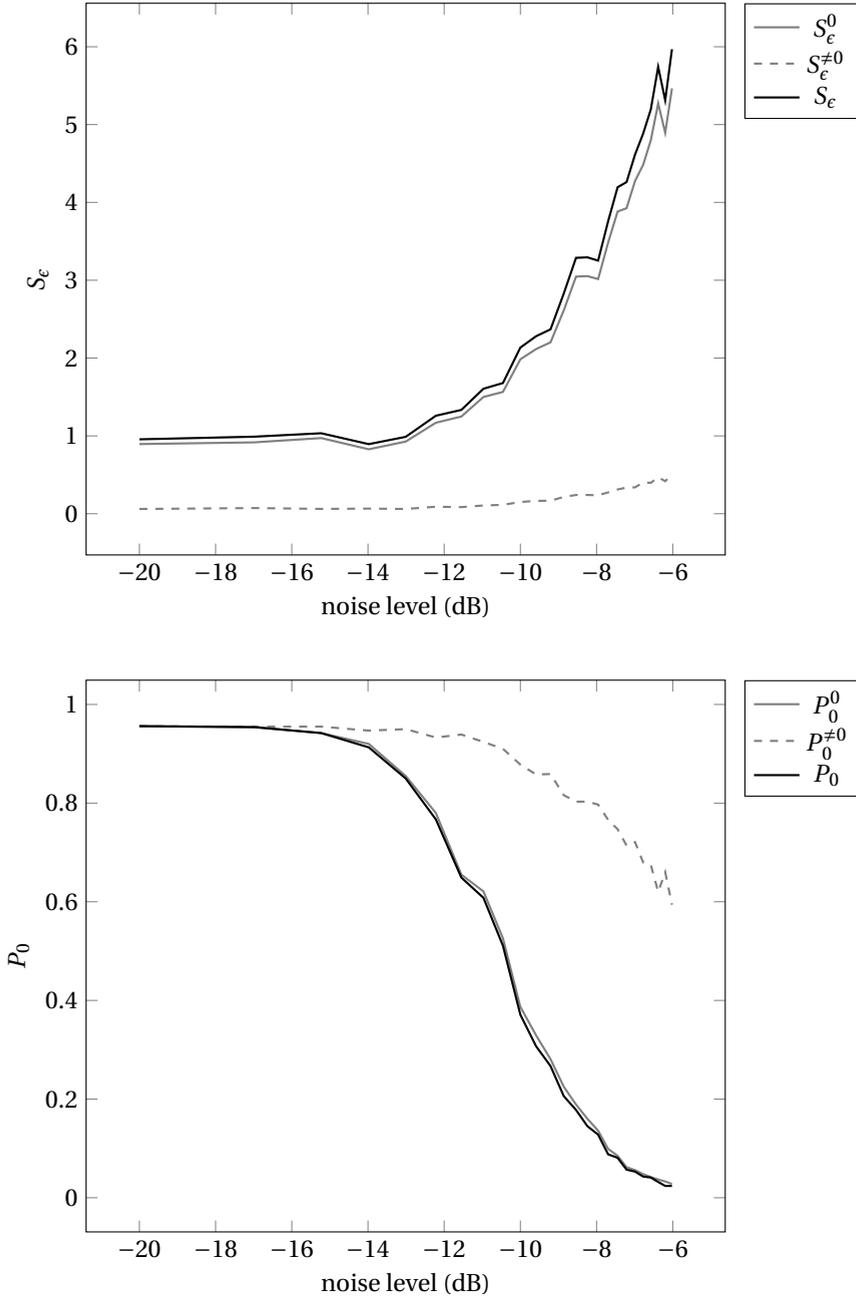


Figure 3.13: Effect of the noise level on the number of errors S_ϵ and the probability P_0 that no parameters are incorrectly classified. $N = 100$, $k = 5$ and $M = 50$.

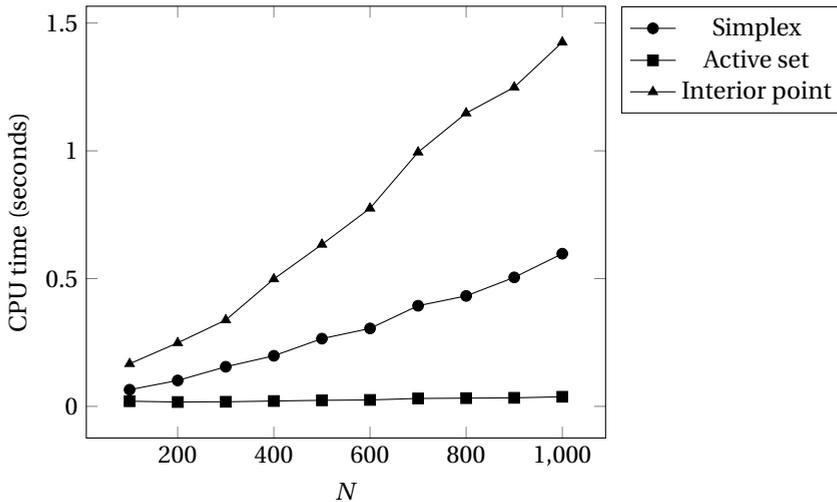


Figure 3.14: Average CPU time of different LP solvers on the ℓ_1 -minimization problem (dual LP formulation, $k = 5$, $M = 50$).

problem (2.15) and the *dual* problem (2.17). Solving the primal problem immediately gives one the solution for the parameter vector with minimal ℓ_1 -norm. Solving the dual problem instead, requires one to translate the solution to the dual problem back to a solution to the primal problem via the relation given in Section 2.2.2. The performance of these two approaches is compared using the active set LP solver. The result can be seen in Figure 3.15. Solving the dual problem first, clearly is much faster than solving the primal problem directly.

3.3.3 Performance of the least absolute shrinkage and selection operator (LASSO)

As mentioned in Section 3.2, another method that favors sparse solution in the linear regression setting, is the lasso method. This method is also applicable in the underdetermined setting investigated here. The formulation in Equation 3.13 can be formulated as a regularized linear regression:

$$\min_{\theta} \|\Phi\theta - Y\|_2 + \lambda \|\theta\|_1, \quad (3.29)$$

and is solved for several values of the regularization parameter λ , using for instance a technique called least angle regression (LARS) [32]. In the same experimental setup as in item 4, the performance of the lasso method was compared to the performance established for the sparse maximization algorithm. Lasso performance was calculated in a slightly different way, since the lasso tends to underestimate the parameter values θ , depending on the choice for λ . Therefore a lasso solution was considered correct if the nonzero parameters *indices* were correctly determined. Figure 3.16 shows the result of the comparison in terms of the minimal number of measurements M_{\min}/M needed

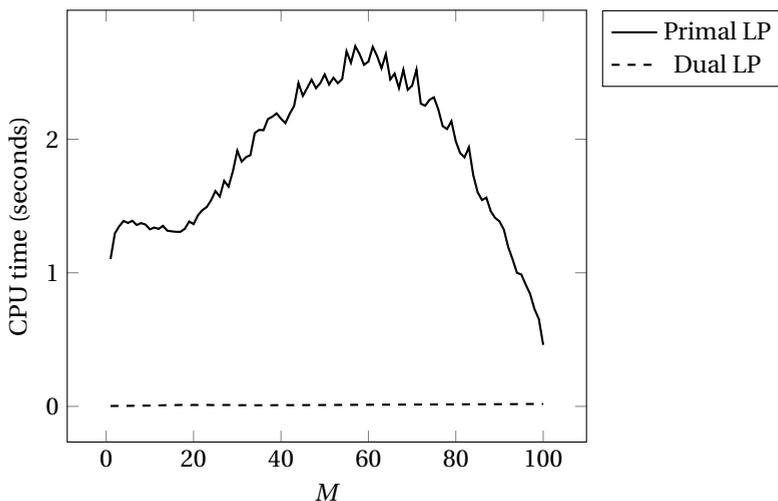


Figure 3.15: Average CPU time for solving the sparse maximization problem with the primal LP formulation and the dual LP formulation ($N = 100$, $k = 5$, active set method).

to correctly identify the data generating model parameters, depending on the model sparsity. In this setup the sparse maximization algorithm is better able to correctly identify the parameters at a lower amount of available measurements. At the same level of sparsity, the lasso method in general requires a larger number of measurements for successful estimation than the sparse maximization algorithm. The elastic net algorithm is left out of the comparison, as its strengths lie mainly in the ability to handle correlated parameters, which are not expected to be present given the experimental setup.

3.3.4 Overdetermined problems and subset selection problems

A common problem is the case where the number of available measurements is (much) higher than the number of parameters ($M \gg N$), but the measurements are corrupted by noise, which reduces the chance to retrieve the original data generating model parameters. In this case, finding a sparse solution corresponds to the *subset selection* problem where one tries to find the *dominant* regressors in the set of available regressors that still provides a good fit to the measurements. In this case there is no clear equivalence space: the equivalence space at the least squares solution θ_{LS} is a singleton, because θ_{LS} is the unique solution. One could try to search within a sub-optimal space for a sparse solution, as is done in the lasso approach. An extension of the current sparse estimation to the overdetermined case is however beyond the scope of this.

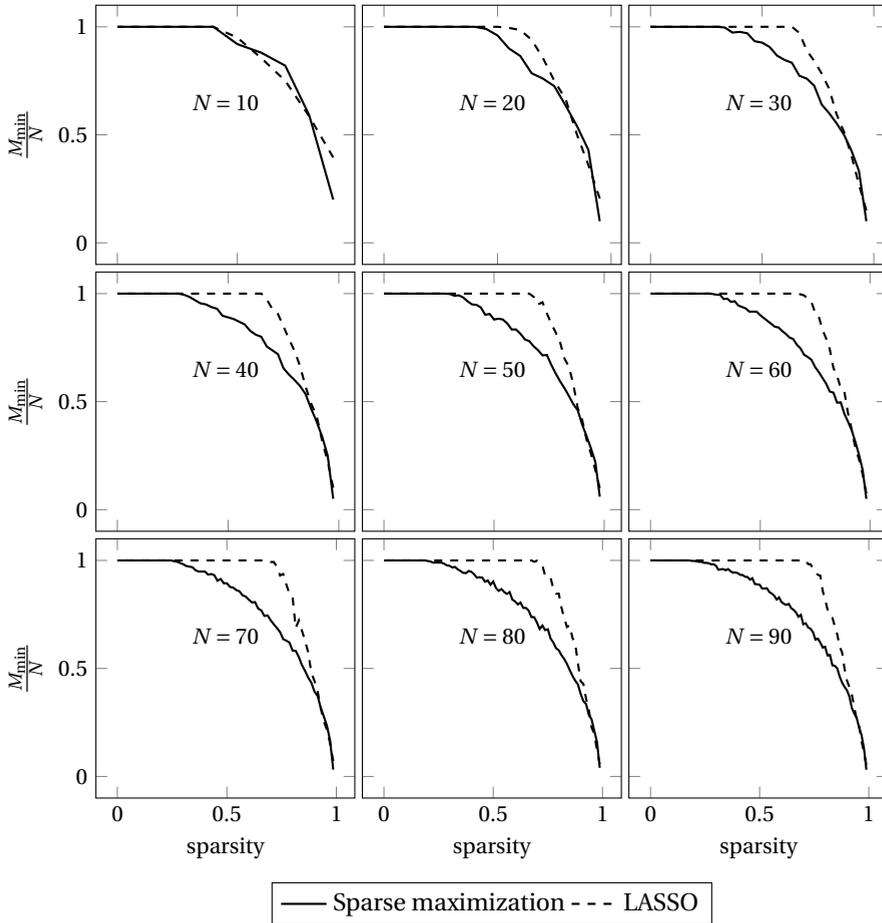


Figure 3.16: Comparison of the minimum relative number of measurements $\frac{M_{\min}}{N}$ needed to estimate the data generating model parameters correctly between the sparse maximization algorithm and the LASSO method, as a function of the sparsity of the parameter vector $(1 - \frac{k}{N})$ for $N \in \{10, 20, \dots, 100\}$ and $k = (1, 2, \dots, N)$. The number of trials for each combination is 20.

3.4 Conclusions

In a linear regression setting the equivalence space that can be exploited for sparse maximization of the parameter vector, is given by the kernel of the regressor matrix Φ . The examples in Section 3.3.2 show that in an underdetermined situation, where the number of available observations M is lower than the number of parameters N that are to be estimated, ℓ_1 -minimization of the parameter vector θ can lead to accurate reconstruction of the maximally sparse data generating parameters. However, it is also possible to design an example where a maximally sparse data generating parameter vector (the vector with minimal ℓ_0 -norm) is not found by ℓ_1 -minimization. Scaling of the columns of the regressor matrix Φ to unit length improves the conditions under which ℓ_1 -minimization leads to correct reconstruction of the maximally sparse solution.

The experiments in Section 3.3 indicate that the number of non-zero parameters k , the number of measurements M and the number of parameters N together determine the accuracy of the reconstruction of the data generating parameters. In general, for a reconstruction to be successful, the number of available measurements needs to exceed the number of non-zero parameters ($M > k$). At high sparsity (k relative to N), the number of measurements needed to get a high probability of accurate reconstruction is relatively low. This is in contrast to the strict bound imposed by a measure based on the mutual coherence of the regressor matrix. Moreover, there is a consistent relationship between k , M and N that can be helpful in determining the likelihood of computing an accurate reconstruction, given the dimensions of the linear regression problem and the assumed sparsity. The main effect of adding measurement noise is that at increasing noise levels the originally zero-valued parameters are assigned a non-zero value, which eventually makes it infeasible to correctly classify non-zero and zero-valued data generating parameters. More measurements are then needed to counteract this effect. When compared to the lasso method, an alternative approach to sparse parameter vector reconstruction, the sparse maximization algorithm shows superior performance.

Chapter 4

Application: State-space models

4.1 Introduction

In this chapter the possibilities and limitations of applying the sparse estimation procedure to the class of linear time-invariant (LTI) state-space models are investigated. State-space models are used in system identification and control for a wide range of applications. One important feature of a state-space model is that it allows one to define a state interaction matrix that describes the strength of the relations between the state components of a system. Here the state components are regarded as nodes in a network and the interaction matrix as the structure of this network. In many networks the number of direct interactions is limited, meaning that every node only interacts with a small number of other nodes, making the interaction matrix sparse. When the number of available measurements is too low for conventional identification methods to be able to estimate the network uniquely, this creates a relevant case to apply the sparse estimation algorithm. In the case of full parameterization and sufficient measurements, the sparse estimation algorithm is also applicable, as in this case there still exists a natural equivalence space of input/output equivalent models that can be searched for parameter sparsity.

The sparse estimation algorithm is employed to identify the network interactions in a number of different settings.

1. *Discrete-time model*: in this case the problem can be formulated as a linear regression problem and solved by sparse linear regression. (Section 4.6)
2. *Discretized continuous-time model*: a continuous state-space model can be transformed to a discrete model under certain conditions. The discrete model can be estimated by sparse linear regression and transformed back again to a continuous model. (Section 4.7)
3. *Continuous-time model*: the parameters of the continuous model are estimated directly, leading to a non-linear optimization problem, that can be solved by an iterative version of the sparse estimation algorithm. (Section 4.8)

The performance and applicability of the sparse estimation algorithm is evaluated for each of these settings. To investigate the properties of the *iterative* sparse estimation algorithm, a different, but related setting is treated first, where the state-space model is fully parameterized (abandoning the network interpretation) and the number of available measurements is sufficient. This is a setting worth investigating, as in this case the algorithm can be tuned in such a way that it stays within the current equivalence space at each sparse maximization step. This property is achieved by relating the space to search for sparsity to the search space in an identification procedure called data-driven local coordinates (DDLCL). When the state-space model is not fully parameterized, the sparse maximization step has to be bounded in some way to stay sufficiently close to the current model equivalence space. The influence of the choice of bound on the algorithm performance is assessed by comparing it to the case in which it is guaranteed that the sparse maximization step stays within the equivalence space. (Section 4.5).

4.2 The model class

The model class is the well-known class of discrete time LTI *state-space models*, described by the equations

$$\begin{aligned}x[k+1] &= A_d x[k] + B_d u[k] + w[k], \\y[k] &= Cx[k] + Du[k] + v[k].\end{aligned}\tag{4.1}$$

Here, at each time instant $k \in \mathbb{Z}$, the n -vector $x[k]$ denotes the state, the m -vector $u[k]$ denotes the exogenous input and the p -vector $y[k]$ denotes the output. The p -vectors $w[k]$ and $v[k]$ are independent, identically distributed random variables, representing process or measurement noise sources respectively. Models from the model class are assumed to be stable, i.e. the matrix A_d has eigenvalues that lie within the open unit disk. It can be shown that every model in (4.1) has an associated model representation in *innovations form*:

$$\begin{aligned}x[k+1] &= A_d x[k] + B_d u[k] + K_d e[k], \\y[k] &= Cx[k] + Du[k] + e[k].\end{aligned}\tag{4.2}$$

where the p -vector $e[k]$ is also an independent and identically distributed (i.i.d.) variable, called the innovations input. In this case the innovations input $\{e[k]\}$ is assumed to constitute a zero mean white noise stationary process with constant covariance $\Sigma_d > 0$; this is the innovations process from which the model representation derives its name. The matrix K_d is the Kalman gain matrix. It is assumed that a record of input-output observations is available with respect to the exogenous input signal $\{u[k]\}$ and the output signal $\{y[k]\}$. This i/o data record can be used to identify the state-space matrices (A_d, B_d, C, D, K_d) . As usual (see for instance [58]), it is further assumed that:

- I minimality holds: the number n of state components used to describe the i/o behavior is as small as possible; equivalently $(A_d, [B_d, K_d])$ is controllable and (C, A_d) is observable,

- II A_d is asymptotically stable: all its eigenvalues are in the open unit disk, ensuring that the noise part of the model is stationary,
- III $(A_d - K_d C)$ is asymptotically stable: the noise part of the model is strictly minimum phase, so that the associated one-step-ahead predictor is asymptotically stable,
- IV for each k the innovation $e[k]$ is independent from the state $x[\ell]$ and the input $u[\ell]$ at time instants $\ell \leq k$, and also from the output $y[\ell]$ at time instants $\ell < k$.

For a further background on this model class and its use in system identification, see [58], [43] and [92]; see also [57] and [72].

4.3 Prediction error identification

To identify the discrete-time system (A_d, B_d, C, D, K_d) and the noise covariance Σ_d from an available record of i/o data, a prediction error method (PEM) can be used. The one-step-ahead predictor associated with an estimate $(\hat{A}_d, \hat{B}_d, \hat{C}, \hat{D}, \hat{K}_d)$ of (A_d, B_d, C, D, K_d) is given by the equations

$$\hat{x}[k+1] = (\hat{A}_d - \hat{K}_d \hat{C})\hat{x}[k] + (\hat{B}_d - \hat{K}_d \hat{D})u[k] + \hat{K}_d y[k], \quad (4.3)$$

$$\hat{y}[k] = \hat{C}\hat{x}[k] + \hat{D}u[k]. \quad (4.4)$$

It gives rise to a *prediction error process* $\{\hat{e}[k]\}$ according to

$$\hat{e}[k] = y[k] - \hat{y}[k] = -\hat{C}\hat{x}[k] - \hat{D}u[k] + y[k]. \quad (4.5)$$

When the one-step-ahead predictor is driven by i/o data obtained from a system (A_d, B_d, C, D, K_d) and a noise covariance Σ_d , it yields a prediction error process of which the covariance $\hat{\Sigma}_d$ satisfies the inequality $\hat{\Sigma}_d \geq \Sigma_d$ (to be understood in the sense of positive semi-definite matrices). Equality occurs if and only if $(\hat{A}_d, \hat{B}_d, \hat{C}, \hat{D}, \hat{K}_d)$ coincides with (A_d, B_d, C, D, K_d) up to a similarity transformation. This implies that the system matrices can be identified when the trace of the prediction error covariance is minimized. This is the rationale behind least squares prediction error methods for system identification, since $\text{trace}\{\hat{\Sigma}_d\} = \mathbb{E}(\|\hat{e}[k]\|_2^2)$. In a practical situation, the latter expectation is conveniently estimated by the trace of the error covariance resulting from a given i/o data set. Equivalently, it involves minimization of the nonlinear least squares criterion of fit

$$V = \frac{1}{M} \sum_{k=1}^M \|\hat{e}[k]\|_2^2 \quad (4.6)$$

over the set of matrices $(\hat{A}_d, \hat{B}_d, \hat{C}, \hat{D}, \hat{K}_d)$, where M denotes the size of the data record. (Here, with some slight abuse of notation, the quantity $\hat{e}[k]$ now denotes the actual prediction error at time k resulting from the available data record, using some initial conditions to start up the prediction error filter.) The criterion V is popular for various other reasons too:

1. it corresponds to maximum likelihood estimation in a situation where the innovations process is Gaussian,
2. local minimization of V can be achieved in a relatively cheap fashion using a Gauss-Newton method (or a more sophisticated variant such as Levenberg-Marquardt, see e.g. [25], [33], [55]),
3. recursive (approximate) methods for minimizing V can be designed, which can be used in an online or an adaptive identification context (see [58], [92]).

One important drawback of the approach is that the prediction error criterion V is notorious for the many local minima it may possess. When a local search method is employed to minimize V , this makes it necessary either to come up with a good initial estimate (such as may be provided by subspace identification methods; cf. [5], [6], [54], [74], [98]) or to use a large number of different initial estimates.

4.4 Parameterization and identifiability issues

Depending on the chosen parameterization of the model class, there are several factors which contribute to identifiability problems that may arise in the prediction error identification framework described above. One important source of unidentifiability is caused by the freedom to choose a basis for the state-space. As is well known, two minimal systems (A, B, C, D, K) and $(\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}, \tilde{K})$ are input-output equivalent (yielding the same transfer function) if and only if there exists a nonsingular $n \times n$ matrix T such that $(\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}, \tilde{K}) = (TAT^{-1}, TB, CT^{-1}, D, TK)$. The dimension of the system manifold is $nm + 2pn + pm$, which is n^2 less than the total number of entries in the system matrices. There are three main ways to deal with this source of unidentifiability:

1. by factoring out the matrix T by using a (local) canonical form (see, e.g., [44], [47], [76], [81]),
2. by using a regularization technique which extends the prediction error criterion V by penalizing the norm of the parameter vector as well (see, e.g., [91]),
3. by using a full parameterization and dealing with the selection of a specific model from its i/o equivalence class only after convergence of the criterion V has been achieved (see, e.g., [62], [63]).

Each of these approaches has its drawbacks when the aim is to identify a sparse model which fits the data well. The use of a local canonical form imposes structure on the system matrices, which need not correspond to a relevant sparse structure. A regularization approach may compromise the quality of the fit between the model and the data, or run into numerical problems when the two parts of the extended criterion become unbalanced. When a full parameterization is used, convergence of the criterion V may become very slow as parameters may drift without affecting V . A recent technique which attempts to deal with these problems concerns the use of *data-driven local coordinates* (DDLC), see for instance [64], [84], [83] and the references given there. Here

the idea is to compute the tangent space to the i/o-equivalence class of the system at hand, and to use a local parameterization with respect to its orthogonal complement (corresponding to the tangent space to the system manifold at the system at hand). As remarked in [102], the search directions generated by the DDLC approach using a Gauss-Newton type method to minimize V , correspond precisely to the search directions obtained in a fully parameterized setting when using a so-called *robust* Gauss-Newton method. Indeed, the tangent space to the i/o-equivalence class of the system at hand is entirely contained in the kernel of the Jacobian matrix J of the error vector $\hat{e} = (\hat{e}[1]^T, \hat{e}[2]^T, \dots, \hat{e}[M]^T)^T$ at that system, because the error vector does not depend on the choice of basis of the state-space. This type of unidentifiability will be examined in Section 4.5.

Another important source of unidentifiability concerns the amount (and the quality) of the available input-output measurement data. In the ideal situation the available data (locally) admits a unique transfer function model (or input-output behavior) which optimally fits the data according to the criterion V . In that case the concept of (local) structural identifiability applies and there are various techniques available to investigate this (see for instance [8], [24] and [59]). In a situation where the number of measurements M is too small, or when the input signals are not sufficiently exciting, there will be a subset of dimension > 0 of models that are all consistent with the data. The focus of this chapter is to investigate whether it is possible to deal with this issue by using additional prior information, in particular information about the sparsity of a model, to select a relevant model from this subset.

4.5 Sparse state-space estimation

The fact that the least squares prediction error criterion of fit in (4.6) is nonlinear in the case of state-space models in innovations form, means that an iterative version of the sparse estimation algorithm has to be applied. Adapted to this situation, at a given estimate with prediction error vector \hat{e} , the two main ingredients of the sparse estimation algorithm are:

1. The orthogonal complement of the kernel of the Jacobian matrix J of \hat{e} is employed as the subspace in which to select a search direction for improving the nonlinear least squares prediction error criterion V .
2. The kernel of J is employed as the subspace in which to select a search direction for improving the ℓ_1 -norm of (a part of) the parameter vector θ . If θ is a (local) optimum of V , the kernel of the Hessian matrix H of the prediction error criterion V is chosen instead.

In the case of state-space models, the parameter vector θ consists of selected entries from the state-space matrices (A, B, C, D, K) . In this section only full parameterizations are considered, but structured parameterizations (as discussed in e.g. [26]) can be dealt with in the same way. The vector θ is defined as

$$\theta = \text{vec}(A, B, C, D, K), \quad (4.7)$$

using the $\text{vec}(\cdot)$ operator to denote the vectorization of a matrix argument by stacking the columns on top of each other. As mentioned before, the two ingredients can be used to define a number of different algorithms, depending on the order in which steps in the two search directions are taken and the number of steps in each search direction. In this chapter the setup of the iterative algorithm is as follows:

1. Minimize the prediction error $V(\theta)$ starting from an initial parameter vector θ_0 , producing a locally optimal vector θ_{ℓ_2}
2. Compute the Hessian of the prediction error criterion V at θ_{ℓ_2} , $H(\theta_{\ell_2})$, and minimize the ℓ_1 norm of the parameter vector within the kernel of $H(\theta_{\ell_2})$. This produces a vector $\tilde{\theta}_{\ell_1}$ with minimum ℓ_1 -norm.
3. Take a bounded step in the direction of $\tilde{\theta}_{\ell_1}$ from θ_{ℓ_2} : $\theta_{\ell_1} = \theta_{\ell_2} + \beta(\tilde{\theta}_{\ell_1} - \theta_{\ell_2})$, with $0 \leq \beta \leq 1$.
4. Set $\theta_0 = \theta_{\ell_1}$ and repeat steps (1), (2) and (3) until the algorithm has converged or a maximum number of iterations has been reached.

This procedure is visualized in Figure 4.1. The choice of β in step (3) is crucial for the performance of the algorithm (see Chapter 2, Section 2.3.2), as it determines the rate of convergence and the monotonicity of the algorithm with respect to the ℓ_1 -norm of θ_{ℓ_1} in the successive steps of the algorithm.

4.5.1 Jacobian and Hessian

To be able to minimize the prediction error of a general discrete-time state-space system in innovations form as defined in (4.2), the Jacobian matrix J has to be computed. The one-step-ahead predictor using the matrices (A_d, B_d, C, D, K_d) and corresponding parameter vector $\theta = \text{vec}(A_d, B_d, C, D, K_d)$ is

$$\hat{x}[k+1] = (A_d - K_d C)\hat{x}[k] + (B_d - K_d D)u[k] + K_d y[k], \quad (4.8)$$

$$\hat{y}[k] = C\hat{x}[k] + Du[k]. \quad (4.9)$$

and the prediction error process:

$$\hat{e}[k] = y[k] - \hat{y}[k] = -C\hat{x}[k] - Du[k] + y[k]. \quad (4.10)$$

The Jacobian $J(\theta)$ of $\hat{e}(\theta)$ is

$$J(\theta) = \begin{bmatrix} J(\theta)_{1,1} & J(\theta)_{1,2} & \dots & J(\theta)_{1,N} \\ J(\theta)_{2,1} & J(\theta)_{2,2} & \dots & J(\theta)_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ J(\theta)_{M,1} & J(\theta)_{M,2} & \dots & J(\theta)_{M,N} \end{bmatrix}, \quad (4.11)$$

with $J(\theta)_{k,i}$ a $p \times 1$ vector containing the partial derivatives of the error vector with respect to the parameter θ_i at time k :

$$J(\theta)_{k,i} = \frac{\partial \hat{e}[k]}{\partial \theta_i} = -\frac{\partial C}{\partial \theta_i} \hat{x}[k] - C \frac{\partial \hat{x}[k]}{\partial \theta_i} - \frac{\partial D}{\partial \theta_i} u[k], \quad (4.12)$$

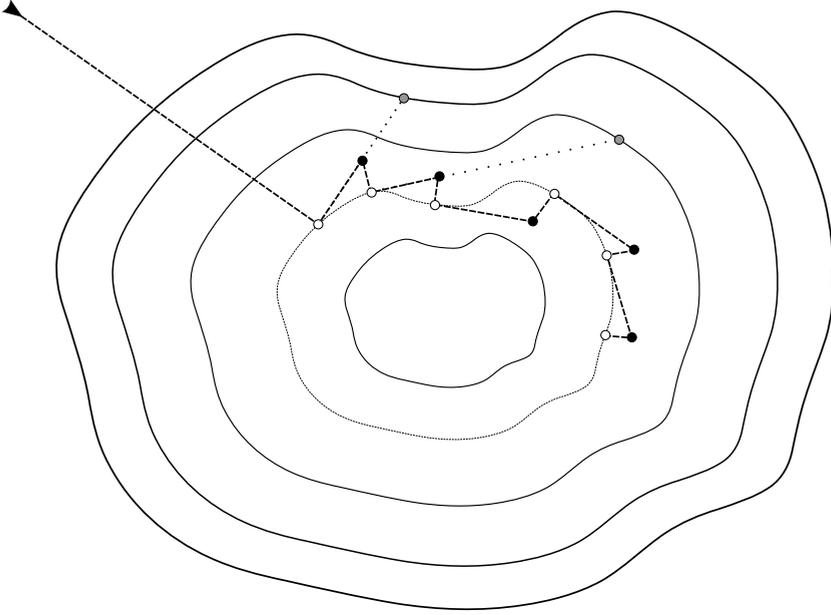


Figure 4.1: An illustration of the sparse estimation iteration scheme. The contours of the optimization criterion V are shown. The contour at the optimal value of V is marked by a dotted line. The white dots correspond to θ_{ℓ_2} solutions, the gray dots correspond to $\tilde{\theta}_{\ell_1}$ solutions, and the black dots correspond to θ_{ℓ_1} solutions. For simplicity, not all $\tilde{\theta}_{\ell_1}$ solutions are visible.

where $\frac{\partial \hat{x}[k]}{\partial \theta_i}$ are the partial derivatives of the state vector with respect to θ_i at time k :

$$\begin{aligned} \frac{\partial \hat{x}[k+1]}{\partial \theta_i} &= \frac{\partial (A_d - K_d C)}{\partial \theta_i} \hat{x}[k] + (A_d - K_d C) \frac{\partial \hat{x}[k]}{\partial \theta_i} \\ &\quad + \frac{\partial (B_d - K_d D)}{\partial \theta_i} u[k] + \frac{\partial K_d}{\partial \theta_i} y[k]. \end{aligned} \quad (4.13)$$

The computation of the *Hessian* matrix $H(\theta)$ of $V(\theta)$ is required to determine the search space for sparsity $\text{Ker}(H(\theta))$:

$$H(\theta) = \begin{bmatrix} H(\theta)_{1,1} & H(\theta)_{1,2} & \dots & H(\theta)_{1,N} \\ H(\theta)_{2,1} & H(\theta)_{2,2} & \dots & H(\theta)_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ H(\theta)_{N,1} & H(\theta)_{N,2} & \dots & H(\theta)_{N,N} \end{bmatrix}, \quad (4.14)$$

with $H(\theta)_{i,j}$ the second order partial derivative of the prediction error criterion $V(\theta)$ with respect to the parameters θ_i and θ_j :

$$H(\theta)_{i,j} = \frac{\partial^2 V}{\partial \theta_i \partial \theta_j} = 2J(\theta)_j^T J(\theta)_i + 2 \sum_k^M \hat{e}[k]^T \frac{\partial^2 \hat{e}[k]}{\partial \theta_i \partial \theta_j}, \quad (4.15)$$

with $\frac{\partial^2 \hat{e}[k]}{\partial \theta_i \partial \theta_j}$ a $p \times N$ matrix consisting of the second order partial derivatives of the error vector with respect to θ_i and θ_j at time k :

$$\frac{\partial^2 \hat{e}[k]}{\partial \theta_i \partial \theta_j} = -\frac{\partial^2 C}{\partial \theta_i \partial \theta_j} \hat{x}[k] - \frac{\partial C}{\partial \theta_i} \frac{\partial \hat{x}[k]}{\partial \theta_j} - \frac{\partial C}{\partial \theta_j} \frac{\partial \hat{x}[k]}{\partial \theta_i} - C \frac{\partial^2 \hat{x}[k]}{\partial \theta_i \partial \theta_j} - \frac{\partial^2 D}{\partial \theta_i \partial \theta_j} u[k]. \quad (4.16)$$

Finally, the second order derivatives of the state vector, $\frac{\partial^2 \hat{x}[k]}{\partial \theta_i \partial \theta_j}$, are given by:

$$\begin{aligned} \frac{\partial^2 \hat{x}[k+1]}{\partial \theta_i \partial \theta_j} &= \frac{\partial^2 (A_d - K_d C)}{\partial \theta_i \partial \theta_j} \hat{x}[k] + \frac{\partial (A_d - K_d C)}{\partial \theta_i} \frac{\partial \hat{x}[k]}{\partial \theta_j} \\ &\quad + \frac{\partial (A_d - K_d C)}{\partial \theta_j} \frac{\partial \hat{x}[k]}{\partial \theta_i} + (A_d - K_d C) \frac{\partial^2 \hat{x}[k]}{\partial \theta_i \partial \theta_j} \\ &\quad + \frac{\partial^2 (B_d - K_d D)}{\partial \theta_i \partial \theta_j} u[k] + \frac{\partial^2 K_d}{\partial \theta_i \partial \theta_j} y[k]. \end{aligned} \quad (4.17)$$

In a discrete-time setting the parameter vector θ consists of entries of the matrices (A_d , B_d , C , D , K_d), which means that second order derivatives of single matrices with respect to θ are always zero, and (4.16) simplifies to

$$\frac{\partial^2 \hat{e}[k]}{\partial \theta_i \partial \theta_j} = -\frac{\partial C}{\partial \theta_i} \frac{\partial \hat{x}[k]}{\partial \theta_j} - \frac{\partial C}{\partial \theta_j} \frac{\partial \hat{x}[k]}{\partial \theta_i} - C \frac{\partial^2 \hat{x}[k]}{\partial \theta_i \partial \theta_j}, \quad (4.18)$$

and (4.17) can be reduced to

$$\begin{aligned} \frac{\partial^2 \hat{x}[k+1]}{\partial \theta_i \partial \theta_j} &= \left(-\frac{\partial K_d}{\partial \theta_i} \frac{\partial C}{\partial \theta_j} - \frac{\partial C}{\partial \theta_i} \frac{\partial K_d}{\partial \theta_j} \right) \hat{x}[k] + \frac{\partial (A_d - K_d C)}{\partial \theta_i} \frac{\partial \hat{x}[k]}{\partial \theta_j} \\ &\quad + \frac{\partial (A_d - K_d C)}{\partial \theta_j} \frac{\partial \hat{x}[k]}{\partial \theta_i} + (A_d - K_d C) \frac{\partial^2 \hat{x}[k]}{\partial \theta_i \partial \theta_j} \\ &\quad + \left(-\frac{\partial K_d}{\partial \theta_i} \frac{\partial D}{\partial \theta_j} - \frac{\partial D}{\partial \theta_i} \frac{\partial K_d}{\partial \theta_j} \right) u[k]. \end{aligned} \quad (4.19)$$

If the state-space system is in fact a continuous-time system but only a discrete-time i/o data record with sample time T is available, the Jacobian and Hessian have to be computed in a different way to be able to relate the change in the discrete-time prediction error to the continuous-time model parameters θ^c . The continuous-time state equations are assumed to be of the form

$$\dot{x}(t) = A_c x(t) + B_c u(t), \quad (4.20)$$

deliberately omitting the stochastic part of the model to avoid having to model continuous-time disturbances. Process noise is therefore assumed to be absent in continuous-time. Under the zero order hold assumption for the input u ($u(t) = u[k]$, $kT \leq t \leq kT + T$ for given sample time T), the corresponding discrete-time state equations are given by

$$x[k+1] = e^{A_c T} x[k] + \left(\int_{kT}^{kT+T} e^{A_c(kT+T-\tau)} B_c d\tau \right) u[k]. \quad (4.21)$$

The conversion from continuous-time matrices A_c and B_c to discrete-time matrices A_d and B_d and vice versa can be performed using a relation described by Van Loan [97]:

$$e^{M_c T} = M_d, \quad (4.22)$$

where

$$M_c = \begin{bmatrix} A_c & B_c \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad \text{and} \quad M_d = \begin{bmatrix} A_d & B_d \\ \mathbf{0} & I_m \end{bmatrix}. \quad (4.23)$$

The first and second order partial derivatives of A_d and B_d with respect to the continuous-time parameter θ_i^c and θ_j^c can be computed relatively straightforward using the following definitions:

$$\begin{aligned} M_d^{(i)} &= \frac{\partial M_d}{\partial \theta_i^c} = \frac{\partial e^{M_c T}}{\partial \theta_i^c}, & M_c^{(i)} &= \frac{\partial M_c}{\partial \theta_i^c}, \\ M_d^{(i,j)} &= \frac{\partial^2 M_d}{\partial \theta_i^c \partial \theta_j^c} = \frac{\partial^2 e^{M_c T}}{\partial \theta_i^c \partial \theta_j^c}, & M_c^{(i,j)} &= \frac{\partial^2 M_c}{\partial \theta_i^c \partial \theta_j^c}. \end{aligned} \quad (4.24)$$

and the general formula for the computation of the directional derivative for any analytical function F on a block triangular matrix (see [69] for a detailed description of the computation of such a derivative):

$$F\left(\begin{bmatrix} P & V \\ \mathbf{0} & P \end{bmatrix}\right) = \begin{bmatrix} F(P) & D_V(F(P)) \\ \mathbf{0} & F(P) \end{bmatrix}, \quad (4.25)$$

where P is a square matrix and $D_V(F(P))$ denotes the first directional derivative evaluated at P in the direction of V . The first order partial derivatives $M_d^{(i)}$ can therefore be computed by setting P to M_c and V to the derivative of M_c with respect to θ_i^c :

$$\exp\left(\begin{bmatrix} M_c & M_c^{(i)} \\ \mathbf{0} & M_c \end{bmatrix}\right) = \begin{bmatrix} M_d & M_d^{(i)} \\ \mathbf{0} & M_d \end{bmatrix}, \quad (4.26)$$

where $\exp(\cdot)$ is the exponential function. The second order derivative $M_d^{(i,j)}$ can be computed by taking the left triangular block matrix in (4.26) as the matrix P in (4.25) and setting V to the derivative of P with respect to θ_j^c :

$$\exp\left(\begin{bmatrix} M_c & M_c^{(i)} & M_c^{(j)} & M_c^{(i,j)} \\ \mathbf{0} & M_c & \mathbf{0} & M_c^{(j)} \\ \mathbf{0} & \mathbf{0} & M_c & M_c^{(i)} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & M_c \end{bmatrix}\right) = \begin{bmatrix} M_d & M_d^{(i)} & M_d^{(j)} & M_d^{(i,j)} \\ \mathbf{0} & M_d & \mathbf{0} & M_d^{(j)} \\ \mathbf{0} & \mathbf{0} & M_d & M_d^{(i)} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & M_d \end{bmatrix}. \quad (4.27)$$

Note that $M_c^{(i,j)}$ is actually a zero matrix in this case, but it is shown here for completeness.

The ability to compute Jacobian and Hessian matrices, discrete-time or continuous-time, is an advantage when applying the sparse estimation algorithm, as opposed to having to approximate them, for two reasons: a) the matrices are accurate (up to numerical precision), and b) the computation can be done more efficiently, which is especially relevant when the problem size increases.

4.5.2 Relation to DDLC

Data driven local coordinates propose a way of parameterizing state-space systems to overcome the first source of unidentifiability mentioned in Section 4.4, that originates from the ‘built-in’ equivalence class present in fully parameterized state-space systems: given the class $\mathbb{M}(n)$ consisting of rational and causal $p \times (p \times m)$ transfer functions with MacMillan degree n , and the parameter space $S_{\min}(n)$ consisting of matrix entries of minimal state-space representations of order n and identical p and m , every transfer function from $\mathbb{M}(n)$ corresponds to an analytical manifold of dimension n^2 in $S_{\min}(n)$. This means that in a system identification procedure at each stage there are at least n^2 possible local parameter vector coordinates that do not change the value of the prediction error criterion V . The DDLC approach avoids this equivalence space by defining a parameterization from a given estimate (A, B, C, D, K) that only considers the orthogonal complement to the tangent space to the i/o equivalence class of the current estimate. The approach is briefly explained here to illustrate the connection to the sparse estimation algorithm; a more detailed discussion and applications to system identification can be found in [81] and [82].

The equivalence class for a given minimal realization (A, B, C, D, K) is given by

$$\{(TAT^{-1}, TB, CT^{-1}, D, TK)\}, \quad (4.28)$$

with T from the set of non-singular $(n \times n)$ matrices. The tangent space to the equivalence class is given by

$$\{(\dot{T}AT^{-1} - TAT^{-1}\dot{T}T^{-1}, \dot{T}B, -CT^{-1}\dot{T}T^{-1}, 0, \dot{T}K), \dot{T} \in \mathbb{R}^{n \times n}\}, \quad (4.29)$$

which at the current minimal realization (A, B, C, D, K) reduces to

$$\{(\dot{T}A - A\dot{T}, \dot{T}B, -C\dot{T}, 0, \dot{T}K), \dot{T} \in \mathbb{R}^{n \times n}\}, \quad (4.30)$$

by setting $T = I$. The tangent space to the i/o equivalence class can be written in vectorized form:

$$\{(\dot{T}A - A\dot{T}, \dot{T}B, -C\dot{T}, 0, \dot{T}K), \dot{T} \in \mathbb{R}^{n \times n}\} = \left\{ \underbrace{\begin{pmatrix} A^T \otimes I_n - I_n \otimes A \\ B^T \otimes I_n \\ -I_n \otimes C \\ 0_{pm \times n^2} \\ K^T \otimes I_n \end{pmatrix}}_Q \cdot \text{vec}(\dot{T}), \dot{T} \in \mathbb{R}^{n \times n} \right\} \quad (4.31)$$

where \otimes stands for the Kronecker product. The tangent space is the span of the columns of Q . The data driven local coordinates are a mapping onto the orthogonal complement of the tangent space, a space that is the span of the columns of the orthogonal complement Q^\perp of the matrix Q .

For the case of a full parameterization and ‘rich’ data, the sparse maximization space $\text{Ker}(J(\theta))$ coincides with the tangent space to the i/o-equivalence class of the

current estimate (this is the space that is deliberately disregarded in the DDLC approach). One may then attempt to stay within the i/o -equivalence class (and thus leave V unchanged) by associating a state-space transformation matrix T with the ℓ_1 -minimization search direction vector s^* . Note that if one puts $T = I_n + \dot{T}$ where \dot{T} represents a ‘first order change’ to the identity matrix I_n , then the corresponding first order change in (A, B, C, D, K) is given by $(\dot{A}, \dot{B}, \dot{C}, \dot{D}, \dot{K}) = (\dot{T}A - A\dot{T}, \dot{T}B, -C\dot{T}, 0, \dot{T}K)$. The latter expression corresponds to the search direction s^* in the same way that the entries of (A, B, C, D, K) are collected in the parameter vector θ . From a given value of s^* a corresponding matrix T can then be computed in a straightforward way by solving the corresponding linear equations for \dot{T} . The matrix \dot{T} can be computed by solving:

$$s^* = Q \cdot \text{vec}(\dot{T}). \quad (4.32)$$

Here it is noted that it is not difficult to design examples where $I + \dot{T}$ becomes singular (and hence cannot be used as a state-space transformation matrix). To resolve this and to avoid singularity one may instead consider the invertible state-space transformation matrix $\exp(\dot{T})$ which exhibits the same first order effects as $I_n + \dot{T}$.

4.5.3 Example: small non-sparse models and full parameterization

To illustrate the iterative version of the sparse estimation algorithm, this example features a discrete-time model \mathcal{M}_0 with 4 states, 1 input and 1 output. This model only contains non-zero parameters: all parameter values are within the interval $[-2\alpha, -\alpha] \cup [\alpha, 2\alpha]$, with α set to 0.25. Starting from a zero initial state, 100 samples are generated using a random Gaussian input signal. No process or measurement noise is added. Table 4.1 shows the state-space matrices of the original model \mathcal{M}_0 and the estimation result $\hat{\mathcal{M}}$ after 100 iterations. The algorithm started at the data generating model \mathcal{M}_0 . The bound β was chosen in such a way that the least squares criterion deteriorated at most by 0.01 when minimizing the ℓ_1 criterion:

$$\max \beta \quad \text{s.t. } V(\theta_{\ell_1}) - V(\theta_{\ell_2}) \leq 0.01, \quad 0 \leq \beta \leq 1. \quad (4.33)$$

The estimated state-space matrices in Table 4.1 contain more zero parameters than the original matrices. Note that the estimate of the matrix K is effectively a zero matrix, which is to be expected since no noise was present in the simulation. The trajectory of the ℓ_1 -norm of the estimated parameter vector θ during the iterations of the algorithm is shown in Figure 4.2. The mixing of ℓ_2 -norm optimization of the error vector and ℓ_1 -norm optimization of the parameter vector causes the ℓ_1 -norm to fluctuate during iterations as depicted in the lower plot in Fig. 4.2. Often, the prediction error minimization step of the algorithm causes an increase in the ℓ_1 -norm of the parameter vectors, as can be seen in Figure 4.3. To stop or decrease the size of these fluctuations at the end of the optimization process and thus ensure convergence, an appropriate damping mechanism has to be constructed.

Fig. 4.4 demonstrates the need for monitoring what happens to V when taking steps to minimize the ℓ_1 -norm of θ . The first plot in Fig. 4.4 displays the values for the ℓ_1 -norm of θ in a situation where all the system estimates are in the same optimal

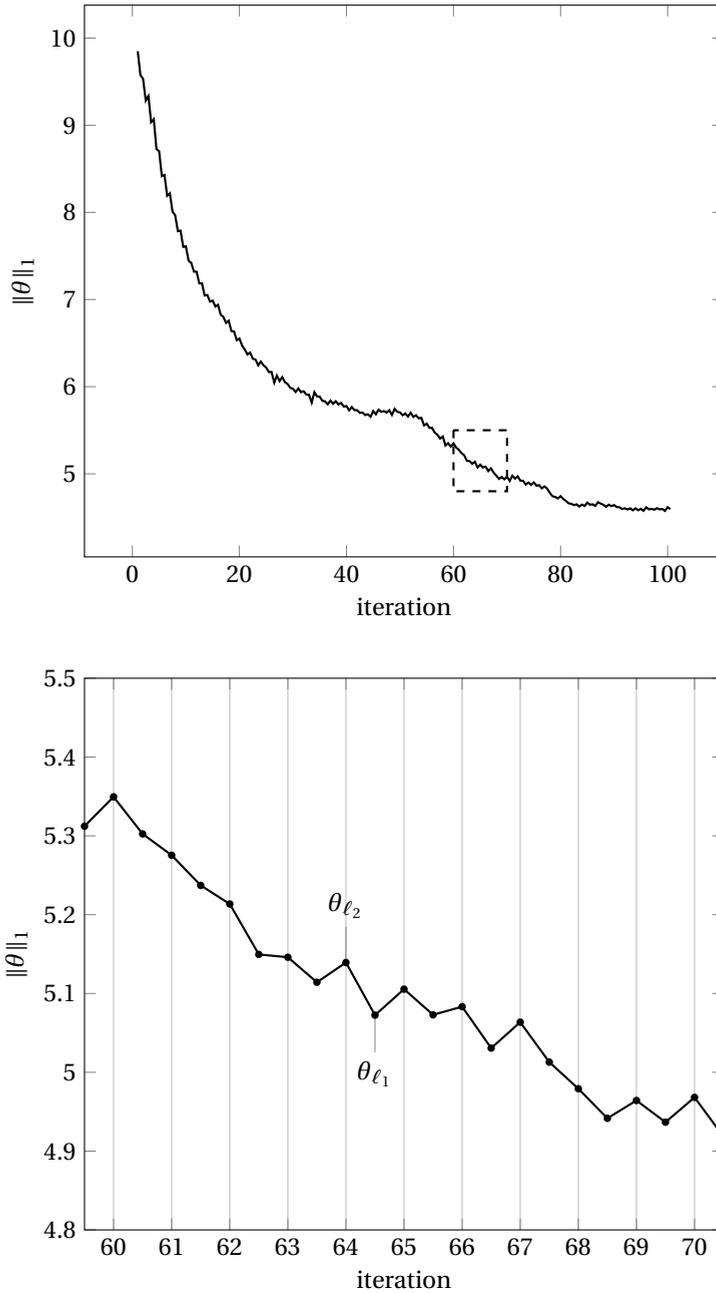


Figure 4.2: ℓ_1 -norm trajectory of the parameter vector estimate θ in the upper plot and a detail of the trajectory in the lower plot, indicating the consecutive ℓ_1 -norm values of θ_{ℓ_2} and θ_{ℓ_1} . It can be seen that minimization of the prediction error criterion V may increase the ℓ_1 -norm of the estimated parameter vector.

Matrix	\mathcal{M}_0	$\hat{\mathcal{M}}$
A	$\begin{pmatrix} 0.29 & -0.31 & 0.33 & -0.36 \\ -0.37 & 0.37 & 0.43 & -0.38 \\ 0.38 & -0.27 & 0.47 & 0.25 \\ 0.33 & 0.25 & 0.30 & 0.34 \end{pmatrix}$	$\begin{pmatrix} 0.67 & -0.14 & 0.01 & -0.04 \\ -0.14 & 0.69 & -0.01 & 0.00 \\ -0.01 & -0.00 & 0.06 & 0.69 \\ 0.02 & -0.01 & -0.27 & 0.04 \end{pmatrix}$
B	$\begin{pmatrix} 0.31 \\ 0.49 \\ 0.30 \\ -0.27 \end{pmatrix}$	$\begin{pmatrix} -0.04 \\ 0.16 \\ -0.06 \\ -0.49 \end{pmatrix}$
C	$(-0.29 \quad -0.33 \quad 0.34 \quad -0.37)$	$(0.00 \quad -0.13 \quad 0.53 \quad -0.01)$
D	-0.28	-0.28
K	$\begin{pmatrix} 0.38 \\ 0.34 \\ 0.39 \\ -0.34 \end{pmatrix}$	$\begin{pmatrix} 0.05 \\ 0.01 \\ 0.00 \\ -0.02 \end{pmatrix}$

Table 4.1: Example of the state-space matrices estimated by the sparse estimation algorithm starting from a non-sparse discrete-time state-space model with 4 states, 1 input and 1 output.

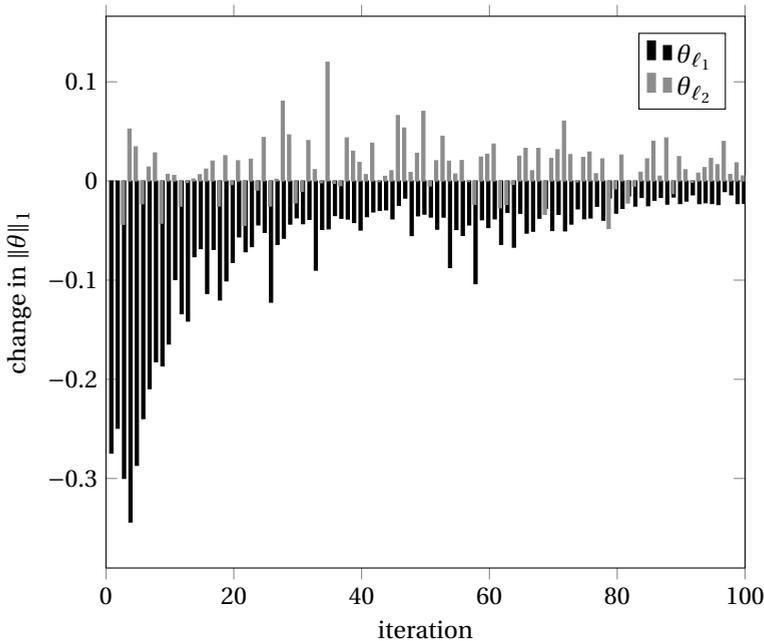


Figure 4.3: Change in the ℓ_1 -norm of the parameter vector θ for 100 iterations. Black bars denote the change after the minimization of the ℓ_1 -norm of the parameter vector, gray bars indicate the ℓ_1 -norm change after prediction error minimization.

i/o-equivalence class. In this case the step s^* in the direction of ℓ_1 -norm minimization is used to compute an associated state-space transformation matrix T as discussed in (4.5.2). This setup is used to minimize the ℓ_1 -norm of θ via state-space transformations. This ensures that the least squares criterion V remains unchanged. The figure shows that when the bound on the least squares criterion value is increased, the speed of ℓ_1 -norm minimization also increases. When the bound on the maximum change in V is too loose, the value of the ℓ_1 -norm of θ cycles between a number of values and does not converge. What happens is that the sparsity maximization step results in a new estimate that does not lie within the i/o-equivalence space of the current estimate. It can easily lead to a non-minimal system, which is explained by the fact that the ℓ_1 -norm minimization tends to create zeros in the parameter vector. The consecutive least squares minimization step takes a relatively large step to return to a local optimum of the least squares error criterion, potentially causing a large deterioration of the ℓ_1 norm of the parameter vector estimate. In this way the algorithm keeps traversing different i/o equivalence spaces and does not converge. Note on the other hand that a tight bound may cause the algorithm to become slow. The trade-off is to choose a bound that forces ℓ_1 -norm minimization steps to remain close to the manifold of i/o equivalent state-space systems, so that the ℓ_2 -step will return to the same equivalence manifold, while maintaining a certain speed of convergence in the ℓ_1 criterion.

4.6 Discrete-time network models

The dynamics of a *discrete-time* linear interaction network can be described by a system of difference equations

$$\begin{aligned} x[k+1] &= A_d x[k] + B_d u[k], \\ y[k] &= x[k] + e[k] \end{aligned} \quad (4.34)$$

with node states x and network inputs u . The network states can be directly observed in the output y , possibly disturbed by a noise term e . This is equivalent to setting $C = I_n$, $K = 0$ and $D = 0$ in the innovations form (4.2). The interaction matrix A_d is assumed to be sparse, the input matrix B_d not. The output equation in Equation (4.34) can be written as:

$$Y = A_d X + B_d U + E, \quad (4.35)$$

with $X = (x[0], x[1], \dots, x[M-1])$, $U = (u[0], u[1], \dots, u[M-1])$, $Y = (x[1], x[2], \dots, x[M])$ and $E = (e[1], e[2], \dots, e[M])$. Each row a_i of A_d and b_i of B_d can be estimated separately by sparse linear regression by solving the *partial* sparsity problem:

$$\begin{aligned} \min_{(a_i, b_i)} \|a_i\|_1 \\ \text{s.t. } (X^T U^T) \begin{bmatrix} a_i^T \\ b_i^T \end{bmatrix} &= \text{proj}_{(X^T U^T)} y_i^T, \end{aligned} \quad (4.36)$$

where y_i is the i -th row of Y and $\text{proj}_{(X^T U^T)} y_i^T$ denotes the orthogonal projection of the vector y_i^T on the column space of $(X^T U^T)$. In the noiseless case, this projection

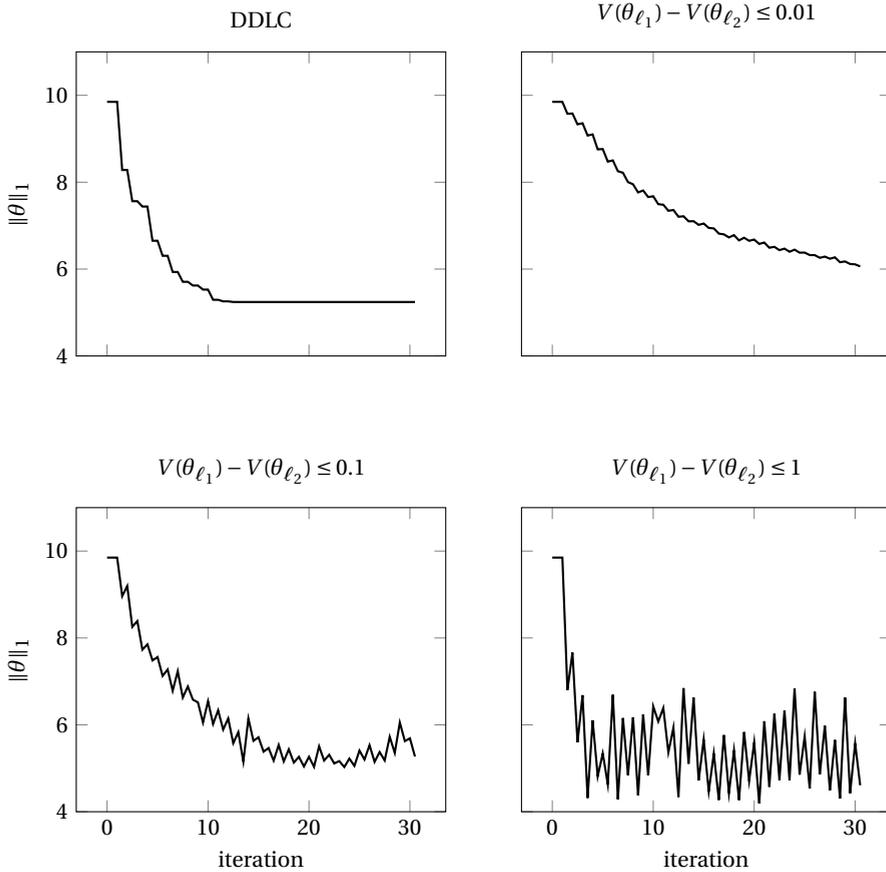


Figure 4.4: ℓ_1 -norm trajectory for different settings for the maximum allowed change in the prediction error criterion value when maximizing the parameter vector sparsity.

can be omitted. Note that not all network structures with a sparse connectivity matrix can be solved in parallel. For this approach to be meaningful, every *row* of A_d has to be sparse. A similar setup, but then in a continuous-time setting with known state derivatives, has been investigated in [77], and the main results for numerous simulations can be found there. In the following experiments this approach is extended to *structured* discrete-time networks that are stable and minimal. These networks are the data generating systems and the estimation problem is to retrieve the correct network interaction matrices from a limited set of available measurements.

4.6.1 Experiments

Experiments are conducted using a directed network with a regular network structure, in this case a ring structure, where each node is connected to itself and at least one neighbor.

The node connectivity matrix has a tridiagonal structure with one additional entry

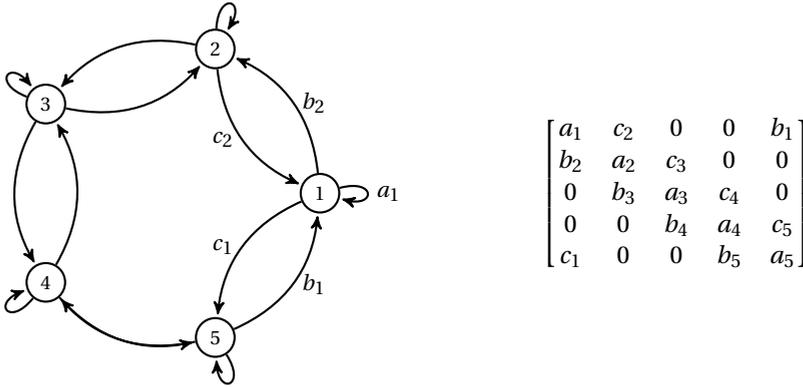


Figure 4.5: An example of the ring shaped network structure of order 5 and the accompanying node interaction matrix

in the upper right corner (b_1) and one in the lower left corner (c_1) to ensure the closure of the ring structure as depicted in Figure 4.5. To generate a stable node interaction matrix from this prescribed structure the following strategy was applied:

1. For a network of order N , the pole locations p_1, p_2, \dots, p_N of the transfer function of the desired system are the roots of the characteristic polynomial of the interaction matrix A . In the case of a discrete-time system they are chosen within the complex unit circle, and in case of a continuous-time data generating system in the left half of the complex plane. The characteristic polynomial of the interaction matrix A can be written as

$$P(x) = (x - p_1)(x - p_2) \cdots (x - p_N). \quad (4.37)$$

2. The polynomial $P(x)$ defined by these pole locations is then disturbed by a (small) number ϵ to alter the roots of the characteristic polynomial. The roots of the disturbed polynomial $\tilde{P}(x) = P(x) + \epsilon$ are denoted as $\tilde{p}_1, \tilde{p}_2, \dots, \tilde{p}_N$. $\tilde{P}(x)$ is factored into linear factors (real roots) and quadratic factors (pairs of complex conjugated roots). These factors are used to construct the intermediate matrix \tilde{A} that has the characteristic polynomial $\tilde{P}(x)$.

$$\tilde{P}(x) = P(x) + \epsilon = (x - \tilde{p}_1)(x - \tilde{p}_2) \cdots (x - \tilde{p}_N). \quad (4.38)$$

As an example, for $N = 5$, with one real root \tilde{p}_1 and two pairs of complex roots (\tilde{p}_2, \tilde{p}_3) and (\tilde{p}_4, \tilde{p}_5), the matrix \tilde{A} is constructed as follows:

$$\tilde{A} = \begin{bmatrix} \boxed{a_1} & -0 & 0 & 0 & 0 \\ 0 & \boxed{a_2} & c_3 & 0 & 0 \\ 0 & \boxed{b_3} & a_3 & 0 & 0 \\ 0 & 0 & 0 & \boxed{a_4} & c_5 \\ 0 & 0 & 0 & \boxed{b_5} & a_5 \end{bmatrix} \begin{cases} a_1 = \tilde{p}_1 \\ (x - a_2)(x - a_3) - c_3 b_3 = (x - \tilde{p}_2)(x - \tilde{p}_3) \\ a_2 + a_3 = \tilde{p}_2 + \tilde{p}_3, \quad c_3 = \frac{\tilde{p}_2 \tilde{p}_3 - a_2 a_3}{b_3} \\ (x - a_4)(x - a_5) - c_5 b_5 = (x - \tilde{p}_4)(x - \tilde{p}_5) \\ a_4 + a_5 = \tilde{p}_4 + \tilde{p}_5, \quad c_5 = \frac{\tilde{p}_4 \tilde{p}_5 - a_4 a_5}{b_5} \end{cases}$$

The values of (a_2, a_3) and (a_4, a_5) are to be chosen in such a way that their sum equals $(\tilde{p}_2 + \tilde{p}_3)$ and $(\tilde{p}_4 + \tilde{p}_5)$. The values of c_3 and c_5 depend on the choices for b_3 and b_5 respectively.

3. The vector b of length N is defined in such a way that $\prod_{i=1}^N b_i = -c$. This vector forms the sub diagonal of the matrix A (with b_2, b_3, \dots, b_N) and the entry $A_{1,N}$ (with b_1). The unknown values in the matrix \tilde{A} can be computed based on the values in b . Setting the remaining variables in the matrix A to zero ($c_1 = 0$ and $c_i = 0$ if it is not involved in a quadratic factor), results in a matrix with characteristic polynomial $P(x)$. In the case of $N = 5$, the matrix will have the following form:

$$A = \begin{bmatrix} a_1 & \boxed{0} & 0 & 0 & b_1 \\ b_2 & a_2 & c_3 & 0 & 0 \\ 0 & b_3 & a_3 & \boxed{0} & 0 \\ 0 & 0 & b_4 & a_4 & c_5 \\ \boxed{0} & 0 & 0 & b_5 & a_5 \end{bmatrix}$$

A network with this prescribed structure will attain a high level of sparsity, even at relatively small network sizes N . Figure 4.6 illustrates this property. At $N = 20$ a sparsity of 0.84 is achieved, and at $N = 50$ the sparsity level is already at 0.93

This procedure will generate a stable network in theory. One consideration however is the observation made by Wilkinson [103], that the location of the roots of a given polynomial can be sensitive to small changes in the coefficients of the polynomial. The implications of this observation are twofold in this case. First of all, the roots of the polynomial $P(x)$ have to be chosen carefully to create a well-conditioned polynomial in which the small disturbance ϵ will cause a small change in the root locations. The approach taken here is to place the roots of the polynomial $P(x)$ within equidistance on the complex circle with radius 0.9. This ensures that the roots are well-separated and of equal magnitude. A second implication is that the computation of the characteristic polynomial needs to be as precise as possible. To avoid changes in the coefficients of the polynomial due to roundoff errors, computations are carried out with 64 decimal digits numeric precision.

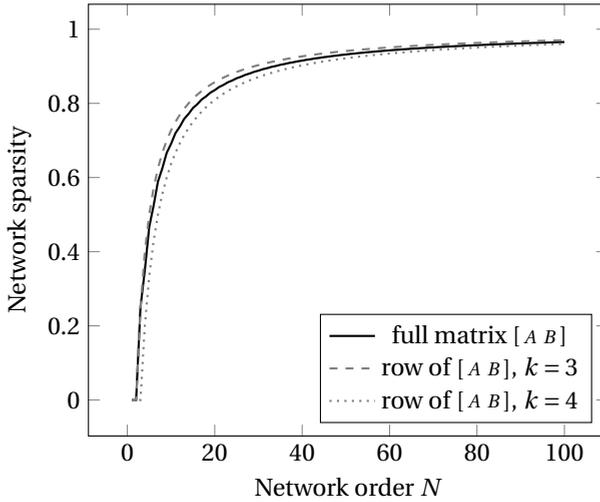


Figure 4.6: Sparsity of the network with the tridiagonal interaction matrix structure for increasing size N of the network. The computation of the network sparsity includes the (non-sparse) matrix B . The sparsity of the system described by the matrices A and B is defined by $\lfloor 3\frac{1}{2}N - 1 \rfloor / (N^2 + N)$. Sparse estimation of the *rows* of A and B , applicable in the case of discrete-time networks, can lead to a slightly sparser (with 3 non-zero parameters) or less sparse setting (with 4 non-zero parameters), with sparsity equal to $3/(N + 1)$ and $4/(N + 1)$ respectively.

The inputs u are Gaussian signals with zero mean and unit variance. The initial states $x[0]$ are also drawn from a Gaussian distribution with zero mean and unit variance. The motivation behind choosing this type of inputs and initial state is that with sufficiently exciting inputs and non-zero initial states the dynamics of the model are more likely to be well represented in the output vector y . The performance of the sparse estimation algorithm is measured by the number of errors and the probability that the correct network model (with a correct A and B) is retrieved.

4.6.2 Results

The example in Figure 4.7 demonstrates that the sparse linear regression algorithm is able to reconstruct the correct network structure in an underdetermined setting. The performance of the sparse linear regression algorithm for a range of network sizes is shown in Figures 4.8 and 4.9. Figure 4.8 illustrates the effect of increasing network size N and available number of measurements M on the sparse network estimation performance, measured in false negatives and false positives. The amount of false negatives (parameters that have a non-zero value, but are estimated as zero) starts off high at a low amount of available measurements M . Typically at $M = 1$, the number of false negatives is almost equal to the number of non-zero parameters. The number of false negatives steadily drops to zero as more measurements are available. The number of false positives (parameters that have a zero value, but are estimated as having a non-

zero value), starts out low but then quickly increased with increasing M , before dropping again to zero. Both the number of false negatives and positives produced by the sparse estimation algorithm clearly outperform a conventional least squares solution. The effect of increasing M on the two error types and the probability of a completely correct estimate P_0 is depicted in Figure 4.9.

4.7 Discretized continuous-time network models

The description of *continuous-time* network consists of a system differential equations

$$\begin{aligned}\dot{x}(t) &= A_c x(t) + B_c u(t) \\ y(t) &= x(t) + e(t).\end{aligned}\tag{4.39}$$

If the values of $\dot{x}(t)$ are available, the sparse estimation problem again reduces to sparse linear regression. However, in general these values are not known and can be approximated at best. An alternative approach is to try and estimate a discrete-time network first, based on the available measurements, and convert the result to a continuous-time description. The discrete-time zero order hold (ZOH) equivalent model of (4.39) is

$$\begin{aligned}x[k+1] &= e^{A_c T} x[k] + \left(\int_{kT}^{kT+T} e^{A_c(kT+T-\tau)} B_c d\tau \right) u[k] \\ y[k] &= x[k] + e[k],\end{aligned}\tag{4.40}$$

with sample time T so that $x[k] = x(kT)$. Note that although A_c is assumed to be sparse, $A_d = e^{A_c T}$ is in general not sparse for all values of T . To deal with this issue, the value of T has to be sufficiently small, as can be appreciated from the first order approximation of A_d :

$$A_d = e^{A_c T} = \sum_{k=0}^{\infty} \frac{1}{k!} (A_c T)^k \approx I + A_c T.\tag{4.41}$$

From this approximation can also be concluded that the diagonal elements of A_d will typically have non-zero values. They should therefore not be selected for sparsity maximization. The adapted sparse linear regression problem, where sparsity is not enforced in the diagonal of the interaction matrix A , is given by

$$\begin{aligned}\min_{(a_i, b_i)} & \|\tilde{a}_i\|_1 \\ \text{s.t.} & (X^T U^T) \begin{bmatrix} a_i^T \\ b_i^T \end{bmatrix} = \text{proj}_{(X^T U^T) y_i^T},\end{aligned}\tag{4.42}$$

where \tilde{a}_i denotes the i -th row of the matrix A_d without the entry A_{ii} . After estimation the discrete-time matrices A_d and B_d need to be translated to their continuous-time counterparts by the inverse of the transformation described in Equation 4.5.1:

$$M_c = \frac{\log M_d}{T},\tag{4.43}$$

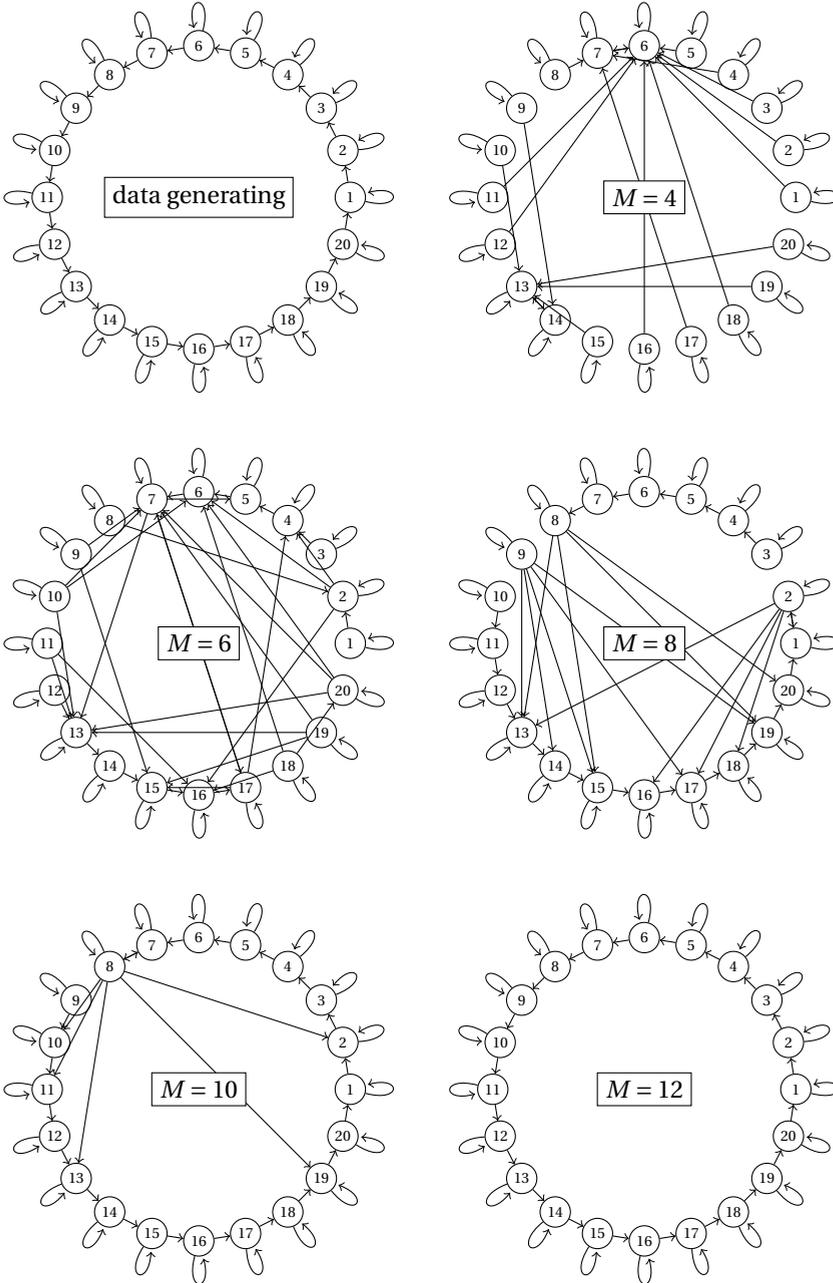


Figure 4.7: Data generating network structure and estimated network structure based on $M = \{4, 6, 8, 10, 12\}$ measurements.

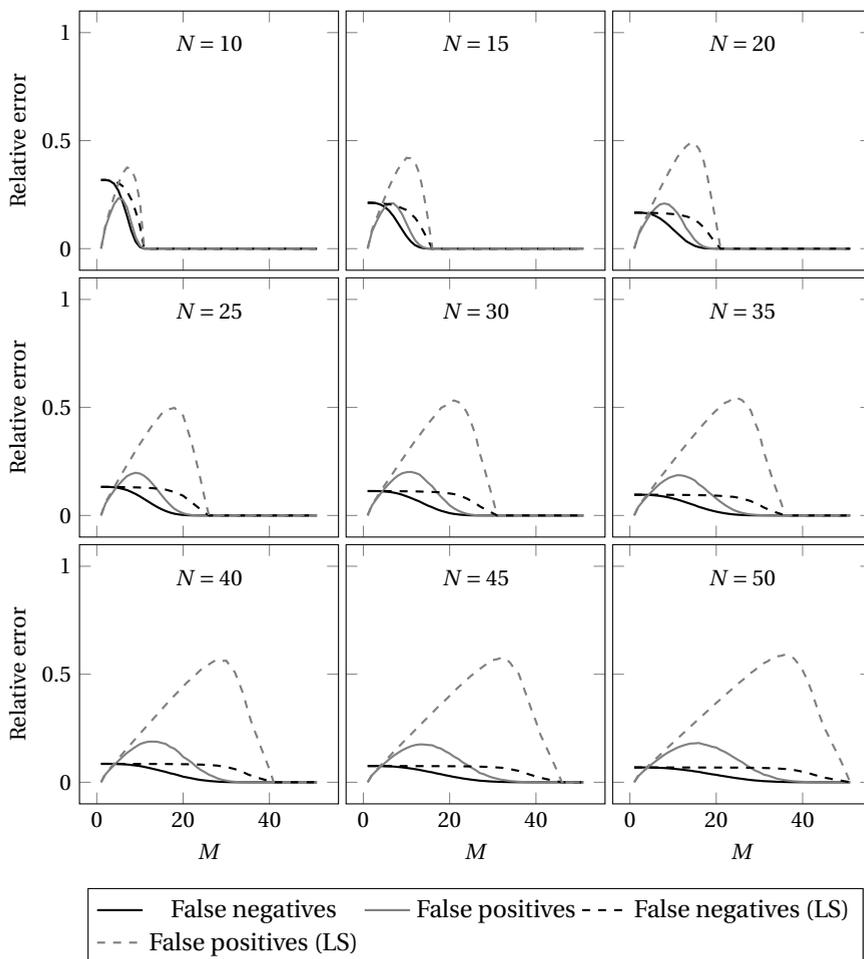


Figure 4.8: Average relative error when retrieving the discrete-time network interaction matrices A and B using sparse linear regression without prior information on sparsity structure of the matrices, for networks of size $N = \{10, 15, \dots, 50\}$. The solid black lines indicate the fraction of false negatives (non-zero parameters that are estimated as zero) and the solid gray lines represent the fraction of false positives (zero parameters that are estimated as non-zero parameters). The dotted lines show the relative error made when computing a standard least squares (LS) solution. Results are averaged over 100 networks for each size N , with 10 random input vectors for each input size M .

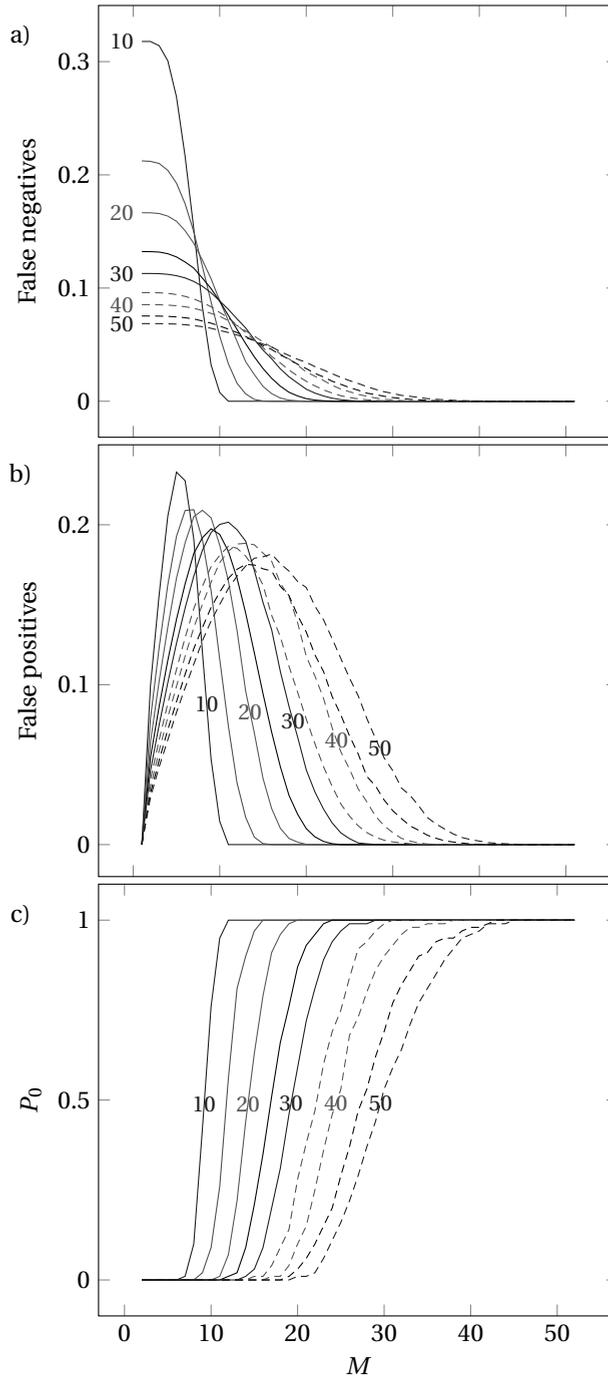


Figure 4.9: Average fraction of a) false negative and b) false positive errors, and c) the probability of a correct estimate P_0 for different network sizes N , ranging from 10 to 50.

where

$$M_c = \begin{bmatrix} A_c & B_c \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad \text{and} \quad M_d = \begin{bmatrix} A_d & B_d \\ \mathbf{0} & 1 \end{bmatrix}. \quad (4.44)$$

An potential issue with this approach of first estimating the parameters of the network in discrete time and then translating those parameters to a continuous-time settings, is that there is no guarantee that the matrix logarithm featured in Equation 4.43 will yield a real-valued matrix M_c and consequently, real-valued matrices A_c and B_c . First of all, the matrix M_d needs to be invertible for $\log M_d$ to exist. Secondly, a real solution exists if and only if each Jordan block of M_d associated with a negative eigenvalue occurs an even number of times. A unique solution M_c exists if and only if all the eigenvalues of M_d are real and positive and no Jordan block belonging to any eigenvalue appears more then once [23]. It is clear that the proposed estimation procedure does not adhere to any of these conditions.

4.7.1 Experiments

The way the continuous-time networks are generated is similar to the way described in 4.6.1, with the only difference being the location of the poles of the generated system. The poles are generated in the left side of the complex plane to ensure that the continuous-time system is stable. The choice of the exact location of the poles is guided by the desired location of the poles in the transformed discrete-time system with sampling time T . The relation between pole locations in discrete-time and continuous time is given by the bilinear transform:

$$z_{1,2} = e^{s_{1,2}T}, \quad (4.45)$$

where $s_{1,2} = x \pm yi$ denotes the pair of complex conjugated continuous-time poles. This means that the magnitude of the discrete-time poles $z_{1,2}$ is determined by the real part of $s_{1,2}$ and the sampling time T :

$$|z_{1,2}| = e^{xT}, \quad (4.46)$$

and the angle of the discrete-time poles is determined by the imaginary part of $s_{1,2}$ and the sampling time T :

$$\angle z_{1,2} = \pm e^{yiT}. \quad (4.47)$$

Figure 4.10 shows the relationship of the choice for the location of the continuous-time poles and the transformed discrete-time poles at different sampling intervals T . It is clear that in this example the value of T is critical for the location of the discrete-time poles and consequently, the ability of any algorithm to correctly estimate the original continuous-time pole locations. A choice of a large value for T will reduce the discrete-time pole magnitude and potentially give rise to an discrete-time pole angle outside of the $\langle -\pi, \pi \rangle$ range, which in turn can lead to a pole multiplicity larger than one. On the other hand, a small value of T will push the discrete-time pole magnitude closer to 1, but may also reduce the separation between the discrete-time poles to such an extent that they become indistinguishable for practical purposes.

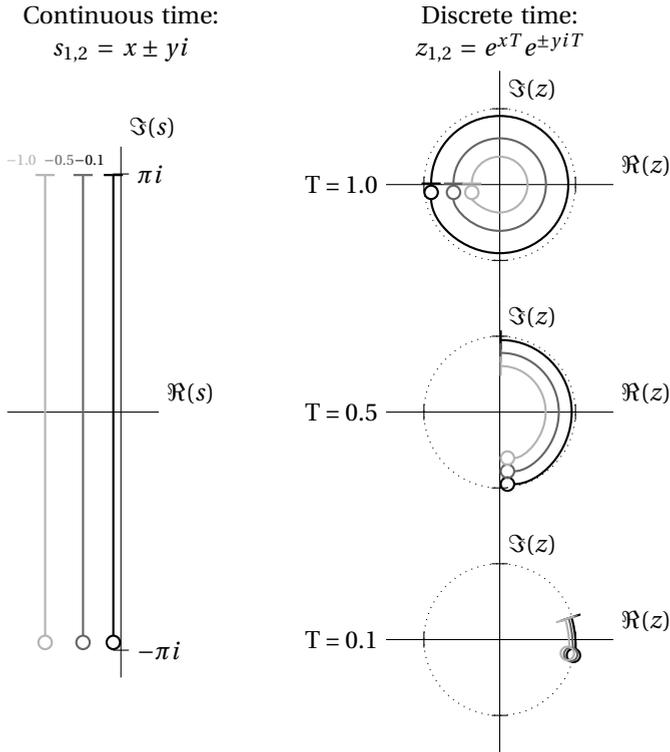


Figure 4.10: Bilinear transform of the continuous- and discrete-time pole locations for sampling time $T = \{0.1, 0.5, 1\}$. The (stable) continuous time pole locations $s_{1,2} = x \pm yi$ with real parts $x = \{-0.1, -0.5, -1\}$ and imaginary parts yi in the interval $\langle -\pi, \pi \rangle$ are transformed to discrete-time pole locations $z_{1,2} = e^{xT} e^{\pm yiT}$.

4.7.2 Results

The results of the estimation of a continuous-time network parameters via a discrete-time representation are depicted in Figures 4.11 and 4.12. It is clear that the performance deteriorates compared to the result in discrete-time as shown in Figure 4.9. There are several explanations for this lower performance:

1. Under- or over-sampling: at lower sampling rates, for instance $T = 1.0$, the dynamics of the continuous-time system may not be sufficiently expressed in the measurement data, especially when the amount of available data is limited. At higher sampling rates the time covered by the measurements may in turn be insufficient to capture the dominant interactions.
2. Sparsity of $e^{A_c T}$: at lower sampling rates, the sparsity of the discrete-time interaction matrix A_d decreases, inhibiting the advantage of the sparse estimation approach. Figure 4.13 illustrates this issue. The sparsity of the discretized matrix A_d decreases quickly with increasing sampling time T . On the other hand, the

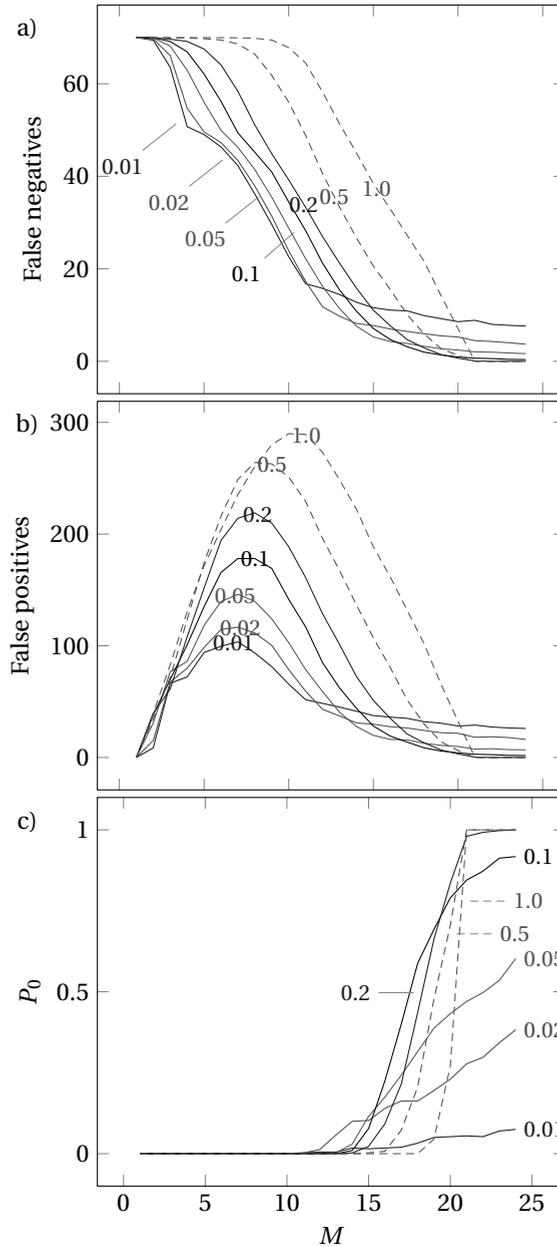


Figure 4.11: Discretized continuous-time networks estimation: average number of a) false negative and b) false positive errors, and c) the probability of a correct estimate P_0 depending on the number of measurements M and sampling time T (between 0.01 and 1). Networks of order 20 were generated with pole locations $-1.0 \pm yi$, with y in the interval $(0, \pi/2)$. Curves were averaged over 20 networks, with 20 experiments for each number of measurements.

ℓ_1 -norm of A_d is lower at higher sampling rates, which means that dominant interactions as defined in the continuous-time network, will not manifest itself as strongly in the discrete-time matrix A_d .

3. Bilinear transformation issues: as mentioned before, the transformation from the estimated (sparse) matrices A_d and B_d to the corresponding continuous-time matrices is very likely to be problematic, again especially when the amount of available data is limited.

These observations can help to interpret the results in Figure 4.11. The effect of the sampling time T on the type of error made in the estimation procedure are visualized in Figure 4.11a) and b). At lower sampling rates $T = 1.0$ and $T = 0.5$ the number of errors is higher, most likely because the discrete-time model is not sparse enough. The number of errors decreases with higher sampling rates, but starting from $T = 0.05$ a different phenomenon appears: the performance at low amounts of measurements M still (slightly) increases, but as M increases, the number of errors becomes higher than at lower sampling rates. Even at a sufficient amount of available measurements, the estimation is incorrect. This can be explained by the fact that at these higher sampling rates, the time covered by the measurements is simply too short to fully capture the dominant interactions of the network. Figure 4.12 further illustrates this point and also shows that the pole locations of the data generating network are in part determining the appropriate choice for the sampling time T and the number of measurements required to correctly estimate the data generating network parameters. Starting from the same network analyzed in Figure 4.11 with continuous-time pole locations at $-1.0 \pm yi$ (shown in Figure 4.12a)), the real parts of the pole locations are shifted to -0.5 and -0.1 (see Figure 4.12b) and c)). The effect of this shift is that at high sampling rates, for instance at $T = 0.01$ and $T = 0.02$, the probability of a correct estimate (P_0) decreases, as the less negative pole locations lead to a slower response to the input to the network, whereas at lower sampling rates, for instance at $T = 0.5$ and $T = 1.0$ the overall probability of a correct estimate increases, which can be explained by the fact that at these sampling times the slower dynamics of the generated response are better detected.

4.8 Continuous-time network models

In the previous section the network model was either defined in discrete-time or converted to discrete-time, to be able to find a sparse representation using sparse linear regression. Both approaches are a somewhat flawed in the sense that they ignore the fact that 1) the data generating network probably is continuous-time in nature and 2) the continuous time interaction matrix is sparse, while the corresponding discrete-time interaction matrix does not necessarily possess this property. A conceptually sounder approach is to hold on to the continuous-time definition of the model and the corresponding parameterization and to try and match the continuous-time model parameters to the discrete-time measurements and inputs while striving for maximum sparsity of those parameters. In this way the model dynamics are really taken into account. A downside to this approach is that sparse linear regression is no longer ap-

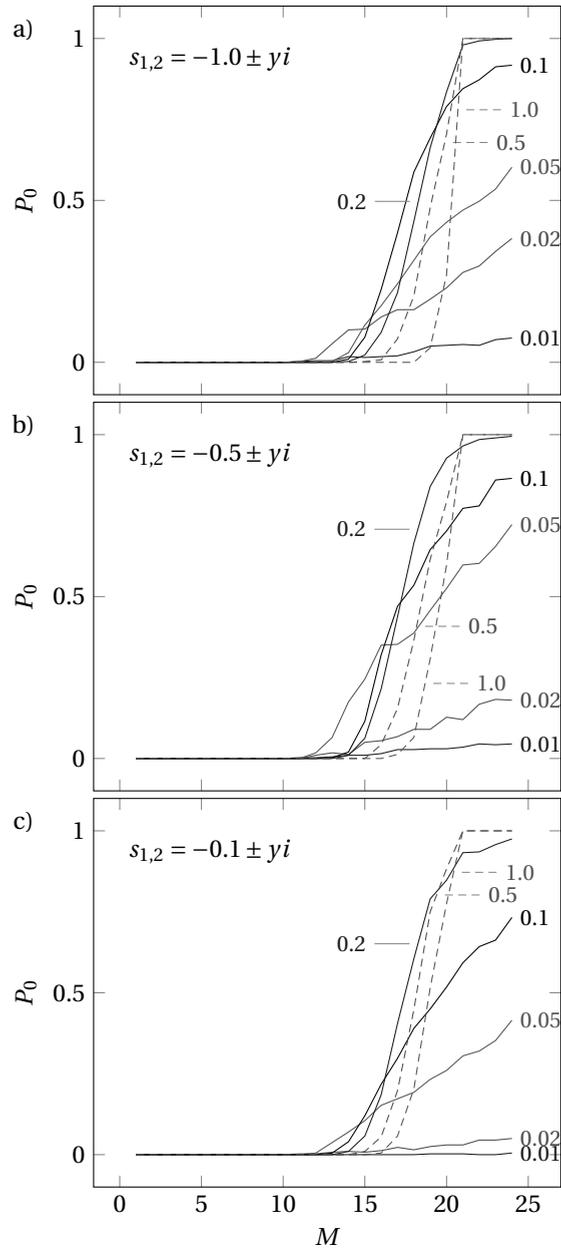


Figure 4.12: Discretized continuous-time networks estimation: effect of pole locations $x \pm yi$ of the data generating network on the probability of a correct estimate P_0 for a) $x = -1.0$, b) $x = -0.5$, and c) $x = -0.1$. Networks of order 20 were generated and estimated using a varying number of measurements M and sampling time T .

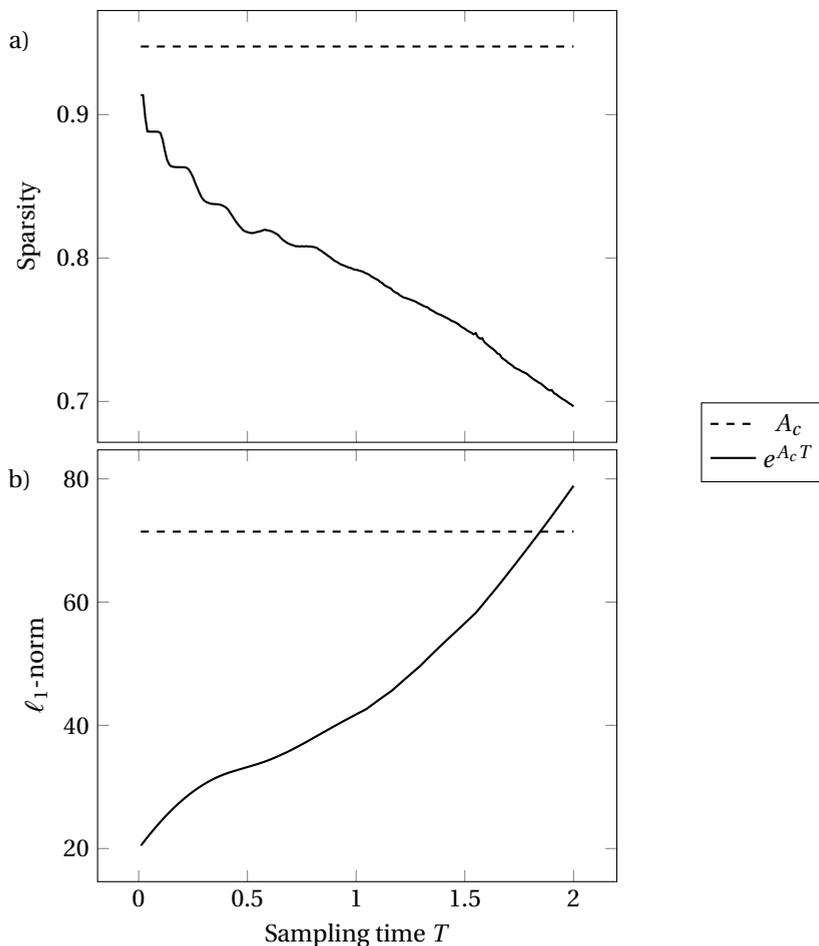


Figure 4.13: Effect of the sampling time T on a) the sparsity and b) the ℓ_1 -norm of the discretized network interaction matrix $A_d = e^{A_c T}$ with the tridiagonal interaction matrix structure. Sparsity is defined as the relative number of zero-valued parameters and is averaged over 100 networks of order $N = 20$.

plicable, since the continuous-time model parameters feature in a non-linear way in the prediction error. This means that the mixed ℓ_2/ℓ_1 optimization procedure has to be applied here.

The discrete-time network formulation in (4.34) has a simplified corresponding Jacobian and Hessian of the prediction error criterion compared to the full parameterization in Section 4.5.1. The Jacobian reduces to

$$\begin{aligned} J(\theta)_{k,i} &= \frac{\partial e[k]}{\partial \theta_i} = -\frac{\partial x[k]}{\partial \theta_i} \\ \frac{\partial x[k+1]}{\partial \theta_i} &= \frac{\partial A_d}{\partial \theta_i} x[k] + A_d \frac{\partial x[k]}{\partial \theta_i} + \frac{\partial B_d}{\partial \theta_i} u[k], \end{aligned} \quad (4.48)$$

and the Hessian of $V(\theta)$ reduces to

$$\begin{aligned} V(\theta) &= e(\theta)^T e(\theta), \\ \frac{\partial V}{\partial \theta_i} &= 2e(\theta)^T J(\theta)_i, \\ \frac{\partial^2 V}{\partial \theta_i \partial \theta_j} &= 2J(\theta)_j^T J(\theta)_i + 2e(\theta)^T \frac{\partial^2 e(\theta)}{\partial \theta_i \partial \theta_j}, \\ \frac{\partial^2 e[k]}{\partial \theta_i \partial \theta_j} &= -\frac{\partial^2 x[k]}{\partial \theta_i \partial \theta_j}, \\ \frac{\partial^2 x[k+1]}{\partial \theta_i \partial \theta_j} &= \frac{\partial^2 A_d}{\partial \theta_i \partial \theta_j} x[k] + \frac{\partial A_d}{\partial \theta_i} \frac{\partial x[k]}{\partial \theta_j} + \frac{\partial A_d}{\partial \theta_j} \frac{\partial x[k]}{\partial \theta_i} + A_d \frac{\partial^2 x[k]}{\partial \theta_i \partial \theta_j} + \frac{\partial^2 B_d}{\partial \theta_i \partial \theta_j} u[k]. \end{aligned} \quad (4.50)$$

The Jacobian and Hessian resulting from a continuous-time network (with sampling time T) can now be derived using the equation 4.24 in Section 4.5.1.

4.8.1 Experiments

The data generating continuous time networks are identical to the networks investigated in Section 4.7. Contrary to the discretized continuous-time network approach, trying to directly estimate a sparse continuous-time network requires an initial parameter estimate as a starting point for the iterative sparse optimization procedure. Here, a trial-and-error strategy is applied, where random non-sparse initial network models are generated, followed by prediction error minimization on the generated data. An estimated network is then accepted as an initial solution if it is stable. All stable initial solutions enter the mixed ℓ_2/ℓ_1 optimization procedure. During the optimization procedure, the stability of the network given a certain parameter estimate is monitored to ensure that the estimate at an iteration step does not wander off to an unstable region within the i/o equivalence space. This is not an unlikely scenario, as the limited amount of available i/o data may very well fit to the dynamics of an unstable network. Choosing an appropriate bound β on the increase of the prediction error in the ℓ_1 -step of the algorithm is crucial to ensure a smooth trajectory over the manifold of i/o equivalent network models. Only network estimates that are stable after the optimization procedure are taken into account, and the most sparse solution from this set of stable estimates is ultimately selected as the final solution.

4.8.2 Results

Examples of continuous-time sparse network estimation

An prerequisite for the mixed ℓ_2/ℓ_1 optimization procedure to be applicable, is that the data generating network represents a local minimum in terms of the ℓ_1 -norm of the parameter vector, given the limited amount of available measurement data. The optimality in terms of the prediction error is guaranteed in the noiseless case, as investigated here. Figure 4.14 illustrates what happens when the data generating network is *not* an optimum of the mixed ℓ_2/ℓ_1 criterion space. Starting from the data generating system as the initial solution, the algorithm moves away from the initial solution to an estimate with a lower ℓ_1 -norm, while retaining the optimality of the ℓ_2 -norm of the prediction error. Eventually the algorithm converges to an unstable network estimate. The choice of the bound β on the increase of the prediction error criterion in the ℓ_1 -step has limited effect, as shown in Figure 4.14. In this case, the i/o equivalence space allows for a sparser network configuration than the data generating network.

If however the data generating network is in fact a local optimum in terms of prediction error and ℓ_1 -norm of the network parameters, the mixed ℓ_2/ℓ_1 optimization algorithm can converge to this network, starting from a different network configuration. Figure 4.15 shows an example in which the original data generating network is retrieved in a underdetermined setting, even after the first prediction error minimization step leads to an unstable network estimate. The number of measurements needed to create such a setting is however already quite close to the number that would be sufficient to find a unique solution based on the prediction error criterion alone (in principle $M = 11$, disregarding potential effects by the choice of sampling time T).

Stability and sparsity of continuous-time network models

The phenomena observed when estimating continuous-time network models, call for a more systematic analysis of the properties networks that correspond to (local) optima given a specific set of input-output measurement data. In Section 4.4 it was postulated that searching for a sparse model could potentially aid in selecting a relevant model from the set of all models that are consistent with the measurement data. Intuitively, one could even imagine that minimizing the ℓ_1 -norm of the network model parameters, while safeguarding the local optimality of the least-squares fit to the measurement data, will favour stable network estimates. However, this seems to be a erroneous assumption, as can be observed from the results in the previous sections. The relationship between the measurement data and the stability of the sparsest fitting model appears to be more complex, especially in the case of continuous-time network models. In Figure 4.16 this relationship is visualized for a low order continuous-time network ($N = 2$), where an estimate is computed based on 2 measurements generated by this network. There are 6 parameter to be estimated (4 entries of the interaction matrix A_c and 2 of the input matrix B), and 4 data points (2 for each time instant), leading to 2 degrees of freedom. The effect of this freedom in choosing parameter values on the stability of the estimated network and the ℓ_1 -norm of the corresponding parameter vector was investigated by fixing the value of the estimated entries \hat{a}_{11} and

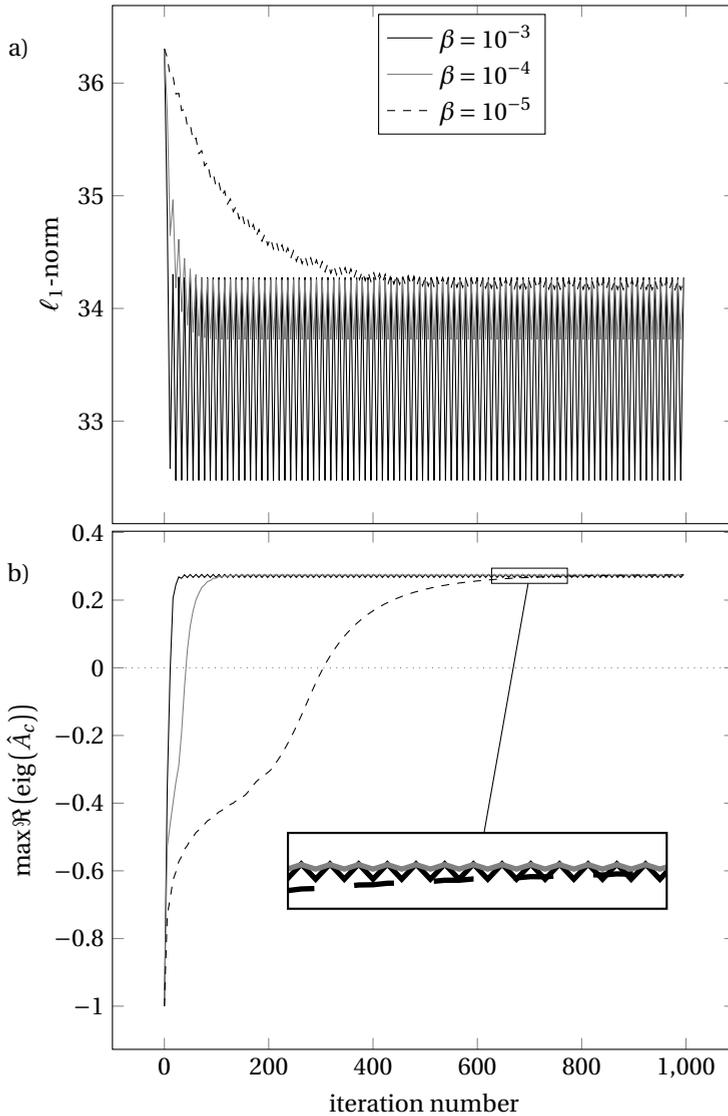


Figure 4.14: Example of the trajectories during the mixed ℓ_2/ℓ_1 optimization algorithm of a) the ℓ_1 -norm of the network parameters and b) the largest real part of the eigenvalues of the estimated network matrix \hat{A}_c . Trajectories are shown for different choices for the bound β on the increase of the prediction error criterion in the ℓ_1 -step. The algorithm was applied on measurement data ($M = 8$, sampling time $T = 1$) generated by a network of order $N = 10$, with the real part of the poles of the system located at -1 . The initial solution was the data generating network. To improve interpretability of the trajectories, not all iterations are shown.

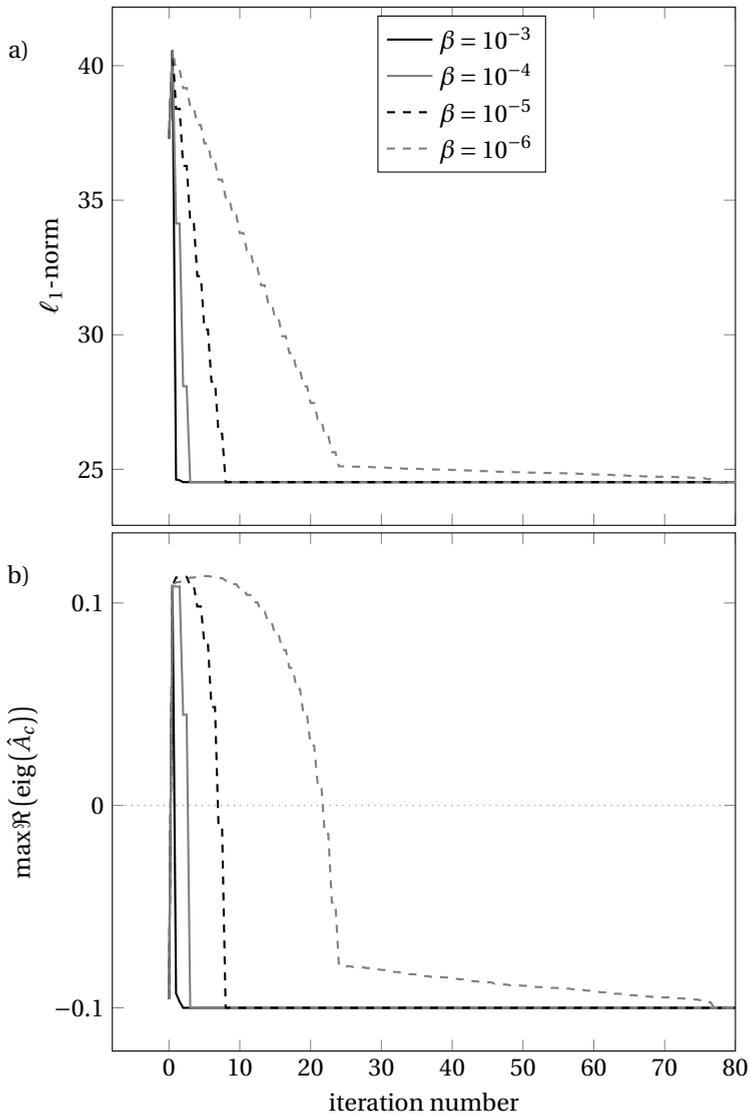


Figure 4.15: Example of the trajectories during the mixed ℓ_2/ℓ_1 optimization algorithm of a) the ℓ_1 -norm of the network parameters and b) the largest real part of the eigenvalues of the estimated network matrix \hat{A}_c . Trajectories are shown for different choices for the bound β on the increase of the prediction error criterion in the ℓ_1 -step. The algorithm was applied on measurement data ($M = 9$, sampling time $T = 0.5$) generated by a network of order $N = 10$, with the real part of the poles of the system located at -0.1 . The initial solution was a randomly generated network.

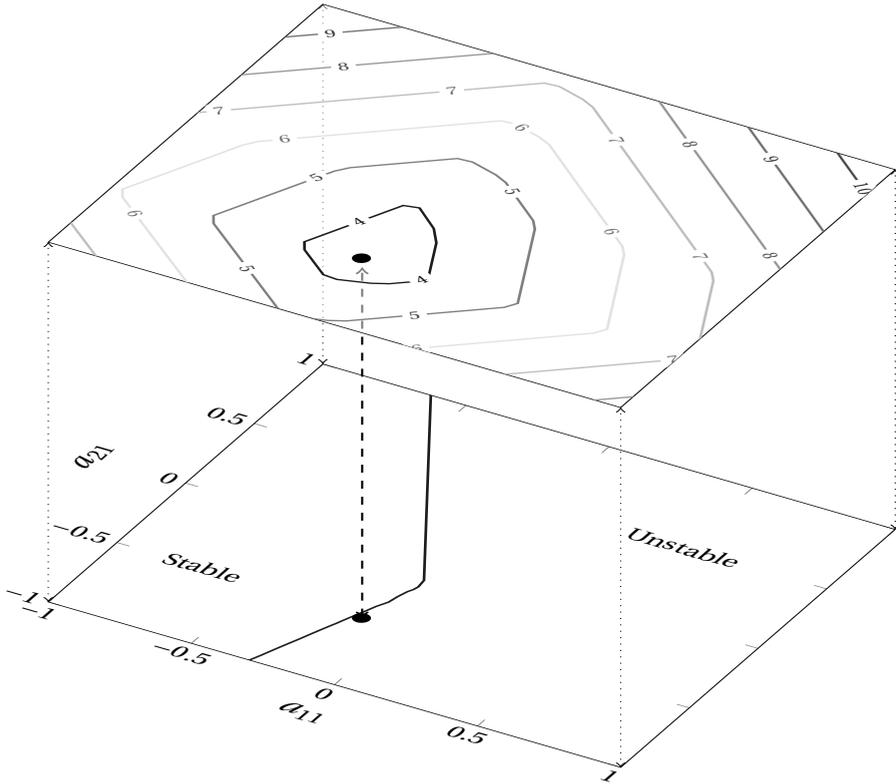


Figure 4.16: Stability (lower plane) and ℓ_1 -norm (upper plane) of networks estimates of order $N = 2$ as a function of the entries a_{11} and a_{21} of the network interaction matrix A_c . Network parameters are estimated on 2 measurements. The black dot indicates the pair of (a_{11}, a_{21}) producing the minimal ℓ_1 -norm (upper plane), which in this setting corresponds to an unstable network estimate (lower plane).

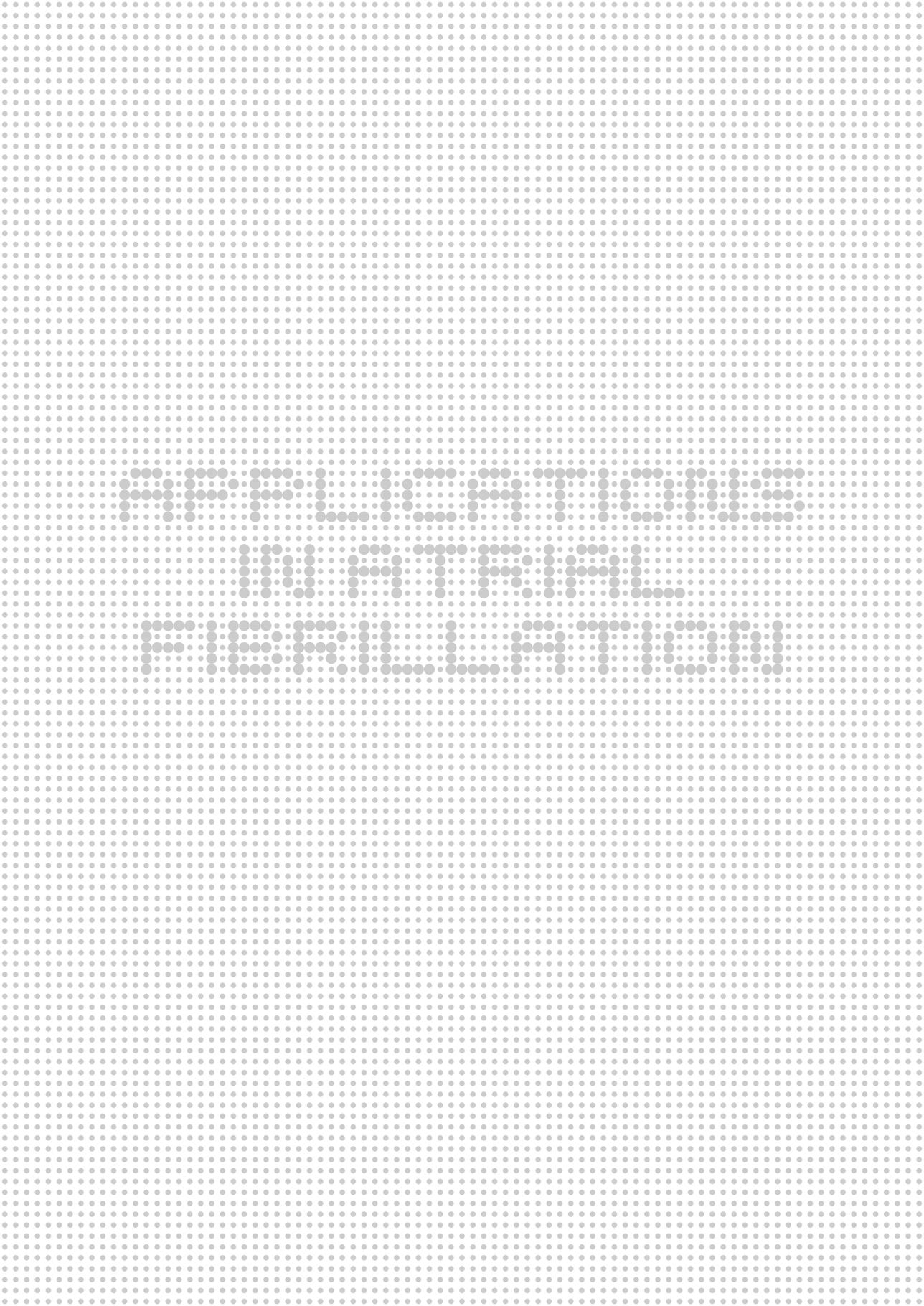
\hat{a}_{12} of the matrix \hat{A}_c and computing the (now unique) prediction error minimization solution for the remaining parameters. As is clear from Figure 4.16, a local minimum for the ℓ_1 -norm of the parameter vector does not necessarily correspond to a stable solution.

Based on this observation it can be concluded that the mixed ℓ_2/ℓ_1 optimization algorithm has major limitations in the setting of continuous-time networks with a limited amount of measurement data, which render it unsuitable for practical applications. The fact that both steps in the algorithm do not take the stability of the current estimate into account makes it likely that the algorithm will converge to an unstable solution, as the i/o-equivalence space of network models can contain unstable sparse optima in the case of a limited amount of measurement data.

4.9 Conclusions

In this Chapter the sparse estimation framework has been applied to the class of LTI state-space models, and specifically to a subclass of network models. Taking into account the nonlinear least squares prediction error criterion, the mixed ℓ_2/ℓ_1 -optimization procedure is applicable here. In a fully parameterized setting, with sufficient available measurement data, model equivalence is still present in the form of natural equivalence via a nonsingular state-space transformation matrix T . In this case the ℓ_1 -minimization step in the algorithm can be taken in such a way that the updated model estimate stays within the optimal i/o equivalence class, by associating the ℓ_1 -minimization search direction with a state-space transformation matrix T . In structured parameterizations, a choice for the bound β has to be made on the maximum deterioration of the prediction error criterion value V of the current estimate, based on a trade-off between convergence speed of the algorithm and the risk of stepping outside of the optimal i/o equivalence manifold.

Sparse maximization of discrete-time network model interactions can be stated as a sparse linear regression problem, and solved accordingly, under the condition that every row of the interaction matrix A_d is sparse. Experiments using a ring-shaped network topology, resulting in a such a sparse interaction matrix, show that also in this case sparse linear regression is able to correctly reconstruct the network structure in an underdetermined setting. Moving to continuous-time interaction networks, it is still possible to apply sparse linear regression by assuming sparsity in the discrete-time interaction matrices and transforming the estimated solution back to a continuous-time state-space representation. Here, the performance of the algorithm is very much dependent on the sampling time T and the sparsity of the discretized network interaction matrix. Estimating the sparse continuous-time network interaction matrix directly is feasible, but here a limited amount of available measurements increases the risk of finding a solution with minimal ℓ_1 -norm that corresponds to an unstable network model.



“This whole world’s wild at heart and weird on top.”

— *Wild at Heart*, dir. David Lynch

Chapter 5

Introduction to atrial fibrillation

5.1 Definition and treatment of atrial fibrillation

Atrial fibrillation (AF) is the most common cardiac arrhythmia, that affects 1.5-2% of the general population [15]. The normal sequence of activation of the heart - sinus rhythm (SR) - is disrupted and instead of organized waves originating from the sinoatrial node in the right atrium and traveling through the atrial conduction system to the atrioventricular (AV) node, the atria are activated in a disorganized way by multiple wavelets. This in turn can cause fast and irregular activation of the ventricles. Patients suffering from AF have a increased risk of stroke, because of potential thrombus formation in the atria as a result of the disturbed blood flow. If untreated, AF is a progressive disease, starting out as short intermittent episodes of silent AF. When eventually diagnosed, it is classified as *paroxysmal AF* when AF episodes are still short and self-terminating, lasting no more than 48 hours. It is classified as *persistent AF* when either an episode lasts longer than 7 days, or when it is no longer self-terminating and needs to be stopped using cardioversion, either pharmacological cardioversion or electrical cardioversion. More severe stages of AF are *long-standing AF*, where a patient is in AF for more than a year, and *permanent AF*, when patient and physician agree to no longer try to restore sinus rhythm.

Treatment and management of AF depend on the progression of the disease and can be divided into rhythm control and rate control strategies (see Figure 5.1). Rhythm control means that one tries to restore sinus rhythm and rate control means that one tries to bring down the heart rate, but not necessarily convert to SR. Cardioversion (CV) and ablation are the main forms of rhythm control. Pharmacological cardioversion is achieved by administering an anti-arrhythmic drug, such as flecainide, amiodarone or vernakalant, and is used in patients with paroxysmal AF, and most effective in patients with recent onset AF (< 48 hours of AF). Electrical or direct current cardioversion works by delivering a synchronized shock to the patient at the right moment in the cardiac cycle. Electrical CV is applied when pharmacological CV is ineffective and in patients with persistent AF. Catheter or surgical ablation is often chosen in patients with symptomatic AF that do not respond to anti-arrhythmic medication. Catheter ablation of

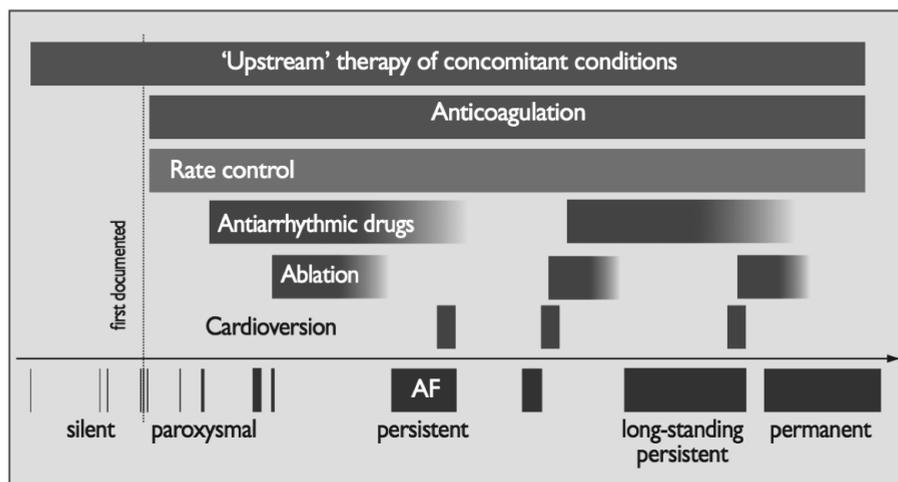


Figure 5.1: Progression and treatment of AF. Taken from the ESC guidelines for the management of atrial fibrillation [15].

AF works by inserting a catheter through a vein in the groin or the neck and guiding this catheter to the atria to create lesions in atrial tissue by means of energy, for instance radiofrequency or cryothermic energy. Surgical ablation is more invasive and relies on open heart surgery to deliver the energy needed to create lesions. At all stages anticoagulation is usually advised to reduce the risk of thrombus formation and stroke.

Although management of AF has been subject to extensive research for many decades, the reasons a patient develops AF and the mechanism(s) behind the perpetuation of AF are still not fully understood. Figure 5.2 from the 2012 expert consensus statement on catheter and surgical ablation of AF [17] illustrates the idea that the initiation and perpetuation of AF is caused by a combination of triggers from the autonomic nervous system and wave reentry. The complexity of the AF substrate in an individual patient is however difficult to assess. The recurrence rate of AF after for instance direct current CV is high (typically around 50% recurrence within 4-6 weeks), indicating that a large amount of patients have not benefitted from this treatment. Better quantification of the complexity of the pattern of AF may help in personalizing patient treatment.

5.2 Quantitive description of atrial fibrillation

The characteristics of AF are studied in many ways, both invasively and noninvasively. The most direct way of measuring the electrical signals in an atrium is direct contact mapping, where electrodes are placed on the epicardial or endocardial surface of an atrium. The signals measured here are called electrograms. Electrograms can be unipolar signals, that contain the voltage difference to some reference signal outside of the atrium, or bipolar signals, the voltage difference between two electrodes placed within a short distance (typically a couple of millimeters) from each other on the atrium. These invasive measurement can be done on the inside of an atrium (en-

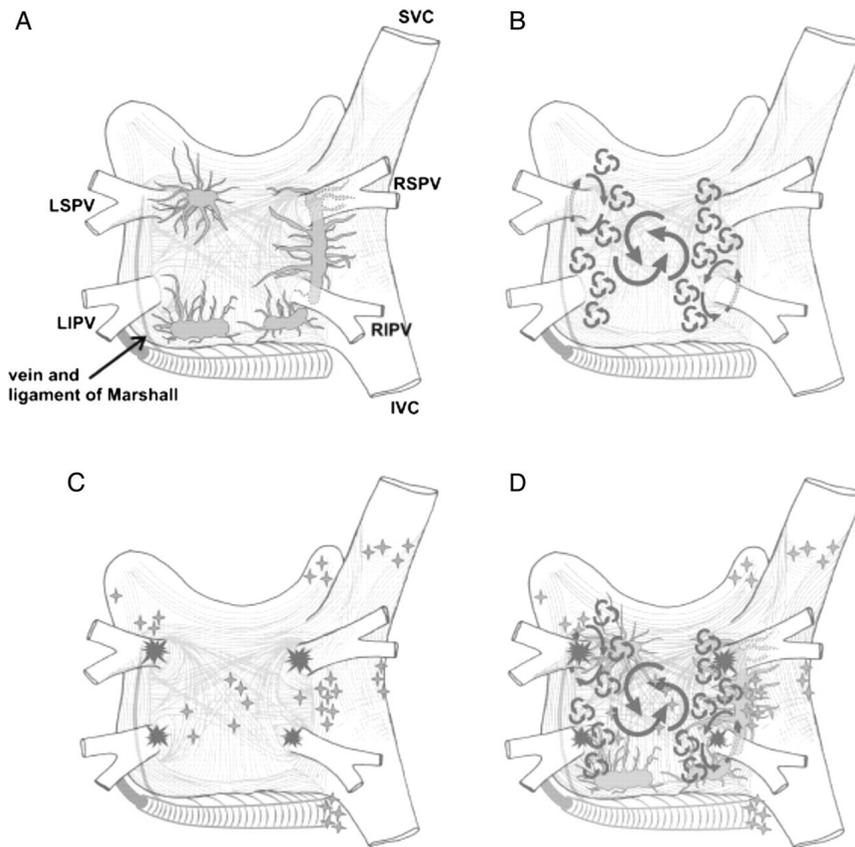


Figure 5.2: Atrial structure and autonomic nervous system (A), and potential mechanisms of AF, with B) large and small reentrant waves, C) common locations of triggers (focal sources) for AF located around the pulmonary veins (PV) and other locations, and D) superposition of anatomical and pro-arrhythmic mechanisms. Figure adapted from [17].

docardial), using for instance commercially available electroanatomic mapping systems such as Ensite™NaVx™(St. Jude Medical) or CARTO®(BioSense Webster), or on the outside of an atrium during cardiac surgery. Signals measured on the body surface provide a noninvasive way to assess cardiac function in general and atrial function in specific. The standard 12-lead electrocardiogram (ECG) is designed to evaluate cardiac function as seen from different angles on the body surface. More recently, high-coverage body surface potential maps (BSPM) have been introduced that try to capture more information than the 12-lead ECG by covering the body surface with many electrodes. In combination with anatomical information obtained by for instance a CT-scan, one can try to reconstruct the electrical signals on the heart. This technique is then called ECG imaging (ECGI).

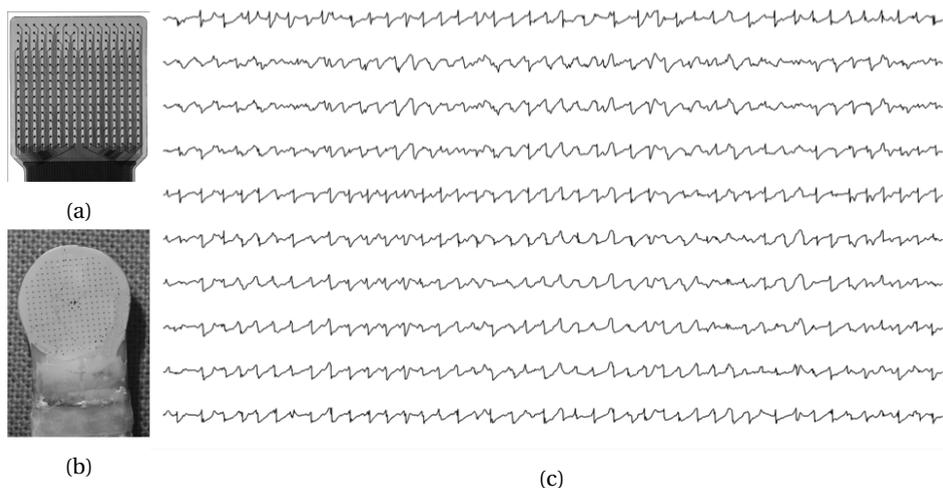


Figure 5.3: Grids of electrodes with a) a square configuration (16×16 grid of electrodes, electrode distance 1.5mm), b) a circular configuration (234 electrodes, electrode distance 2.4mm), and c) examples of unipolar direct contact mapping electrograms.

5.2.1 Invasive analysis

The invasive measurements analyzed in this thesis were acquired using regular grids of electrodes as depicted in Figure 5.3 that were placed on the epicardial atrial wall during open-chest surgery. In this way not only unipolar electrograms were recorded, but also spatial patterns of conduction could be reconstructed based on the configuration of the grid of electrodes. These invasive measurements provide a detailed description of the complexity of the propagation patterns in AF in different models (see [52],[99],[41],[4] for examples of invasive AF complexity analysis based on high-density contact mapping of AF in goat and humans). In Chapter 6 a novel probabilistic algorithm is presented to automatically detect the moment of local atrial depolarization (so-called deflections) in direct contact electrograms and to reconstruct two-dimensional atrial propagation patterns in terms of fibrillation waves. In Chapter 7 an alternative approach to propagation pattern identification is introduced that applies the concept of sparse linear regression to identify time-delayed interactions between electrograms at different electrodes, circumventing the need for manual or automatic annotation of atrial deflections. Here the focus is more on the detection of recurrent propagation patterns than on a detailed wave-based description of the AF substrate.

5.2.2 Noninvasive analysis

Noninvasive analysis of AF relies on the assumption that the ECG contains information on the complexity of the AF substrate. If this is true (at least to some extent), noninvasive assessment of AF complexity could aid in the management of AF, by predicting treatment outcome. In Chapter 8 a systematic comparison of existing noninvasive AF complexity parameters is performed when it comes to the prediction of successful

pharmacological cardioversion. From the large number of candidate predictors, dominant predictors are selected using a sparse logistic regression approach.

Chapter 6

Probabilistic electrogram analysis and atrial fibrillation wave reconstruction

6.1 Introduction

Atrial fibrillation (AF) is an arrhythmia where the electrical activity in the atria is irregular instead of well organized. Multiple wavelets wander throughout the atria, instead of a single coordinated wave [87]. High-density atrial contact mapping of AF provides the most direct information on the spatiotemporal complexity of AF. It allows one to describe the process of AF in its most elementary form, the separate fibrillation waves [4]. From these wave propagation patterns it is possible to quantify the complexity of AF, for example in terms of the number of waves, the wave size, the wave conduction velocity or the wave source (peripheral or transmural breakthrough). Complexity of the AF activation pattern is a strong determinant of responsiveness to AF therapy. Assessment of the AF activation pattern might therefore be used for decision-making in the management of AF patients [51]. Current analysis methods still involve labor-intensive manual annotation of atrial deflections and waves, which limits the amount of fibrillation data that can be analyzed within a reasonable timeframe. Manual editing also increases the risk of subjective editing, which can lead to lower inter-observer consistency. To overcome these limitations we developed a novel method that identifies atrial deflections and fibrillation waves in a rapid and fully automated way, based on estimated probabilistic properties of the recorded fibrillation process. The details of this automatic procedure are presented in this Chapter, as well as the results of a validation study. In the design of the new deflection detection and wave mapping method, we aimed to incorporate electrophysiological knowledge to be able to compute wave map solutions that both visually and intellectually reflect the way electrophysiologists would construct them.

6.2 Methods

6.2.1 Data acquisition

Unipolar atrial fibrillation electrograms were recorded in 15 patients during cardiac surgery using a 16×16 square grid of electrodes with an inter-electrode distance of 1.5mm. Acute AF was induced in 8 patients who were in sinus rhythm, 7 patients were already in AF during surgery (either paroxysmal or persistent AF). Signals were acquired from the epicardium of the right atrial free wall (RA) ($n = 15$) and the posterior left atrium (LA) ($n = 11$) with a sampling frequency of 1kHz. Segments of 4 seconds of AF were manually annotated by three experienced electrophysiologists to determine local atrial deflections and to identify clusters of deflections that form separate fibrillation waves, following the algorithm described in [4].

6.2.2 Electrogram pre-processing

The first step in processing the atrial measurements is to eliminate electrograms that exhibit a bad signal-to-noise ratio. To enable a valid comparison between the new method and manually annotated signals, the same electrograms were eliminated in both methods. Signals were then filtered with a third order zero-phase Chebyshev 0.5Hz high-pass filter to remove any baseline drift. Ventricular far-field disturbances in the atrial signal were removed by ventricular R-wave detection in a synchronously recorded ventricular signal [75], followed by single-beat QRST-template cancellation based on the adaptive singular value decomposition cancellation method by Alcaraz et al. [1]. This method determines the morphology of the QRST complex as the most significant principal component of a set of QRST windows in a single lead. We adapted this method to compute the QRST complex for a single beat in all electrograms to account for beat-to-beat QRST complex variability. In atrial electrograms these ventricular far-fields are usually relatively low in amplitude compared to the local atrial deflections, but nonetheless they can cause false positives when detecting deflections and constructing waves.

6.2.3 Intrinsic deflection detection

By interviewing several electrophysiologists, a single predefined deflection shape was constructed to detect all candidate atrial deflections in the electrograms. This template could vary in duration (5 – 50ms). At each time-point in an electrogram the maximum correlation was determined between the electrogram and all template durations. The local maxima in the resulting template correlogram were marked as candidate deflection positions, as shown in Figure 6.1a, 6.1b and 6.1c. These candidate deflections can be either true intrinsic (or local) deflections or far-field deflections - fluctuations of the electrograms caused by activations remote from the electrode - or deflections caused by external disturbances. To distinguish between the three types of deflections, we exploited the underlying distribution of the AF cycle length (AFCL). This distribution was estimated by iteratively increasing the minimally allowed deflection amplitude until

the distribution of the remaining deflection intervals showed a clear peak. The position and shape of this peak in the interval distribution of deflections with higher amplitude reflect the interval distribution of the true intrinsic deflections. The distribution estimation procedure was automated by introducing a maximum deflection interval threshold (default value 250ms). The minimally allowed deflection amplitude was increased until a maximum percentage (default value 10%) of the remaining deflection intervals were larger than the maximum deflection interval threshold. A normal distribution with parameters $\theta_{\text{CL}} = (\mu_{\text{CL}}, \sigma_{\text{CL}})$ was fitted on the resulting deflection intervals. The procedure is illustrated in Figure 6.2. Given the sequence of candidate deflections $\{c_n\}_{n \in \{1, 2, \dots, N\}}$ and the AFCL distribution estimate, a deflection type assignment problem is formulated, where the goal is to select a subsequence of intrinsic deflections $\{c_{n_r}\}_{r \in \{1, 2, \dots, k\}, n_k \leq N}$ with maximum interval probability. Assuming a sequence of deflection intervals is i.i.d., the joint interval probability of a subsequence $\{c_{n_r}\}$ can be expressed as

$$P(\{c_{n_r}\}|\theta_{\text{CL}}) = \prod_{i=1}^{k-1} f(t_{r_{i+1}} - t_{r_i}|\theta_{\text{CL}}), \quad (6.1)$$

where t_{r_i} is the central time of deflection c_i . The interval between the time t_{r_1} of the first deflection in a subsequence and the beginning of the recording t_0 and the interval between the time t_{r_k} of last deflection in a subsequence and the end of the recording t_{end} has to be included in the joint probability of the subsequence to include the constraint that intrinsic deflections are to be found in all parts of the recording, forming a chain of deflections that are linked by probable deflection intervals.

$$P((t_0, \{c_{n_r}\}, t_{\text{end}})|\theta_{\text{CL}}) = \tilde{f}(t_{r_1} - t_0|\theta_{\text{CL}}) \cdot \left(\prod_{i=1}^{k-1} f(t_{r_{i+1}} - t_{r_i}|\theta_{\text{CL}}) \right) \cdot \tilde{f}(t_{\text{end}} - t_{r_k}|\theta_{\text{CL}}), \quad (6.2)$$

where

$$\tilde{f}(t_j - t_i|\theta_{\text{CL}}) = \begin{cases} f(\mu_{\text{CL}}|\theta_{\text{CL}}) & \text{if } t_j - t_i \leq \mu_{\text{CL}} \\ f(t_j - t_i|\theta_{\text{CL}}) & \text{if } t_j - t_i > \mu_{\text{CL}} \end{cases} \quad (6.3)$$

A heuristic greedy algorithm finds a solution to the sequence selection problem by starting with the complete sequence of candidate deflections and trying to improve the joint deflection probability in Equation 6.2.3 by removing deflections, starting with deflection with low amplitude and slope, until the solution converges. This order of deflection deletion in this algorithm is based on the tendency of electrophysiologists to mark steep deflections with high amplitude as intrinsic deflections and flat deflections with low amplitude as far-field deflections. An example result of an intrinsic deflection assignment solution can be seen in Figure 6.1d, 6.1e and 6.1f. The strength of this intrinsic deflection detection method is that it is able to adapt to substrate-specific deflection properties, such as amplitude, slope and interval distribution.

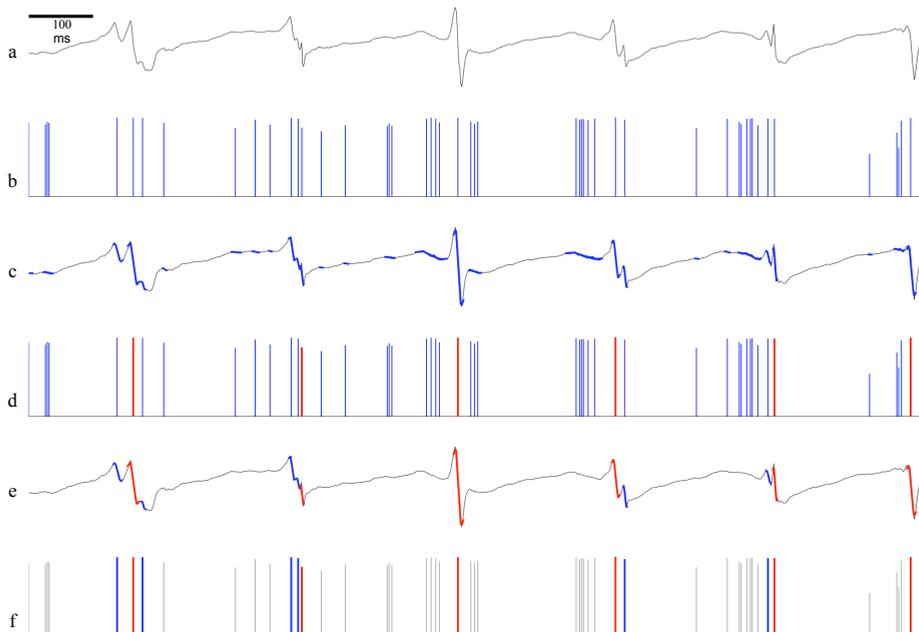


Figure 6.1: Deflection detection procedure. The pre-processed signal (a) is analyzed by a sensitive template-matching algorithm. The peaks in the correlogram are depicted in (b) and the corresponding template matches in (c). The intrinsic deflection assignment algorithm identifies the location of the intrinsic deflections (red) (d). Far-field deflections (blue) are determined as deflections that have a minimal amplitude of 10% of the median intrinsic deflection amplitude and a minimal slope of 10% of the median intrinsic deflection slope. The final result is depicted in (e) and the corresponding correlogram peak locations in (f).

6.2.4 Fibrillation wave construction

The intrinsic deflection detection step determines the sequences of intrinsic deflections that are used to construct the fibrillation waves. The center of an intrinsic deflection is taken as the moment of local activation. Wave construction is divided into three phases. First, partial waves are created based on a minimum conduction velocity criterion (default value 20 cm/s) between two neighboring activations in the electrode grid. In experimental studies this threshold was identified as reasonable cut-off value for the occurrence of conduction block [4]. Activations that can be linked to two or more partial waves are not yet assigned. In the second phase statistical conduction properties of these partial waves are determined. The distribution of the wave conduction velocity is estimated by computing the conduction velocity in partial waves containing at least 9 activations. The local conduction velocity for each wave activation is determined by fitting a tangent plane onto the surface formed by the activation and the activations within the same wave at the directly surrounding electrodes. The conduction

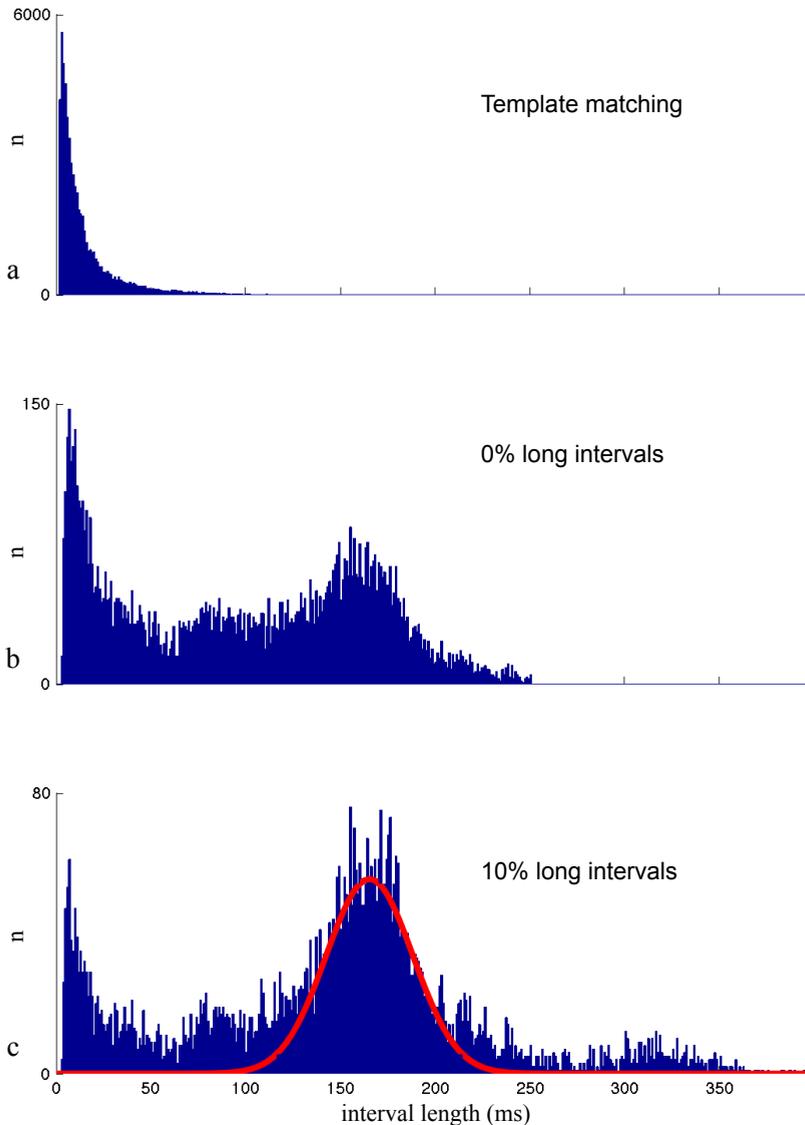


Figure 6.2: Estimating the intrinsic deflection interval distribution. A histogram of the intervals between the deflections found in the template matching procedure is shown in (a). Increasing the minimal amplitude of a deflection without causing long intervals, produces the histogram in (b). If a maximum of 10% long intervals are allowed, a clear peak appears, as can be seen in (c). A normal distribution is fitted onto this peak.

velocity is then computed as the reciprocal of the plane gradient vector length and the direction of conduction as the plane gradient angle. The resulting velocity distribution is approximated by a gamma distribution with parameters $\theta_{CV} = (\alpha_{CV}, \zeta_{CV})$. Besides the velocity distribution also the distribution of the conduction deviation or tortuosity within a wave is determined by computing the mean conduction direction difference at each activation compared to the activations within the same wave at the directly surrounding electrodes. This tortuosity distribution is approximated by a normal distribution with parameters $\theta_{TO} = (\mu_{TO}, \sigma_{TO})$. The third phase consists of assigning the unassigned activations to adjacent waves based on a maximum local conduction velocity and conduction direction probability. Given an activation $a_{e,t}$ at electrode e at time t and a set of candidate waves W , this probability is defined as

$$P(CV = cv_{a,w}, TO = to_{a,w}) = P(cv_{a,w} | \theta_{CV}) \cdot P(to_{a,w} | \theta_{TO}), \quad (6.4)$$

where $cv_{a,w}$ denotes the conduction velocity that results from adding activation $a_{e,t}$ to wave $w \in W$, and to a,w the mean conduction tortuosity.

6.2.5 Validation

The results of the novel automated deflection detection and wave map construction procedure were validated by comparing the location of intrinsic deflections to the location manually annotated deflections. Locations were considered equal if the manually annotated deflection was positioned within the descending part of the automatically detected intrinsic deflection. Automatically computed wave maps were compared to manually constructed maps in terms of median wave conduction velocity, median AF cycle length, number of waves per AFCL, number of breakthrough waves (BT) per AFCL and average wave size.

6.3 Results and discussion

The sensitivity of the intrinsic deflection detection algorithm compared to the manual intrinsic deflection annotation is $87 \pm 6.7\%$ (mean \pm SD). The positive predictive value of the automated intrinsic deflection algorithm is $89 \pm 3.8\%$. Figure 6.3 and Table 6.1 contain the comparison between the result of automated wave construction algorithm and the manually created waves. In general, the automated procedure produces very similar results to a manual annotation, most notably the wave conduction velocity and the AF cycle length. The number of waves per AFCL is only slightly overestimated, but the number of breakthrough waves per AFCL is roughly doubled by the automated procedure. An explanation for this phenomenon is that a manual editor tends to minimize the number of breakthroughs by searching for alternative activation pathways originating from the edge of the mapping array. A consequence of the larger number of waves detected in the automated procedure is that the average size of automatically created waves is smaller than the average size of manually created waves. Importantly, correlations are high, which effectively shows that the automated procedure is an adequately and valid substitute for the cumbersome manual annotation of atrial electrograms and manual atrial wave reconstruction.

Table 6.1: Comparison between automated and manual wave construction. Numbers are reported as mean \pm SD. All correlations are significant ($p < 0.01$).

Category	Manual	Automated	r
Number of waves per AF cycle	5.6 \pm 2.7	7.8 \pm 3.3	0.96
Number of BT per AF cycle	1.8 \pm 1.2	3.7 \pm 2.0	0.94
Wave conduction velocity (cm/s)	65 \pm 12	66 \pm 13	0.97
AFCL (ms)	200 \pm 32	204 \pm 29	0.97
Wave size (number of electrodes)	53 \pm 29	36 \pm 18	0.96

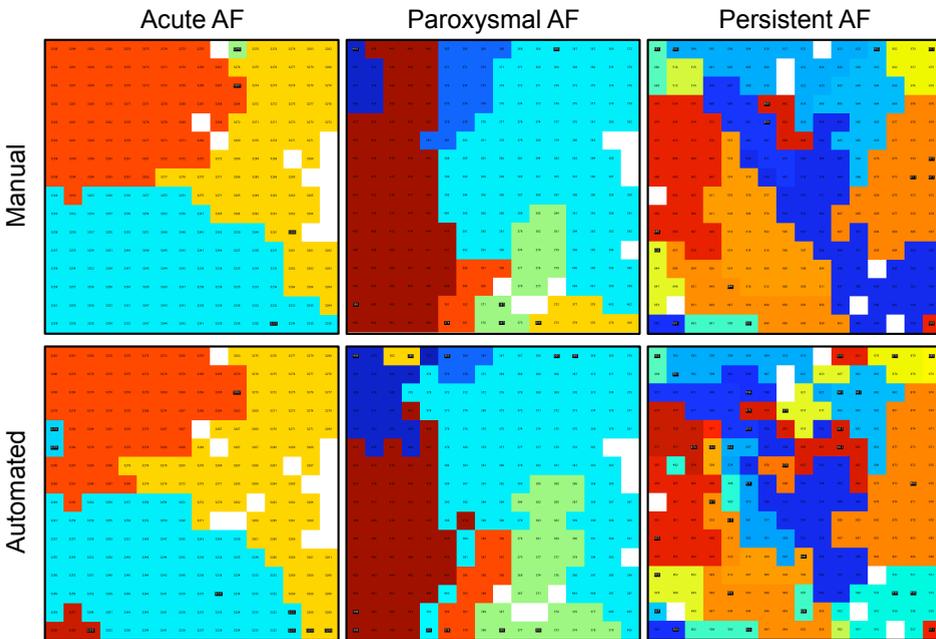


Figure 6.3: Examples of manually edited wave reconstructions versus wave reconstructions computed by the automated procedure. The maps show the wave reconstruction for a patient in acute AF paroxysmal AF and persistent AF. The same wave shapes can be visually identified in both the manually edited maps as well as in the computed maps, although the automated procedure tends to create more and smaller waves. This does not however affect the ranking of AF complexity.

6.4 Conclusion

We developed and validated a novel algorithm for fast and automated spatiotemporal analysis of the substrate of atrial fibrillation. The algorithm identifies the key properties of the substrate with high accuracy. Potential applications of this technique are:

1. *Assessment of spatial and temporal variability of the AF substrate.* Automatic analysis of high-density maps during longer recordings will provide greater insight into the temporal variation in the behavior of the AF substrate and the recording duration required to assess AF complexity in a more reliable way.
2. *Analysis of large amounts of fibrillation data in multicenter trials to establish a new classification of AF.* Using the automated method, the amount of AF fibrillation electrograms that can be feasibly processed and analyzed in a short amount of time will increase, enabling larger scale data studies required to establish a classification of AF.
3. *On-site AF substrate complexity assessment to tailor ablation therapy.* A (quasi) real-time implementation of the automated method can provide direct information on wave conduction patterns to guide the ablation process and can give immediate feedback to assess efficacy of an ablation lesion set.

Chapter 7

Identification of recurring wavefront propagation patterns in atrial fibrillation using basis pursuit

7.1 Introduction

The identification of recurring patterns in high-density recordings of atrial fibrillation (AF) provides valuable insight in the underlying mechanisms that determine the complexity of the AF substrate. Response to treatment of AF, whether this is for instance ablation therapy or pharmacological intervention, is dependent on this AF substrate complexity. High-density contact mapping of AF may give a good first visual impression of recurring wavefront patterns, but detailed substrate analysis requires extensive signal processing of atrial electrograms, atrial deflection detection and fibrillation wave reconstruction. We hypothesize that recurring patterns and signaling pathways can be identified using electrograms and electrode topology alone, employing a sparse multivariate autoregression (MVAR) modeling approach. A recent study showed that a similar approach can lead to meaningful results when identifying propagation patterns between several intracardiac recording sites using bipolar electrograms from a basket catheter [86]. In contrast to this approach, we intend to identify interactions between *unipolar* electrograms within a *short timeframe* to account for the dynamical nature of complex atrial fibrillation patterns. Furthermore, the spatial resolution of these unipolar recordings is higher (interelectrode distance 2.4mm) than in bipolar electrogram acquisition devices such as a basket catheter.

The MVAR model used in our approach allows us to choose 1) the maximum time-delay (the model order), 2) some dead time, and 3) model sparsity. The model order has to be chosen in such a way that sufficient but not too much past electrical activity is used to explain current activity. Including dead time is necessary to prevent (almost) simultaneous activations in a wavefront to be used to explain and detect spatial interaction in the direction of wave propagation. Sparsity is used to highlight dominant

interactions. We are interested in spatial interactions, which is why we use an averaging procedure over time, within a suitably short window to prevent too many different wavefronts to cancel each other out. A key issue here is that otherwise unrelated atrial complexes often have similar morphology, so that any method that maximizes sparsity is prone to identify dominant interactions between unrelated locations. To address this issue we developed a distance-weighted adaptation of the basis pursuit algorithm [19] that maximizes the sparsity of the MVAR interaction matrices, while also regularizing sparsity based on interelectrode distance. The algorithm was then applied to a set of recordings in a goat model of AF that contained multiple recurrent wavefront propagation patterns.

7.2 Methods

7.2.1 Sparse multivariate autoregression model

The MVAR model of order P for a set of N synchronous recorded electrograms $\mathbf{x}[k] = [x_1[k], x_2[k], \dots, x_N[k]]^T$, $k = 1, 2, \dots, M$ with a dead time δ , can be formulated as

$$\mathbf{x}[k] = \sum_{\tau=\delta}^{\delta+P-1} \mathbf{A}_\tau \mathbf{x}[k-\tau] + \mathbf{w}[k], \quad (7.1)$$

where each \mathbf{A}_τ is the $N \times N$ matrix with entries $a_{ij}(\tau)$ quantifying an interaction from $x_j[k-\tau]$ to $x_i[k]$. The vector $\mathbf{w}[k] = [w_1[k], w_2[k], \dots, w_N[k]]$ is a multivariate white noise process. The model in (7.1) can be written in matrix form:

$$X = \Phi\Theta + W, \quad (7.2)$$

where

$$\begin{aligned} X &= [\mathbf{x}[\delta + P], \mathbf{x}[\delta + P + 1], \dots, \mathbf{x}[M]]^T \\ \Phi &= \begin{bmatrix} \mathbf{x}[P] & \cdots & \mathbf{x}[M - \delta] \\ \vdots & \ddots & \vdots \\ \mathbf{x}[1] & \cdots & \mathbf{x}[M - (\delta + P - 1)] \end{bmatrix}^T \\ \Theta &= [\mathbf{A}_\delta, \mathbf{A}_{\delta+1}, \dots, \mathbf{A}_{\delta+P-1}]^T \\ W &= [\mathbf{w}[\delta + P], \mathbf{w}[\delta + P + 1], \dots, \mathbf{w}[M]]^T. \end{aligned}$$

For a set of electrograms with $M \geq NP + \delta$ and Φ with full column rank NP , the unique least squares (LS) solution can be derived for each column $\theta^{(i)}$ of Θ separately

$$\theta_{\text{LS}}^{(i)} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{x}^{(i)}, \quad (7.3)$$

where $\mathbf{x}^{(i)}$ denotes column i of X . In this study we focus on short time-segments of electrograms, which typically causes the number (NP) of parameters that needs to be estimated for each column $\theta^{(i)}$ to be much larger than the available number of observations M . If the matrix Φ has rank $r < NP$, which necessarily happens if $M - \delta - P + 1 < NP$, then there is no unique LS solution, but an optimal least squares

solution space of dimension $NP - r$. This equivalence space with respect to the LS criterion can be exploited to maximize the sparsity of the parameter vector $\theta^{(i)}$. We apply a basis pursuit algorithm to find a parameter vector $\theta^{(i)}$ with minimal ℓ_1 -norm, while retaining the optimal least squares fit:

$$\min_{\theta^{(i)}} \|\theta^{(i)}\|_1 \quad \text{subject to } \Phi\theta^{(i)} = \Phi\theta_{\text{LS}}^{(i)}. \quad (7.4)$$

Minimizing the ℓ_1 -norm has been shown to produce a solution with maximum sparsity under certain conditions [36], even in the presence of noise [29]. Note that the right-hand side vector in the constraint in (4) is a unique vector, formed by the orthogonal projection of $\mathbf{x}^{(i)}$ on the column space of Φ which does not depend on the choice of a solution in the equivalence space for the LS criterion.

An $N \times N$ spatiotemporal weight matrix C_τ can be defined and incorporated into the criterion function in (7.4) which allows to regularize sparsity at corresponding entries of \mathbf{A}_τ . Define C as

$$C = [C_\delta, C_{\delta+1}, \dots, C_{\delta+P-1}]^T, \quad (7.5)$$

and $C^{(i)}$ as column i of C . The regularized problem can now be written as

$$\begin{aligned} \min_{\theta^{(i)}} (C^{(i)})^T \left[|\theta_1^{(i)}|, |\theta_2^{(i)}|, \dots, |\theta_{NP}^{(i)}| \right]^T \\ \text{subject to } \Phi\theta^{(i)} = \Phi\theta_{\text{LS}}^{(i)}, \end{aligned} \quad (7.6)$$

The problem in (7.6) can be solved using linear programming by bringing it into standard form:

$$\begin{aligned} \min_{\theta^+, \theta^-} (C^{(i)})^T (\theta^+ + \theta^-) \\ \text{s.t. } \Phi(\theta^+ - \theta^-) = \Phi\theta_{\text{LS}}^{(i)} \\ \theta_i^+, \theta_i^- \geq 0 \quad i = 1, 2, \dots, NP, \end{aligned} \quad (7.7)$$

where $\theta^{(i)} = \theta^+ - \theta^-$. The resulting column vectors $\tilde{\theta}^{(i)}$ with minimal ℓ_1 -norm are joined to form the estimated matrix $\tilde{\Theta}$. From this matrix the estimated interaction matrices $\tilde{\mathbf{A}}_\tau$ can be constructed. To compute a solution $\theta_{\text{LS}}^{(i)}$ which features in (7.7), several approaches are possible. One is to employ the data directly as indicated in the definitions of Φ and X above and to use classical techniques from linear algebra such as QR-decomposition or an SVD-approach. However, the matrix $\Phi^T\Phi$ in the LS formula (7.3) is known to have a near block-Toeplitz structure which admits highly efficient recursive inversion. This is the basis for the well-known Levinson algorithm and its multivariate generalizations such as the Whittle-Wiggins-Robinson algorithm, see [92], which allow for the computation of an LS solution recursively in the order P of the MVAR model. This may speed up the estimation process, and it also allows one to use information theoretic criteria to select an appropriate value for P .

7.2.2 Dominant pathway identification

The goal of the sparse MVAR model estimation is to identify the dominant interactions between synchronous electrograms within a relatively short time interval of length M ,

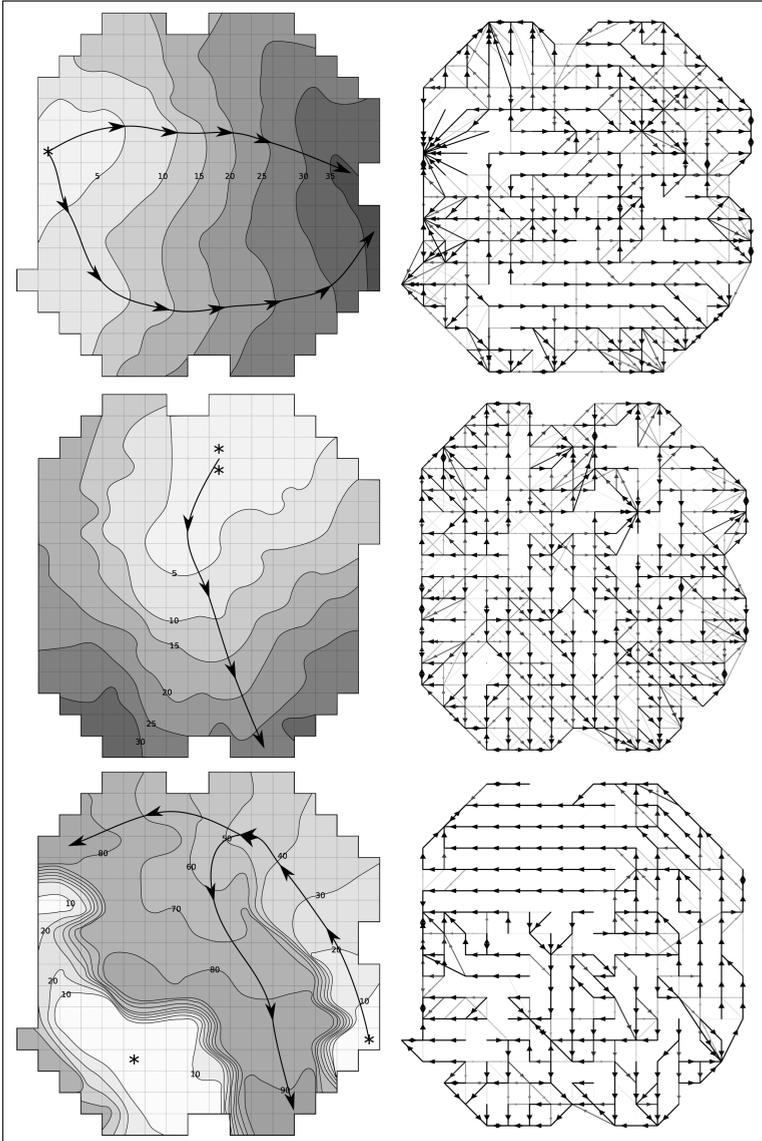


Figure 7.1: Isochrone contour maps depicting a single occurrence of three recurring patterns of AF (left) and graphs showing estimated dominant electrode interactions (right). The asterisk (*) marks the starting point of a wavefront in the contour map. The wavefront trajectories are indicated by the bold directed lines. The directed electrode interaction graphs are constructed by drawing a line between electrodes where dominant interactions occur. The strength of the interaction is indicated by a grayscale and width, ranging from black and thick (strong) to white and thin (weak).

typically shorter than one AF cycle length. A set of electrograms of duration $> M$ is analyzed by extracting L consecutive intervals of length M with overlap $M/2$ and estimating the sparse MVAR model for each of the intervals. The estimated MVAR model matrices $\tilde{\mathbf{A}}_\tau$ are then analyzed to extract information on recurring interaction patterns. We will focus on the regularization of spatial sparsity only and define the spatiotemporal weight matrix C_τ as

$$C_{\tau,ij} = \exp\left(\frac{d_{ij}}{\lambda}\right), \tau = \delta, \dots, \delta + P - 1, \quad (7.8)$$

where d_{ij} is the euclidean distance between electrodes i and j in millimeters and λ is a decay factor. The model coefficients are indicators of the strength of the interaction between two electrodes at a certain time delay, based on electrogram morphology. The mean interaction \bar{a}_{ij} from electrode j to i is defined as

$$\bar{a}_{ij} = \frac{\sum_{l=1}^L \sum_{\tau=\delta}^{\delta+P-1} \tilde{\mathbf{A}}(l)_{\tau,ij}}{PL}, \quad (7.9)$$

where $\tilde{\mathbf{A}}(l)_\tau$ denotes the interaction matrix $\tilde{\mathbf{A}}_\tau$ estimated on the interval l . Pathway maps are constructed from the mean electrode interaction matrix $\bar{\mathbf{A}}$.

7.2.3 High-density contact mapping

A subset of the mapping data presented in [99] was used for analysis. In short, goats were instrumented with an atrial endocardial pacemaker lead and a burst pacemaker. AF was maintained for 3 weeks (short-term AF [ST], $n=10$) or 6 months (long-term AF [LT], $n=7$). In an open-chest follow-up experiment, electrograms during AF were recorded from the left atrial (LA) and right atrial (RA) free walls using a round, high-density electrode array of 4cm in diameter, consisting of 234 unipolar recording electrodes with an interelectrode distance of 2.4mm (sampling rate 1kHz). Three recordings were chosen to investigate the relevance of the sparse multivariate regression approach: 1) an ST goat with a large peripheral wavefront, entering the mapping array at the same location and recurring with the same pattern for several seconds, 2) an ST goat with a recurring breakthrough wavefront pattern, and 3) an ST goat with a recurring rotating wavefront pattern.

7.3 Results

The MVAR model was estimated on 3 second recordings segmented into consecutive 100ms intervals ($M = 100$) with 50ms overlap. Model delay and order were set at $\delta = 1$ and $P = 11$, corresponding to 1ms and 11ms respectively. These values are based on the assumption that local conduction velocity between two horizontally or vertically adjacent electrodes can be as slow as 0.2mm/ms ($2.4/0.2 = 12$ ms) and should not be faster than 1.5mm/ms ($2.4/1.5 = 1.6$ ms). The value of the decay factor λ was set at 10mm. Dominant electrode interaction graphs were reconstructed by normalizing the outgoing mean interactions for each electrode to the interval $[0, 1]$ and selecting only interactions with a normalized value higher than 0.75. Fig. 7.1 shows three examples of identified recurring wave propagation patterns.

Peripheral wave train

The top isochrone map and interaction graph in Fig. 7.1 show the analysis of a recurring pattern of a peripheral wavefront entering from the left, moving from left to right and leaving the mapping area at the right. The contour map on the left shows one such wave, with higher conduction velocity at the top and bottom of the mapping array and slower conduction in the middle. The direction of the wavefront pattern is captured by the dominant interaction graph, where the top and bottom of the mapping array show clear rightward conduction patterns while the pattern in the middle is more diverse. Surprisingly the points where the wavefront enters the mapping area on the left are identified as sinks with many dominant incoming interactions. This can be explained by the observation that at an electrode where a wavefront originates there is no preceding activation at the surrounding electrodes, but they do contain later activations that still correlate with the activation at the originating electrode. At points where a wavefront has just passed, there will be a number of correlating similar activations before and after the front, but for a point where waves originate or come in, only correlations after the front are possible. When the two sides (before and after) are present, the resulting dominant interaction is likely not as strong as when only one side is present. Therefore at sources (just as at sinks) dominant interactions will more easily be pointed towards them. In this case the auto-interaction coefficient a_{ii} can provide additional information.

Repetitive breakthrough wave

The middle map and graph in Fig. 7.1 show a repetitive breakthrough wave pattern, where a wavefront originates from a deeper layer of the atrium. The activation breaks through in the right upper corner of the mapping array and subsequently shows radial spread of activation, leaving the mapping array at the left, bottom and right side. The dominant interaction graph again captures this breakthrough pattern. The upper right part of the interaction graph shows a more complex pattern, again with several electrodes that act as sinks. These electrodes are the locations where most breakthrough waves tend to enter the mapping area.

Rotating wave

The bottom map and graph in Fig. 7.1 depict the pattern of a rotating wave that enters the mapping area in the lower right, moves upward and then turns left. One part of the wavefront leaves the mapping area on the upper left, the other part keeps on turning anticlockwise to return to the point where the wavefront entered the area. The dominant interaction graph clearly shows the first upward movement of the wavefront, the left turn and the leftward exit of the wavefront. The second part of the wavefront is less clear, although several downward paths can be distinguished in the middle and lower part of the mapping area. This can be explained by the fact that the second part of the wavefront is not as recurrent as the first part.

7.4 Discussion and Conclusions

The three examples illustrate the identification of recurrent wave propagation patterns using a sparse multivariate regression model. Without the need for annotation and including only a limited amount of underlying assumptions, the developed method is able to capture the relevant dominant interactions between electrograms located at different recording locations. The distance-weighted version of the basis pursuit algorithm is a fast and promising tool to identify interactions within a short timeframe. The constructed dominant interaction graph can be further analyzed using graph theoretical algorithms to identify sources and sinks related to wave propagation, to compute maximum flow between different locations in the mapping array and to quantify graph connectivity. In a clinical setting a real-time implementation of the pattern identification algorithm might be used to guide the ablation process by identification of specific conduction patterns as ablation targets and for verification of conduction block.

Chapter 8

Systematic comparison of techniques for noninvasive assessment of atrial fibrillation complexity to predict outcome of pharmacological cardioversion in patients with recent onset atrial fibrillation

8.1 Introduction

Atrial fibrillation (AF) is a common cardiac arrhythmia that progresses in complexity over time. Current classification of AF is mainly based on AF stability and episode duration [15]. A decision on a rhythm control strategy is usually taken based on this classification of AF, also taking into account symptoms, and the doctor and patient preference. However, whether a patient will respond to rhythm control therapy is difficult to predict. Moreover, any kind of rhythm control strategy is associated with considerable risks such as ventricular pro-arrhythmia in case of anti-arrhythmic drugs or procedural risks in case of AF ablation. Predicting the acute and long-term success of AF treatment at any stage of the disease is therefore desirable and subject to extensive research [53]. One of the most urgent research questions is whether the degree of the electrophysiological changes in the atria can be assessed using noninvasive techniques. The increasing incidence of conduction block in the atria as a consequence of a progressive structural remodelling process causes an increase in number of fibrillation waves. The standard 12-lead ECG is an attractive choice for noninvasive assessment of this level of AF complexity because of its widespread use in daily clinical practice. However, whether AF complexity quantified from the surface ECG can be employed in a clinical setting to predict treatment outcome and ultimately guide management of AF, still has

to be established. Many complexity parameters have been proposed that are derived from the 12-lead ECG. Although several studies report encouraging results in either classifying AF or predicting treatment outcome, it is not a straightforward task to compare and interpret the reproducibility of these results, because of the large differences in patient populations, the parameters computed on the ECG, and the specific clinical setting. There is a clear need for standardization of ECG-based AF complexity analysis [88, 13] to be able to get closer to answering the question whether ECG-derived complexity parameters can be of value and if so, which parameters computed on which leads are useful in which clinical setting. To try and address these issues, we have performed a comparison of a large set of ECG-derived AF complexity parameters and their ability to predict successful outcome of pharmacological CV using flecainide in a population of patients with recent onset AF. Unlike previous studies, we systematically compared the performance of both time- and frequency-domain parameters, computed on single leads and multiple leads. We assessed prediction models containing a single type of complexity parameter, but also performed a thorough analysis of models containing all possible combinations of different types of parameters, grouped by parameter domain (time or frequency) and number of leads involved (one lead or multiple leads). A direct comparison of best performing ECG parameter models to clinical predictors is also provided, to quantify the added value of ECG-derived complexity parameters in the prediction of successful pharmacological CV. Long-term implications of recent onset AF complexity were investigated by associating patient complexity to the risk of progression to persistent AF.

8.2 Methods

Patient database and electrocardiogram processing

Patient data were retrieved from a patient database at Maastricht University Medical Centre, Maastricht, the Netherlands, which contained the records for patients with recent onset AF (< 48h), who underwent cardioversion with the anti-arrhythmic drug flecainide for the first time between the years 2008 and 2012. From the database a total of 221 patients were selected for this study. Exclusion criteria were the use of anti-arrhythmic drugs prior to the CV attempt, the use of additional medication during the CV procedure, and a missing ECG or an ECG of poor quality. An overview of the patient characteristics is given in Table 8.1. Echocardiographic parameters were only included if an echocardiography was performed within one year before or after the CV attempt. Before the CV attempt, for each patient a standard 10-second 12-lead ECG was recorded during AF using a GE MAC[®] 5500 resting ECG recording device at a sampling frequency of 250Hz. CV success was defined as restoration of sinus rhythm within one hour after the start of the flecainide infusion. Flecainide was dosed at 2 mg/kg with a maximum dose of 150 mg intravenously. Follow-up data on progression to persistent AF within the period 2008-March 2015 was available for 201 (out of 221) patients. Before parameter computation, signals were filtered with a 1-100Hz band-pass filter (3rd order Chebyshev, 20dB stop-band attenuation). To enable analysis of

TQ-segments, the end of the T-wave and the onset of the Q-wave were detected in the unfiltered signals using Woody's improved method [16]. To isolate the atrial signal from each lead, the ventricular QRST complexes were detected and subtracted using a single lead cancellation method based on singular value decomposition of the QRST windows [1]. Large QRS residues after cancellation were replaced by interpolated values. Single ventricular extra-systoles were blanked. The extracted atrial signal was filtered with a 3Hz high-pass filter (3th order zero phase Chebyshev filter, 20dB stop-band attenuation) to remove any remaining T-wave residues. Finally, the first and last second of the recording were truncated to avoid the border effect of filtering procedures, leaving 8 seconds available for analysis.

Noninvasive AF complexity parameters

The list of noninvasive AF complexity parameters included in this study was composed of parameters that appear frequently in the last decade of noninvasive AF complexity literature. These parameters are often computed on the extracted atrial signal (AA), but some of them can also be computed on just the (concatenated) TQ segments. Parameters were computed with algorithms provided by the original author(s) or otherwise as described in the original publication. An overview of these parameters and their domain is shown in Figure 8.1.

Spectral complexity

The frequency content of the each lead was determined by computing the spectrum for each lead using 1), the (fast) Fourier transform of the extracted atrial signal, 2) Welch's power spectral density estimate (3 segments, 1024 points, 50% overlap), and 3) the compressed spectrum (CS) [11] using the original ECG signal. The dominant frequency (DF) was defined as the frequency with the largest power within the 3-12Hz band. The organization index (OI) of the spectrum was defined as the relative contribution of the 2 largest peaks to the total spectral power. Spectral entropy (SE) is the application of Shannon's entropy to the frequency distribution and can be interpreted as a measure of uniformity of the spectrum. A high value of SE indicates high complexity. The single lead spectral analysis can be extended to a multidimensional analysis that incorporates spectral information from multiple leads using the so-called spectral envelope. The spectral envelope describes the shared spectral characteristics of a multidimensional signal [94]. This means that the spectral information from multiple leads is represented in a single spectrum. From the spectral envelope, the same three spectral parameters were derived: multidimensional dominant frequency (MDF), multidimensional organization index (MOI) and multidimensional spectral entropy (MSE) [96].

Fibrillation wave amplitude

The amplitude of fibrillation waves was determined in two ways: automatic annotation of f-waves in each lead by peak detection, followed by fibrillation wave amplitude (FWA) computation, comparable to the manual annotation method used by Nault et al. [70], and - analogous to the computation of spectral complexity - a signal envelope approach that computes a multidimensional fibrillation wave amplitude (MFWA) on

multiple leads [65]. MFWA was computed on both the AA signal and TQ segments. A low value of FWA or MFWA indicates high complexity.

Sample entropy

Sample entropy (SAE) is a time-domain parameter that quantifies the irregularity of a signal by searching for similar segments of a certain length. As proposed by Alcaraz et al. [3] SAE was computed on the main atrial wave (MAW) of each lead. The MAW is the signal resulting from filtering the atrial signal centered around the dominant frequency with a 3Hz bandwidth. A high value of SAE indicates a high complexity. Additional parameters related to the MAW are the f-wave power of the MAW (FWP MAW), with similar interpretation as FWA, and the relative sub-band energy (RHE) [3], computed as the relative energy present in the first and second harmonics of the MAW. A low value of RHE indicates high complexity.

Principal component analysis

Another multidimensional approach to AF complexity quantification is principal component analysis (PCA), which expresses the information from all 12 leads in a number of linearly uncorrelated components that essentially describe the amount of variance between the leads. Complexity measures based on PCA included were spatial complexity $k_{0.95}$, the number of components required to describe 95% of the variance in all 12 leads, and spatiotemporal stationarity (NMSE), the degree in which the three major signal components are varying over time [12]. A high value of $k_{0.95}$ and NMSE indicates high complexity. Additional measures of spatial complexity C and variability of spatial complexity CV were also included. The parameter C defines spatial complexity as the relative signal variance, excluding the three major components [61]. Frequency domain parameters derived from PCA were spectral concentration SC and spectral variability SV [61], where SC quantifies the concentration of the spectral power around the dominant frequency and SV the temporal variation of the SC. A low value of SC or a high value of SV indicates high complexity. PCA parameters were computed on both the AA signal and the TQ segments.

Prediction models and statistical analysis

Cardioversion prediction models were built using logistic regression. Prediction performance was expressed as the area under the curve (AUC) of the receiver operating characteristics (ROC) (also known as the c-statistic). Complexity parameters were first individually scored in terms of predictive performance. Then, parameters were divided into 4 groups, based on their computational domain (time or frequency) and the number of leads involved in the computation (a single lead or multiple leads). The best combination of parameters in terms of predictive performance (expressed as AUC) for each of the 4 groups of parameters was determined by iterating over all possible combinations in that group, if feasible. For larger numbers of possible parameter combinations, for instance in the case of multidimensional spectral complexity, elastic net logistic regression[105] was applied to first select a smaller subset of candidate parameters, for which again all possible combinations were evaluated. Elastic net regression

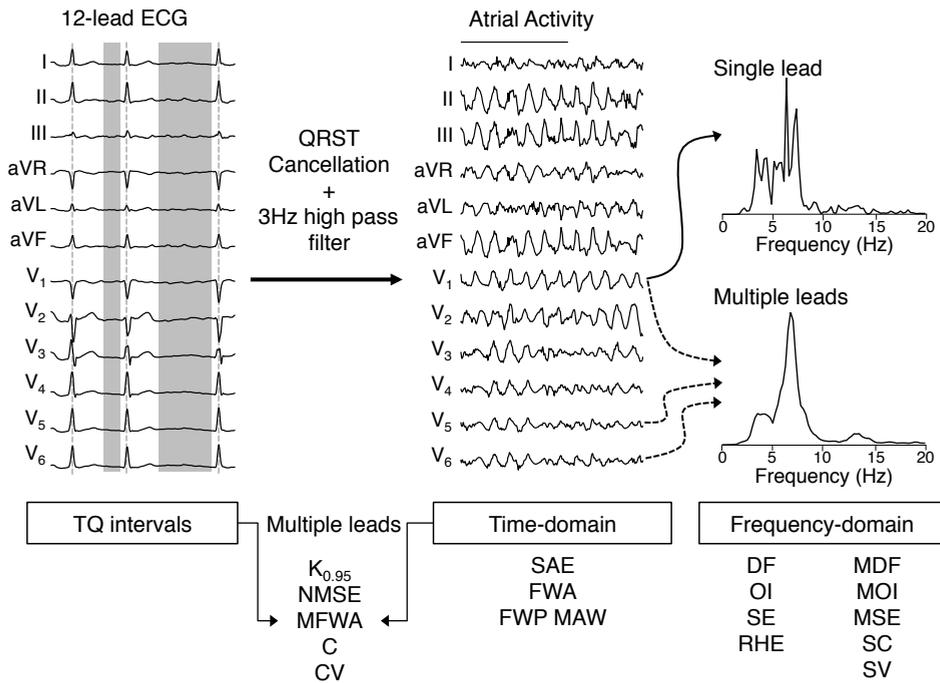


Figure 8.1: Overview of ECG signal processing and complexity parameter computation. In the time domain, multi-lead parameters can be computed on both the extracted atrial activity, as well as on the TQ-segments of the original ECG. In the frequency domain, complexity can be quantified based on a single lead or on a multi-lead spectrum.

is a regularized form of regression aimed at parameter selection, where all parameters can be entered into the model at once and it is especially appropriate in the case where there are many correlated predictors, overcoming the limitations of classical stepwise feature selection algorithms. This latter approach was also employed to select the best model from the set of all parameters. See Appendix 8.A for a more detailed description of the parameter selection procedure. Best performing parameter combinations were compared to the prediction performance that could be obtained using conventional clinical and echocardiographic predictors. Model prediction performance was cross-validated using 5-fold data partitioning with 20 Monte-Carlo repetitions. Significant differences in model performance were tested with a Student's t-test with a significance threshold of $p=0.05$. Univariate parameter differences between patients with a successful and patients with an unsuccessful CV, both ECG-derived and clinical, were tested for normality with the Lilliefors test and compared using a standard 2-tailed unpaired t-test or a Mann-Whitney U-test if the test for normality failed. The association between ECG complexity parameters and clinical parameters, and the risk of progression to persistent AF was investigated using a Cox proportional hazards model. Hazard models were estimated for each parameter individually, as well as for combinations

Table 8.1: Patient characteristics of the patients with an analyzed ECG. Numbers are given as mean \pm standard deviation or count. Between brackets is the actual number of observations for each parameter. Binary variables were tested using a χ^2 -test for proportions.

Characteristic	Successful CV <i>n</i> = 157(71%)	Unsuccessful CV <i>n</i> = 64(29%)	P-value
Sex (Male \ Female)	93\64	52\12	0.002
Age (years)	61 \pm 13	57 \pm 15	0.170
Height (cm)	174 \pm 10(116)	179 \pm 12(51)	0.004
Weight (kg)	81 \pm 14(116)	91 \pm 19(51)	0.001
BMI (kg/m ²)	26.9 \pm 3.9(116)	28.1 \pm 5.3(51)	0.305
Diabetes	11(141)	5(61)	0.924
Hypertension	66(141)	29(62)	0.996
COPD	8(141)	2(61)	0.471
PVI	2(141)	5(60)	0.014
Left atrial diameter (mm)	40.3 \pm 5.1(113)	43.1 \pm 6.0(49)	0.003
Left atrial volume (ml)	74.2 \pm 20.8(111)	80.8 \pm 19.6(48)	0.067
Right atrial volume (ml)	56.3 \pm 18.0(99)	69.9 \pm 23.3(47)	< 0.001
LVEDD (mm)	49.4 \pm 5.4(117)	51.5 \pm 6.0(51)	0.127
LVESD (mm)	33.6 \pm 4.5(116)	36.3 \pm 7.3(50)	0.064
LVEF (%)	60.1 \pm 5.6(117)	57.0 \pm 10.0(51)	0.171

COPD: chronic obstructive pulmonary disease; PVI: pulmonary vein isolation; LVEDD/LVESD: left ventricular end diastolic/systolic diameter; LVEF: left ventricular ejection fraction

of parameters, again applying an elastic net technique with 5-fold cross-validation to select a smaller subset of candidate parameters [89]. Differences in hazard model fit quality were assessed using a likelihood ratio test with a significance threshold of $p=0.05$. All computations were performed in MATLAB (MATLAB and Statistics Toolbox Release 2014a, The MathWorks, Inc., Natick, Massachusetts, United States), using custom made software and the Glmnet for MATLAB toolbox [80] for elastic net regression.

8.3 Results

Prediction using a single parameter derived from one lead

Single lead parameter results are listed in Table 8.2.

Frequency-domain (DF, OI, SE, RHE) The most significant difference between successful and unsuccessful CV, and best predictive performance was achieved using DF (Welch's spectral density estimate) at leads II, III, aVR, aVF, and V_1 with maximum single lead AUC (0.66) at lead II. Computed on the same spectrum, OI was significantly different at lead III (AUC 0.60) and SE at lead V_6 (AUC 0.59). The other spectral density estimation methods produced fewer or no significant differences and lower predictive performance. RHE showed significance at lead I (AUC 0.59). Overall, differences

Table 8.2: Significant single lead parameters and logistic regression AUC. Parameter values are reported as mean \pm SD or median (interquartile range). AUC values are given as mean \pm SD

Parameter	Lead	Successful CV	Unsuccessful CV	P-value	AUC
Frequency domain					
DF (Hz)	II	5.9(1.0)	6.3(1.1)	< 0.001	0.66 \pm 0.08
Welch	III	6.1(0.7)	6.3(1.2)	0.009	0.61 \pm 0.09
	aVR	5.9(1.2)	6.3(1.2)	0.008	0.61 \pm 0.08
	aVF	5.9(1.0)	6.3(1.5)	0.007	0.61 \pm 0.09
	V ₁	6.3(1.2)	6.6(1.6)	0.027	0.59 \pm 0.09
RHE	I	0.201(0.155)	0.166(0.106)	0.043	0.59 \pm 0.08
OI (%)	III	57.8 \pm 17.7	53.1 \pm 19.7	0.020	0.60 \pm 0.09
SE	V ₆	5.67(0.67)	5.83(0.55)	0.043	0.59 \pm 0.08
Time domain					
SAE	II	0.317 \pm 0.046	0.341 \pm 0.055	0.001	0.64 \pm 0.08
	aVF	0.323 \pm 0.060	0.345 \pm 0.061	0.010	0.61 \pm 0.09
FWP MAW	aVL	0.0064 \pm 0.0018	0.0058 \pm 0.0018	0.021	0.60 \pm 0.09
FWA (mV)	II	0.055(0.024)	0.051(0.017)	0.044	0.59 \pm 0.08
	III	0.059(0.024)	0.053(0.021)	0.007	0.62 \pm 0.09
	aVL	0.044(0.016)	0.040(0.015)	0.023	0.60 \pm 0.10
	aVF	0.053(0.024)	0.048(0.016)	0.016	0.60 \pm 0.08
	V ₆	0.038(0.012)	0.034(0.010)	0.018	0.60 \pm 0.08

between successful and unsuccessful CV were small, but as expected given the interpretation of the frequency-domain parameters. Predictive performance was equally low for all single lead frequency-domain parameters, with the exception of DF at lead II.

Time-domain (FWA, SAE, FWP) Time-domain parameters showed similar predictive performance to single lead frequency-domain parameters. Maximum significance difference between successful and unsuccessful CV was found for SAE at lead II (AUC 0.64). FWP computed on the MAW was significantly different at lead aVL, but predictive performance was low (AUC 0.60). FWA differences between successful and unsuccessful CV were significant at many leads, but differences were too small to achieve better prediction (maximum AUC 0.62 at lead III).

Prediction using a single parameter derived from multiple leads

Multidimensional parameter results are listed in Table 8.3. This analysis focused on the parameters that were computed using the information derived from multiple leads, but expressed the complexity of those multiple leads as a single parameter value.

Frequency-domain (MDF, MOI, MSE, SC, SV) The MDF, MOI and MSE parameters were computed on the spectral envelope of all possible combinations of 2 or more precordial leads (57 for each parameter in total), as opposed to only pairs of leads in

Table 8.3: Significant multidimensional parameters and prediction AUC. Parameter values are reported as mean \pm SD or median (interquartile range). AUC values are given as mean \pm SD

Parameter	Leads \ Signal	Successful CV	Unsuccessful CV	P-value	AUC
Frequency domain					
MDF (Hz)	V _(2,5)	6.0(1.0)	6.3(1.3)	0.008	0.61 \pm 0.09
Top 4 (42)	V _(1,4,5)	6.0(1.3)	6.8(1.5)	0.003	0.63 \pm 0.09
	V _(1,2,4,6)	6.0(1.3)	6.5(1.4)	0.001	0.64 \pm 0.08
MOI (%)	V _(1,2,4,5,6)	6.0(1.3)	6.5(1.3)	0.001	0.65 \pm 0.08
	V _(3,4)	50.6 \pm 8.7	47.6 \pm 6.9	0.015	0.61 \pm 0.08
Top 3 (22)	V _(1,2,4)	53.8(10.2)	51.2(10.0)	0.005	0.62 \pm 0.08
	V _(2,4,5,6)	41.6(8.0)	38.9(6.2)	0.005	0.62 \pm 0.08
MSE	V _(3,4)	6.37 \pm 0.39	6.49(0.34)	0.046	0.60 \pm 0.08
SC (%)	All leads	23.8(1.4)	23.5(1.5)	0.031	0.59 \pm 0.09
SV	All leads	0.51(0.26)	0.66(0.37)	0.001	0.65 \pm 0.08
Time domain					
$k_{0.95}$	AA	4.8(0.8)	5.0(0.8)	0.033	0.59 \pm 0.09
	TQ	3.2(0.6)	3.4(0.4)	0.050	0.58 \pm 0.08
MFWA _{median}	AA	0.049(0.037)	0.040(0.031)	0.025	0.60 \pm 0.09
MFWA _{mean}	AA	0.049(0.033)	0.039(0.024)	0.031	0.59 \pm 0.08
C	AA	9.4(3.5)	10.7(4.3)	0.021	0.60 \pm 0.08
	TQ	4.6 \pm 1.8	5.2 \pm 2.0	0.050	0.59 \pm 0.09
CV	AA	2.5 \pm 1.0	3.0 \pm 1.5	0.005	0.59 \pm 0.08
	TQ	2.8 \pm 1.2	3.2 \pm 1.3	0.076	0.58 \pm 0.09

Uldry et al. [96]. Results show that the predictive power of MDF increases when more leads are included, but only reaches a maximum AUC of 0.65 for the combination of V_(1,2,4,5,6), still lower than single lead DF performance on limb lead II. Many MOI combinations resulted in statistically significant differences between successful and unsuccessful CV, but only reached a maximum AUC of 0.62. MSE was only significant for the combination V_(3,4), with a low AUC of 0.60. SC and SV were both significant, with a higher AUC for SV (0.65 vs. 0.60).

Time-domain ($k_{0.95}$, NMSE, MFWA, C, CV) The time-domain multidimensional parameters $k_{0.95}$, C and CV attained significant differences and predictive performance, both computed on the AA signal as well as on the concatenated TQ segments, but with equally low predictive performance (AUC between 0.58 and 0.60). MFWA was only significant when computed on the full AA signal. Also here predictive performance was low (AUC 0.60). Overall, single lead and multidimensional time-domain predictive performance was very similar, and generally lower than frequency-domain parameter performance.

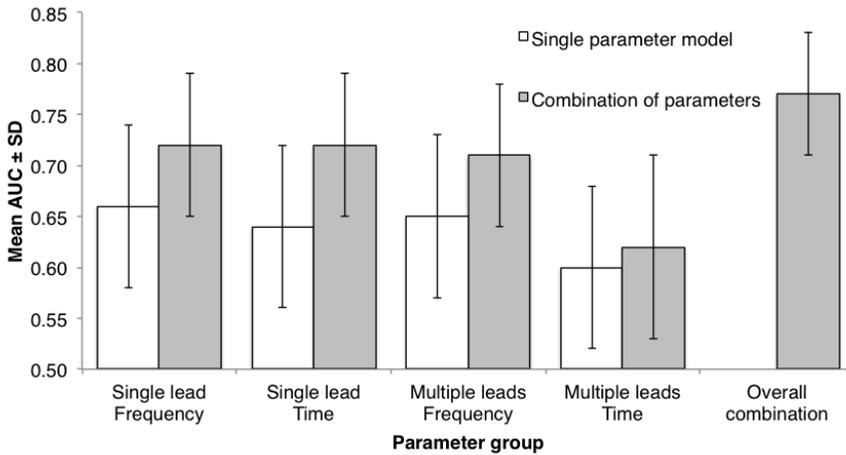


Figure 8.2: Prediction performance of different types of complexity parameters computed on one lead and multiple leads, in the frequency-domain and the time-domain. Performance is expressed as the mean $AUC \pm SD$ for each of the 4 groups (parameter computed in the frequency/time-domain and on one/multiple leads). For every group the best performing single parameter model AUC is given, as well as AUC of the best combination of parameters belonging to the same group.

Prediction using a combination of ECG parameters

Results for the parameter models consisting of a combination of ECG parameters are listed in Table 8.4. In this analysis we determined the combination of parameters in each of the 4 groups (time- or frequency-domain and single lead parameters or multidimensional parameters) that led to the optimal predictive performance. Out of these group results, we also distilled the optimal combination of all ECG parameters under investigation. The best model containing a combination of frequency-domain parameters computed on a single lead improved performance from an AUC of 0.66 to 0.72 ($p < 0.001$), by adding OI (lead III) and SE (lead I) to the best single lead parameter DF (lead II). In the time-domain prediction also improved by extending the best performing single lead parameter SAE (lead II) with FWA (lead aVF and V_1) and FWP (lead V_2), from an AUC of 0.64 to 0.72 ($p < 0.001$). Combining multidimensional parameters produced similar results in the frequency-domain, with a 6-parameter model increasing the AUC from 0.65 (MDF on $V_{(1,2,4,5,6)}$ or SV) to 0.71 ($p < 0.001$). In the time-domain predictive performance remained poor (MFWA and CV, AUC 0.60 vs. 0.62, $p = 0.082$). Combining the best predicting parameters of each group of parameters in a single model further improved prediction performance to an AUC of 0.78, with parameters selected from each group. Figure 8.2 shows the effect of combining parameters within each parameter group on the prediction performance. A selection of cross-validated ROC curves, from which the AUC values were derived, is depicted in Figure 8.3.

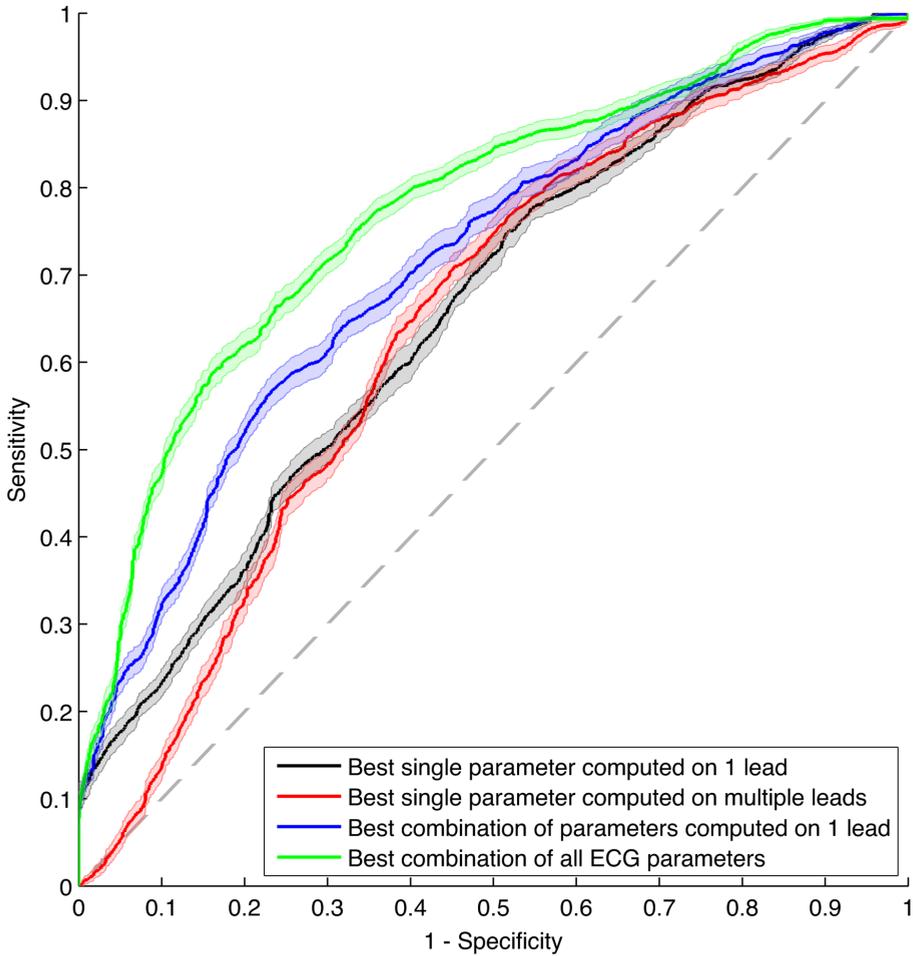


Figure 8.3: Cross-validated ROC curves of various ECG parameter models listed in Tables 8.2, 8.3, and 8.4. The narrow band around each of the curves indicates the 95% confidence interval of the sensitivity for a given specificity. Depicted are the ROC curves for the best prediction performance of a single parameter computed on 1 lead (DF on lead II, AUC 0.66, black line), the best performance for a single multidimensional parameter (SV derived from all leads, AUC 0.65, red line), the best performing combination of parameters computed on 1 lead (DF (II), OI (III) and SE (I), AUC 0.72, blue line), and the best combination of all ECG parameters (DF (II), SE (I), FWA (aVF, V_1), MOI ($V_{(3,4)}$, $V_{(3,5)}$), SV, and MFWA, AUC 0.78, green line).

Table 8.4: Best performing parameter models for single lead and multidimensional parameters. For each group of parameters the selected parameters and leads are given together with the prediction $AUC \pm SD$ of the estimated model.

Group	Parameters	Leads or signal	AUC
Single lead	DF	II	0.72 ± 0.07
Frequency domain	OI	III	
	SE	I	
Single lead	SAE	II	0.72 ± 0.07
Time domain	FWA	aVF, V ₁	
	FWP MAW	V ₂	
Multiple leads	MDF	V _(1,2,4,5) , V _(1,2,4,5,6)	0.71 ± 0.07
Frequency domain	MOI	V _(3,4) , V _(3,5) , V _(2,4,5,6)	
	SV	All leads	
Multiple leads	MFWA _{median}	AA	0.62 ± 0.09
Time domain	CV	AA	
Combined	DF (II), SE (I), FWA (aVF, V ₁), MOI (V _(3,4) , V _(3,5)), SV, MFWA median		0.78 ± 0.06

Prediction using clinical parameters and ECG parameters

An overview of the predictive performance obtained when combining clinical patient characteristics and ECG parameters can be found in Table 8.5. In the smaller subset of patients with complete clinical characteristics and echocardiographic data available ($n = 139$), the predictive capability of clinical parameters alone was limited. The combination of patient weight and right atrial volume (RAV) performed best with a mean AUC of 0.68. The predictive performance in this patient subset of the prediction models based on the 4 ECG parameter groups outperformed the model using the best combination of all available clinical parameters, except for the multidimensional time-domain parameter model. Adding the ECG parameters from the best performing group models to the clinical parameters significantly enhanced predictive performance in all cases, again except for the multidimensional time-domain parameter model. Maximum AUC was reached by combining weight and RAV with single lead frequency-domain parameters (AUC 0.81, $p < 0.001$ compared to clinical). Combining the clinical parameters with the best performing single lead and multidimensional ECG parameter model, as determined on the full ECG data set, did not further improve prediction performance (AUC 0.78). The ROC curves demonstrating the added value of the ECG parameters are shown in Figure 8.4.

8.4 Risk of progression to persistent AF

Out of the 201 patients for whom follow-up was available, 38 (19%) developed persistent AF between the date of CV and March 2015 (median survival time: 408 days, interquartile range (IQR): 171-822 days). Table 8.6 contains the significant hazard ratios (HR) for individual clinical and ECG complexity parameters, and the parameter unit increment used in the computation of the HR. Age, body mass index (BMI), left atrial

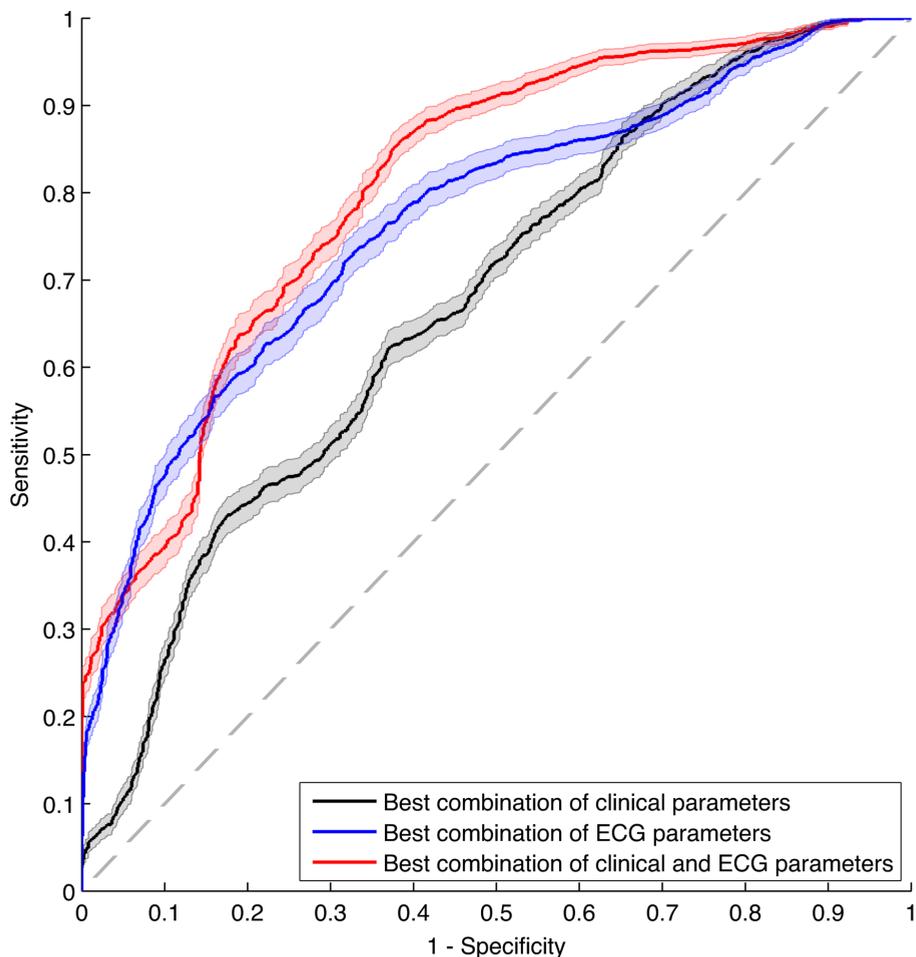


Figure 8.4: Cross-validated ROC curves showing the difference in prediction performance between the best model consisting solely of clinical parameters (weight and RAV, AUC 0.68, black line) and the best model consisting only of parameters computed on the ECG (best ECG parameter model derived from the full data set, AUC 0.78, blue line), and the performance of the best combination of clinical parameters and ECG parameters (weight, RAV and the best single lead frequency-domain parameter model, AUC 0.81, red line). The bands around the curves indicate the 95% confidence interval of the sensitivity for a given specificity.

Table 8.5: Predictive performance of clinical parameters and the added value of the best performing single lead and multidimensional ECG parameters models, as determined on the full data set (see Table 8.4). P-values relate to the comparison between the model consisting of only clinical parameters and the specific combination. P-values between brackets relate to the difference between the model consisting of ECG parameters and the combined model of clinical and ECG parameters.

Parameter model	AUC on subset	AUC clinical & ECG parameters	P-value
Clinical parameters (Weight, RAV)	0.68±0.11	N/A	N/A
Single lead	0.76±0.09	0.81±0.07	< 0.001
Frequency domain			(< 0.001)
Single lead	0.73±0.09	0.77±0.09	< 0.001
Time domain			(< 0.001)
Multiple leads	0.73±0.10	0.77±0.09	< 0.001
Frequency domain			(< 0.001)
Multiple leads	0.60±0.11	0.68±0.09	0.289
Time domain			(< 0.001)
Best ECG model	0.78±0.08	0.78±0.08	< 0.001 (0.473)

diameter (LAD), RAV, left ventricular end systolic diameter (LVESD) and left ventricular ejection fraction (LVEF) showed small, but significant hazard ratios. Unsuccessful CV was not a significant hazard (HR 1.58, 95% confidence interval (CI) 0.82-3.06, $p=0.17$). ECG complexity parameters were only significant for DF (on leads III, aVL, avF, V_4) and FWA (V_1). Both a higher DF and a higher FWA were associated with a larger risk of developing persistent AF. Correcting for age and sex only eliminated DF on lead V_4 as a significant hazard. Figure 8.5 depicts the Kaplan-Meier curves for four dichotomized parameters BMI, LAD, DF (aVL) and FWA (V_1), showing that an increased risk of progression to persistent AF was associated with obesity ($\text{BMI} > 30 \text{ kg/m}^2$), an enlarged left atrium ($\text{LAD} > 41 \text{ mm}$), a faster atrial rate ($\text{DF} > 5.7 \text{ Hz}$), and a higher f-wave amplitude ($\text{FWA} > 0.06 \text{ mV}$). Progression to AF was significantly faster for patients with $\text{FWA} > 0.06 \text{ mV}$ (median survival time 296 days vs. 796 days, $p=0.03$). Multivariate analysis showed that the risk of progression to persistent AF is best described by a model containing LAD ($n=155$), when considering only clinical parameters. ECG complexity parameters modelled progression best using a combination of DF (lead aVL) and FWA (lead V_1) ($n=201$, DF(aVL): HR 1.45, CI 1.08-1.94, $p=0.01$; FWA(V_1): HR 1.16, CI 1.05-1.27, $p<0.01$). Adding DF (aVL) or FWA (V_1) to the best model containing only clinical parameters both improved the model fit ($p=0.05$ or 0.02 respectively).

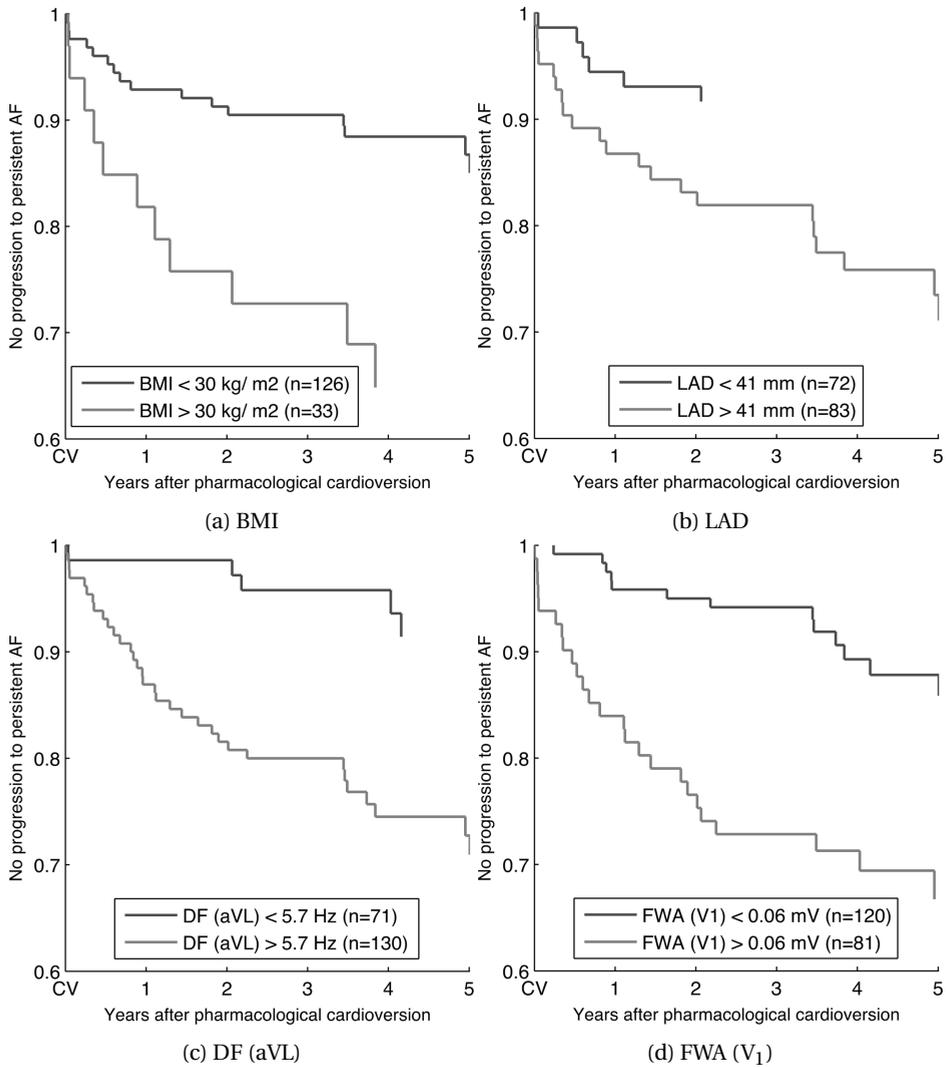


Figure 8.5: Kaplan-Meier curves for the risk of progression to persistent AF after the CV attempt for patients with a) BMI > 30 kg/m² (HR 2.97, CI 1.38-6.40), b) LAD > 41 mm (HR 2.65 CI 1.16-6.06), c) DF > 5.7 Hz on lead aVL (HR 4.16, CI 1.62-10.67), and d) FWA > 0.06 mV (HR 3.25, CI 1.66-6.35). All HR are significant with $p < 0.01$.

Table 8.6: Significant hazard ratios for risk of progression to persistent AF

Parameter	Increment	Hazard ratios (95% CI)	
		Unadjusted	Adjusted for Sex and Age
Age	1 year	1.03 (1.01-1.06)*	N/A
BMI (n=159)	1 kg/m ²	1.09 (1.03-1.17)#	1.10 (1.03-1.18)#
COPD (n=185)	N/A	3.59 (1.26-10.21)*	3.38 (1.19-9.61)*
LAD (n=155)	1 mm	1.12 (1.05-1.19)#	1.11 (1.04-1.18)#
RAV (n=141)	5 ml	1.11 (1.03-1.21)*	1.12 (1.02-1.22)*
LVEDD (n=158)	1 mm	1.06 (1.00-1.11)*	1.08 (1.02-1.14)#
LVEF (n=160)	-1 %	1.06 (1.02-1.11)#	1.06 (1.02-1.10)#
DF (III)	1 Hz	1.57 (1.17-2.10)*	1.65 (1.24-2.19)#
DF (aVL)	1 Hz	1.50 (1.14-1.99)#	1.64 (1.24-2.16)#
DF (aVF)	1 Hz	1.46 (1.11-1.92)#	1.53 (1.17-2.00)*
DF (V ₄)	1 Hz	1.37 (1.03-1.87)*	1.34 (0.99-1.83)
FWA (V ₁)	0.01 mV	1.17 (1.07-1.29)#	1.16 (1.06-1.27)#

CI: Confidence interval; * $p < 0.05$, # $p < 0.01$

8.5 Discussion

ECG parameters as predictors of pharmacological cardioversion outcome

The results of the single and multidimensional ECG lead analysis show that the 12-lead ECG contains valuable information related to the response to pharmacological CV of recent-onset AF. This information may be used to characterize complexity of AF at an early stage. As expected patients with a lower DF are more likely to respond to treatment, an observation that corroborates the findings of Choudhary et al. [21] who showed that recent-onset AF patients with a lower atrial fibrillatory rate (AFR) were more likely to spontaneously cardiovert. Moreover, significant single lead measures of organization of the signal spectrum (OI and SE) indicate that a higher degree of organization favours successful CV. In the time-domain, SAE on lead II provides the most significant difference between groups, with lower SAE associated with higher chance of CV success, which is consistent with the results of Alcaraz et al. in a study on prediction of spontaneous CV [2]. Using multidimensional parameters that compute a single value from multiple leads is a logical extension of single lead analysis. In theory, incorporating spatial differences among leads and capturing inter-lead variability as an additional measure of complexity should provide a more robust estimate of AF complexity. The results from the multidimensional parameter analysis indeed confirm this to some extent. While DF computed on one of the precordial leads only gives a significant result on lead V₁, the multidimensional extension MDF performs better, with significant results for many combinations of the precordial leads (42 significant combinations in total). Maximum MDF AUC is however still lower than the maximum AUC of DF on limb lead II (0.65 vs. 0.66). The same holds true for OI and MOI, apart from the fact that MOI has a slightly higher maximum AUC on a combination of the precordial leads than the maximum AUC of OI on lead III (0.62 vs. 0.60). We did not notice an im-

portant role of left atrial content in this patient population as indicated by Uldry et al. in their study on discriminating persistent and long-standing persistent AF [96]. The most significant multidimensional parameter differences consisted of a mix of right- and left-oriented precordial leads. Although differences were indicating a higher degree of organization in the successful CV group, they were too subtle to use for classification or prediction purposes. The significant result for SV (AUC 0.65) indicates that a higher degree of temporal regularity in the frequency-domain is also a determinant for successful CV. Multidimensional parameters in the time-domain appear to have a lower performance in this dataset, but we do see several significant differences that point to lower complexity of AF in patients with successful CV, such as the higher multidimensional f-wave amplitude and lower number of signal components $k_{0.95}$ needed to describe the signal. Overall however, predictive performance of single parameters in the time-domain is low. This is likely due to the very subtle differences in AF complexity between patients in this stage of recent-onset AF. Combining several complexity parameters in a prediction model significantly improves prediction, regardless of whether these different parameter values are calculated from single lead or multiple leads. Worthwhile noting is that the best combination of frequency-domain parameters computed on a single lead is composed completely of limb leads I, II and III, with a strong role of DF at lead II, again suggesting the need to include leads that contain both right and left atrial activity. The best model combining time-domain parameters computed on a single lead seems to reach a similar performance, but this is partially driven by the correlation between DF and SAE computed on the MAW (correlation DF and SAE at lead II: $r = 0.79$, $p < 0.001$, mean correlation DF and SAE: $r = 0.75 \pm 0.03$). This correlation, as already noted by Platonov et al. [79], can be explained by the fact that the SAE is computed on a signal that has been filtered around the DF, making it likely for a signal with a higher DF to have a higher SAE. The best single lead time-domain parameter model without SAE has a somewhat lower predictive performance (AUC: 0.71 ± 0.08 vs. 0.72 ± 0.07 , FWA(aVF, V_1), FWP MAW(aVL, aVF, V_2)).

Added predictive value of ECG parameters compared to clinical information

The ability of clinical parameters (including echocardiographic parameters) to predict successful outcome of CV was limited. Combinations of ECG parameters performed better on the subset of patients with complete clinical and echocardiographic data records. Combining ECG and clinical parameters further improved prediction. This implies that features extracted from the ECG contain complementary information to the available clinical characteristics in this patient population. In particular, single lead frequency-domain parameters improved prediction.

Noninvasive complexity and risk of progression to persistent AF

Both clinical as well as ECG complexity parameters were associated with risk of progression to persistent AF. Clinical parameters like age and BMI, and echocardiographic parameters like LAD, RAV and LVEF were indicators for an increased risk of progression

to persistent AF, which is largely in line with previous findings [48]. From the set of ECG parameters only parameters computed on a single lead showed significant hazard ratios, namely DF and FWA. The threshold of 5.7 Hz computed for DF on lead aVL to produce the survival curve in Figure 8.5c is very comparable to the AFR threshold of <350 fibrillations per minute (5.8 Hz) found by Choudhary et al. [21] associated with a significant increase in the likelihood of spontaneous cardioversion of recent onset AF within 18 hours. One could argue that patients that are not likely to spontaneously cardiovert have a higher risk to develop persistent AF, due to more electrical and, eventually, structural remodelling caused by prolonged episodes of AF. The role of FWA on V_1 in the development of persistent AF in this patient cohort is more challenging to interpret: a higher FWA was associated with a higher risk for persistent AF, while the inverse relation was found for FWA in the prediction of successful pharmacological CV. To our knowledge, the long-term implications of FWA determined in patients in recent-onset AF have not yet been studied. Within the set of patient with follow-up, the combination of available clinical and echocardiographic parameters that best described FWA on V_1 consisted only of LVEF, but with a low correlation coefficient ($r=-0.20$, $p=0.018$). All other clinical parameters showed no significant (linear) association with FWA on V_1 . Atrial dimensions did not correlate well with FWA on V_1 (LAD $r=0.16$, $p=0.05$; LAV $r=0.04$, $p=0.66$; RAV $r=0.06$, $p=0.50$). FWA on V_1 was also not related to left-right atrial DF or OI differences, as one might hypothesize based on the idea that increased left atrial AF complexity can lead to a reduced cancellation of the f-waves present in the right atrium.

Limitations

The retrospective nature of this study has obvious implications for the availability and quality of both clinical information and ECG signals. Echocardiography was not recorded at the same time as the ECG, but selecting an available echocardiography within a year produced similar results compared to a narrower timeframe (see Appendix 8.B for an analysis of the echocardiography timeframe). ECG signals were not recorded with the intention to analyse AF, meaning that quality was varying and recording length was limited to 10 seconds.

8.A Parameter selection via elastic net logistic regression

In several cases, the number of candidate parameters in the logistic regression model makes it infeasible to iterate over all possible parameter combinations to select the overall best performing model. Parameter selection using stepwise logistic regression has the disadvantage that it is dependent on the order in which parameters are added or removed from the prediction model. Stepwise parameter selection is also affected by parameter correlation. To select dominant parameters from a large set of candidate parameters, and to overcome the limitations of stepwise methods we applied an approach that combines information from classical stepwise logistic regression and elastic net logistic regression. Elastic net regression is based on mixed ℓ_1/ℓ_2 -norm

regularization of the parameter coefficients in the criterion function of the regression model at hand. This regularization aims to minimize the number of non-zero parameter coefficients in the estimated model. Given a certain output data y of length N , in the case of logistic regression the objective is to minimize the model deviance

$$D(y, \theta) = -2(\log(p(y|\theta)) - \log(p(y|\theta_s))), \quad (8.1)$$

where the vector θ contains the parameter coefficients and θ_s denotes the parameter vector of the saturated model. The formulation for the elastic net logistic regression problem is

$$\min_{\theta} \left(\frac{1}{N} D(y, \theta) + \lambda P_{\alpha}(\theta) \right), \quad \text{with} \quad (8.2)$$

$$P_{\alpha}(\theta) = \frac{1 - \alpha}{2} \|\theta\|_2^2 + \alpha \|\theta\|_1.$$

The two regression tuning parameters are λ and α . The parameter λ determines the strength of the regularization of the parameter coefficients, while α (a value between 0 and 1) controls the balance between penalizing either the ℓ_2 - and/or the ℓ_1 -norm of the coefficient vector [34]. Several steps of the parameter selection procedure are outlined in Figure 8.6. In the analysis shown there the set of parameters under investigation was the group of parameters computed on 1 lead in the frequency domain (DF, OI, SE and RHE). Figure 8.6a) and b) show the elastic net estimation result for a fixed value of α ($\alpha = 0.5$). The choice of λ influences the estimated parameter coefficients and the deviance of the estimated model. A commonly accepted choice for λ is the value that corresponds to a model deviation that lies within 1 standard deviation of the cross-validated minimum deviation. These λ -values are indicated with a gray (minimum deviation) and a black line (minimum deviation + 1 standard deviation). The choice of α also determines the number of parameters that are selected. For $\alpha = 1$ the algorithm corresponds the Lasso algorithm, which tends to select one parameter from a group of correlated parameters, but for α values between 0 and 1, the elastic net algorithm will include more correlated parameters. Therefore a range of alpha (between 0.1 and 1) was investigated and for each value of α the non-zero parameter coefficients were stored (see Figure 8.6c). Parameters that appeared in any of the models computed with this range of α were considered potential candidates for the final logistic regression model. In this case the parameters DF(on leads II, aVR and V₄), OI (leads I and III) and RHE (lead I) were selected. As an additional step, parameters were also selected through forward stepwise logistic regression ($p < 0.05$ for significant deviance improvement by adding a parameter). In this case selected parameters were DF (lead II), OI (III) and SE (I). The union of the parameters selected by the two regression methods was then taken to iterate over all possible combinations of parameters to find the model with the best prediction performance. Figure 8.6d) shows the result of this last step. The model performance increased by adding more parameters, but reached a maximum at a model containing 3 parameters (DF (II), OI (III) and SE(I)).

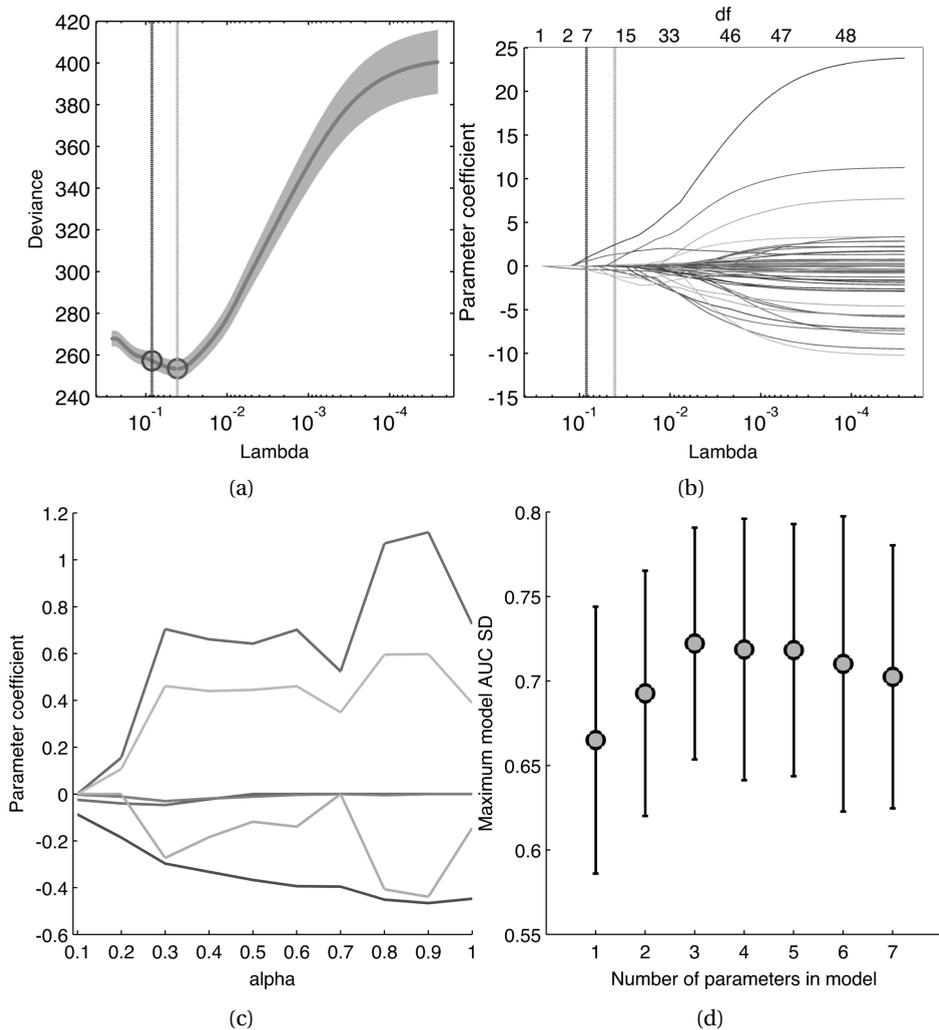


Figure 8.6: Parameter selection via elastic net logistic regression (parameters computed on a single lead, frequency-domain). The upper plots show the result of the analysis for $\alpha=0.5$, with in a) the cross-validated deviance as a function of λ and in b) the parameters coefficients (df indicates the number of non-zero parameter coefficients). The gray vertical line/circle marks the choice of λ that minimizes the deviance, the black vertical line/circle marks the solution that is within 1 standard deviation. Panel c) shows the non-zero parameters coefficients selected by the elastic net regression as a function of α . Panel d) contains the result for the cross-validated maximum AUC for models composed of a specific number of candidate parameters, defined by union of the stepwise regression and elastic net parameter selection.

8.B Effect of time interval between echocardiography and CV attempt

In our analysis we included echocardiographic data that was collected within a year (365 days) of the date of the CV attempt. In this analysis we also included patients without an ECG or a poor quality ECG before the CV attempt (n=198). Results are shown in Figure 8.7. From Figure 8.7a it becomes clear that the number of patients that can be included in the analysis based on their echocardiographic data, initially decreases slowly when we move from 365 days to a narrower timeframe. This decrease accelerates when we reach 100 days as a cut-off value. The performance of the best model containing only clinical parameters (weight and right atrial volume (RAV)), shown in Figure 8.7b, remains relatively stable until 100 days, and then starts to increase, but also becomes more irregular, due to the lower number of patients included in the analysis. This observation is supported by examining the evolution of the two clinical parameters forming the best performing model, as shown in Figure 8.7c and 8.7d.

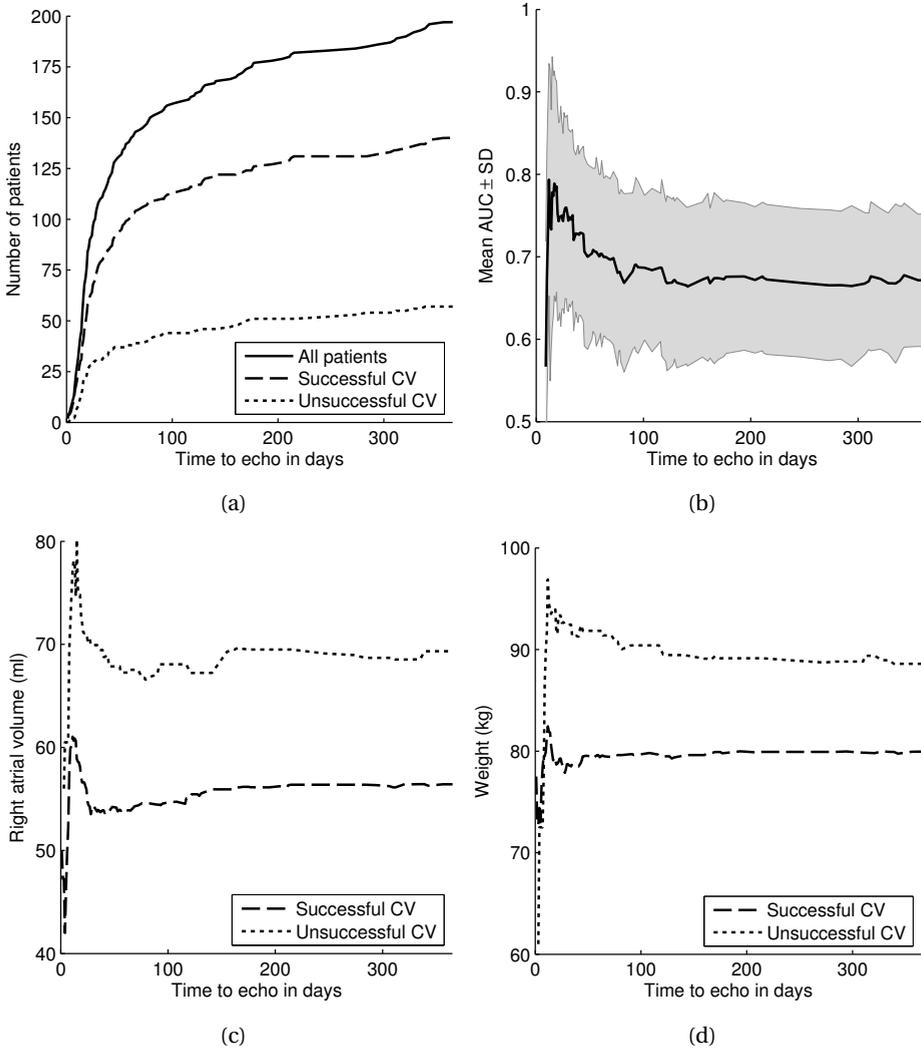


Figure 8.7: The effect of maximum allowed time difference between the date of cardioversion and the closest date of an echocardiography on (a) the number of patients included in the analysis, (b) prediction performance of the best performing model containing clinical parameters weight and RAV, (c) differences in patient RAV (successful and unsuccessful CV), and (d) differences in patient weight.

Chapter 9

General discussion

9.1 Sparse estimation

In this thesis it was investigated to what extent sparse linear, time-independent systems can be estimated given a limited amount of available measurement data. In linear regression minimizing the ℓ_1 -norm of the parameter vector within the input-output equivalence space, provided a good heuristic to correctly estimate the data generating model parameters, even in a strongly underdetermined setting. In system identification, this method was employed to estimate (structured) state-space matrices using an iterative mixed ℓ_2/ℓ_1 minimization procedure, where conventional prediction error minimization is followed by ℓ_1 minimization of the parameter vector within the local input-output equivalence space, defined either by the Jacobian matrix of the error vector or the Hessian matrix of the prediction error criterion. The ability to correctly estimate sparse network interactions was also evaluated, using sparse linear regression in the discrete-time interaction networks, and the iterative version of the algorithm in continuous-time networks.

9.1.1 Sparse linear regression

In Chapter 3 the error rates and probability of a correct parameter estimate were investigated in the class of linear regression models. Given the theoretical bound on the sparsity of the parameter vector needed to ensure that the vector with minimal ℓ_1 -norm corresponds to the sparsest parameter vector, a relative low amount of measurements were needed on average to correctly estimate the data generating parameters. The strictness of this theoretical bound was already noted in [28], where it was shown that especially in large underdetermined systems of equations, the minimum ℓ_1 -norm solution will typically also be the sparsest solution. In a noiseless setting, the method outperformed the lasso, a method that employs a parameter shrinking approach to maximizing sparsity. It was therefore considered a good candidate for sparse maximization in nonlinear optimization as well. The approaches to sparse linear regression mentioned in this thesis, represent only a subset of the recent advances that

have been made in this field of research. The lasso algorithm has been extended to work for grouped variables called the group lasso [104] and sparse-group lasso [90]. Component lasso [46] is another extension that groups variables based on the sample covariance matrix. Sparse approaches have been developed for the class of generalized linear models, to include for instance sparse logistic regression and Cox's proportional hazards analysis [89]. Moreover, software packages have been developed that contain efficient implementations of several sparse regression algorithms (see for instance [80]).

9.1.2 Sparse state-space identification

In Chapter 4 an extension of the sparse linear regression algorithm was proposed. An iterative version of mixed ℓ_2/ℓ_1 minimization was able to increase model sparsity, starting from a nonsparse data-generating model. The choice of the size of the step in the sparse minimization direction was however critical for the convergence of the algorithm to a specific solution. In a fully parameterized state-space model, an approach could be used that was guided by the data-drive local coordinates technique, which guaranteed that the mixed ℓ_2/ℓ_1 minimization parameter trajectory stayed within the space of state-space systems that can be described by a state coordinate transformation. Here, the existence of an input-output equivalence space is inherent to the parameterization of the model and not necessarily related to an undetermined setting. Note also that minimization of a mixed ℓ_2/ℓ_1 -norm of the parameter vector in a fully parameterized state-space model may lead to (continuous-time) balanced realizations, see [40] and [100]. Balanced realizations are well-known to exhibit good numerical conditioning properties in many applications (cf., e.g., [68, 73, 78]) and the approach presented here may also be of practical value in situations where the system order is relatively low and sparsity does not apply. Possible additional applications of sparse state-space identification lie in the field of model reduction by eliminating redundant interactions, in the field of coordination control, by identifying nested structures or groups of uncoupled of states within large-scale systems, which may enable decomposition into a hierarchical structure (see [49] and references given there for an overview of coordination control of linear systems), and time-dependent role identification, by identifying activating or inhibiting components in a system.

9.1.3 Sparse network identification

To be able to generate sparse network structures, a specific ring-like topology was designed, corresponding to a tri-diagonal state interaction matrix. Within this prescribed structure, relatively large, stable and sparse networks could be generated. In networks that are sparse in discrete time, sparse linear regression could be applied for each network node, leading to similar performance as in the simple linear regression setting. For a network with sparse interactions in continuous time, two approaches were investigated. A sparse network in continuous time sampled at a sufficiently high sampling time will still yield a sparse discretized interaction matrix. The estimated discretized interaction matrix (estimated by sparse linear regression) can be transformed

back to continuous time to try and retrieve the original continuous-time interaction structure. Experiments indicated that this is indeed possible at appropriate sampling rates, in combination with a sufficiently exciting input signal and a fast enough impulse response of the generated network. The number of measurements needed to correctly estimate the data generating model parameters is however much larger than in the case of sparse linear regression. The risk of over- or under-sampling is also present and the range of sampling rates for which network interactions were estimated correctly was narrow. A second, and potentially more viable approach was to estimate the continuous-time network parameters directly, given a certain sampling rate. This also proved to be feasible, but here the convergence to a stable sparse network estimate turned out to be problematic. In an underdetermined setting the mixed ℓ_2/ℓ_1 minimization algorithm indeed converged to a sparse solution, but unfortunately this solution often corresponded to an unstable network configuration. When the data generating system was the sparsest solution given the number of measurements at hand, convergence to this solution was achieved, but the number of measurements needed to create such a situation was already close to the number of measurements sufficient to find a unique solution. The application of sparse network identification is feasible, assuming sparsity in discrete time. Continuous-time sparse estimation applicability is limited by the interplay between the time-resolution needed to observe the relevant model behaviour, the number of available measurements, and the optimality of the data generating model parameter sparsity (see Garnier et al. [39, 20] for more details on continuous-time system identification techniques). In this thesis only a specific type of regular network structure was investigated in the evaluation of the applicability of the (mixed) sparse estimation procedure. Although several observations that were made for this type of network are likely equivalent for other network structures, like the dependence on the sampling rate and the time-scale of the model behaviour, other sparse network topologies might possess more (or less) favourable properties when it comes to estimating the network interactions in an underdetermined setting. The main challenge here lies in generating stable and minimal networks with specific properties, like small-world, scale-free or random networks (see for instance [71] for a review of network structures), that are large enough to fulfil the sparsity assumption. The mixed ℓ_2/ℓ_1 minimization algorithm as presented here does not take into account the stability of the estimated solution. One could envision a regularized adaptation of the algorithm to promote stable solutions by for instance including the infinity-norm (H_∞) of the state-space system corresponding to the current parameter vector estimate, as a regularization term in the prediction error minimization criterion. This however will probably lead to estimates that are only marginally stable, as the trajectory followed in the mixed ℓ_2/ℓ_1 minimization procedure when starting from a stable initial solution and converging to maximally sparse unstable solution, contains such a marginally stable solution.

9.2 Applications in atrial fibrillation

Several techniques investigated in the analysis of sparse estimation approaches have potentially interesting applications in the analysis of atrial fibrillation, from electrogram interaction analysis to feature selection in prediction of treatment outcome.

9.2.1 Recurrent propagation pattern identification

In Chapter 7 the sparse linear regression algorithm was adapted to identify time-delayed interactions between direct contact electrograms in goats, recorded with a high-density grid of electrodes. A distance-weighted version of this algorithm identified dominant local electrode interactions within a short time frame. By averaging these sparse electrode interaction matrices recurrent propagation patterns could be extracted that corresponded to patterns reconstructed by manual annotation of local atrial deflections. This approach is promising for several reasons. First of all, reconstructing activation patterns through annotation (manually or automatically as described in Chapter 6), depends on accurate detection of local atrial deflections, a task that can be time-consuming and/or subject to uncertainty, especially in more complex (i.e. fractionated) electrograms. The sparse multivariate autoregression approach has the advantage that it is designed to take into account both instantaneous and time-delayed coupling between electrodes, without having to assign atrial deflections and corresponding activation times. Furthermore, it is able to incorporate spatial information to regularize the sparse estimation procedure to focus more on local interactions, while still taking into account the full set of possible interactions. Future work in this application can provide a more systematic evaluation of ability of the algorithm to identify recurring patterns, by validating the estimated electrode interactions with manually annotated propagation patterns. The question of causal relationships between electrodes in terms of conduction sources and sinks can be further investigated by comparing the reconstructed recurrent patterns to the frequency-domain analysis, as proposed by Richter et al. [85, 86]. This could show its potential to identify (recurrent) drivers of AF, or elucidate underlying structures in the atrial tissue (e.g. endocardial bundle architecture or epicardial fiber orientation [60]) that contribute to electrical dissociation and conduction disturbances, which in turn may lead to conduction block, wave break, and transmural breakthroughs [30, 31]. Finally, although the method was developed for high-density contact electrograms where distance between electrodes is known and constant, it could also be applicable in endocardial basket recordings of AF to identify dominant patterns of conduction.

9.2.2 Feature selection to predict pharmacological cardioversion

A second application in the analysis of AF was presented in Chapter 8. Here a different sparse estimation technique was employed to identify the dominant predictors (ECG complexity parameters and clinical parameters) of successful cardioversion of AF using flecainide. Although the identification problem at hand was not necessarily underdetermined (the number of predictors was typically lower than the number of observa-

tions), a sparse estimation approach was called for because of the need to find a small subset of parameters that possessed a good predictive performance. The collinearity of several predictors, arising from the fact that identical parameters were computed on neighbouring lead locations, or that different parameters had a similar interpretation, motivated the use of elastic net regression in this application. A generalized linear model implementation of elastic net logistic regression was employed to select dominant features, supplemented with features (if different) selected by ordinary forward stepwise logistic regression. This application illustrated the use and advantages of a sparse estimation approach to feature selection, by being able to identify subsets of (possibly correlated) candidate predictors without the need to evaluate all possible combinations of parameters, while overcoming the shortcomings of traditional stepwise methods when it comes to variable selection in the presence of collinearity. The selected predictors indicated that individual ECG-derived AF complexity parameters are capable to detect small, but significant differences in patients with recent onset AF when it comes to successful CV. Predictive performance of ECG-parameters improved by combining different types of parameters computed on different leads. Compared to conventional clinical predictors, ECG-parameters provided better prediction of successful CV, most notably by parameters computed on a single lead in the frequency domain. Although clinical implications of this study may be limited, since pharmacological cardioversion is - at least in patients without contraindication - a low-cost and low-risk procedure, the approach to compare and select relevant predictors in a rigorous and systematic way is a contribution to the standardization of noninvasive AF complexity quantification analysis. Analyzing larger patient cohorts and, as advocated in this thesis, more robust approaches to feature selection will in the near future clarify the added value of the large set of candidate noninvasive AF complexity parameters at various stages of AF management. Analysis of the added value of high-coverage body surface potential maps to improve noninvasive AF complexity quantification will also benefit from this type of robust feature selection.

Bibliography

- [1] R. ALCARAZ AND J.J. RIETA, Adaptive singular value cancelation of ventricular activity in single-lead atrial fibrillation electrocardiograms, *Physiological Measurements* vol. 29 (12), pp. 1351–1369, 2008. [cited at p. 92, 109]
- [2] R. ALCARAZ, J.J. RIETA, Sample entropy of the main atrial wave predicts spontaneous termination of paroxysmal atrial fibrillation, *Medical Engineering & Physics* vol. 31, pp. 917–922, 2009. [cited at p. 121]
- [3] R. ALCARAZ, F. SANDBERG, L. SÖRNMO, J.J. RIETA, Classification of paroxysmal and persistent atrial fibrillation in ambulatory ECG recordings, *IEEE Transactions on Biomedical Engineering* vol. 58, pp. 1441–1449, 2011. [cited at p. 110]
- [4] M.A. ALLESSIE, N.M.S. DE GROOT, R.P.M. HOUBEN, U. SCHOTTEN, E. BOERSMA, J.L. SMEETS, AND H.J. CRIJNS, Electropathological Substrate of Long-Standing Persistent Atrial Fibrillation in Patients With Structural Heart Disease: Longitudinal Dissociation, *Circulation: Arrhythmia and Electrophysiology*, vol. 3, no. 6, pp. 606–615, Dec. 2010. [cited at p. 88, 91, 92, 94]
- [5] D. BAUER, Subspace algorithms, *Proceedings of the 13th IFAC Symposium on System Identification*, Rotterdam, The Netherlands, pp. 1030–1041, 2003. [cited at p. 50]
- [6] D. BAUER, Asymptotic Properties of Subspace Estimators, *Automatica* vol. 41 (3), Special Issue on Data-Based Modeling and System Identification, pp. 359–376, 2005. [cited at p. 50]
- [7] T. BARAN, D. WEI, A.V. OPPENHEIM, Linear Programming Algorithms for Sparse Filter Design. *IEEE Transactions on Signal Processing*. vol. 58(3):1605–17, 2010. [cited at p. 2]
- [8] R. BELLMAN AND K.J. ASTROM, On structural identifiability, *Mathematical Biosciences* vol. 7 (3-4), pp. 329–339, 1970. [cited at p. 51]
- [9] B. BESSELINK, U. TABAK, A. LUTOWSKA, N. VAN DE WOUW, H. NIJMEIJER, D.J. RIXEN, M.E. HOCHSTENBACH, W.H.A. SCHILDERS, A comparison of model reduction techniques from structural dynamics, numerical mathematics and systems and control, *Journal of Sound and Vibration*, vol. 16;332(19):4403–22, 2013. [cited at p. 1]
- [10] P. BLOOMFIELD AND W.L. STEIGER, *Least Absolute Deviations: Theory, Applications, and Algorithms*, Birkhäuser, Boston, 1983. [cited at p. 7]

- [11] P. BONIZZI, O. MESTE, V. ZARZOSO, Spectral analysis of atrial signals directly from surface ECG exploiting compressed spectrum, *Computers in Cardiology*, pp. 221–224, 2008. [cited at p. 109]
- [12] P. BONIZZI, M. DE L.S. GUILLEM, A.M. CLIMENT, J. MILLET, V. ZARZOSO, F. CASTELLS, O. MESTE, Noninvasive assessment of the complexity and stationarity of the atrial wavefront patterns during atrial fibrillation, *IEEE Transactions on Biomedical Engineering*, vol. 57, pp. 2147–2157, 2010. [cited at p. 110]
- [13] P. BONIZZI, S. ZEEMERING, J.M.H. KAREL, L.Y. DI MARCO, L. ULDRY, J. VAN ZAEN, J.-M. VESIN, U. SCHOTTEN, Systematic comparison of non-invasive measures for the assessment of atrial fibrillation complexity: a step forward towards standardization of atrial fibrillation electrogram analysis, *Europace*. vol. 17 (2), pp.318–25, 2015. [cited at p. 108]
- [14] C. BOONE, H. BUSSEY, AND B. J. ANDREWS, Exploring genetic interactions and networks with yeast, *Nature Reviews Genetics*, vol. 8, no. 6, pp. 437–449, Jun. 2007. [cited at p. 2]
- [15] A.J. CAMM, G.Y.H. LIP, R. DE CATERINA, I. SAVELIEVA, D. ATAR, S.H. HOHNLOSER, ET AL., 2012 focused update of the ESC Guidelines for the management of atrial fibrillation: an update of the 2010 ESC Guidelines for the management of atrial fibrillation. Developed with the special contribution of the European Heart Rhythm Association. *European Heart Journal*, pp. 2719–2747, 2012. [cited at p. 85, 86, 107]
- [16] A. CABASSON, O. MESTE, Time Delay Estimation: A New Insight Into the Woody's Method, *IEEE Signal Processing Letters* vol. 15, pp. 573–576, 2008. [cited at p. 109]
- [17] H. CALKINS, K.-H. KUCK, R. CAPPATO, J. BRUGADA, A.J. CAMM, S.-A. CHEN, H.J.G. CRIJNS, R.J. DAMIANO, D.W. DAVIES, J. DIMARCO, J. EDGERTON, K. ELLENBOGEN, M.D. EZEKOWITZ, D.E. HAINES, M. HAÏSSAGUERRE, G. HINDRICKS, Y. IESAKA, W. JACKMAN, J. JALIFE, P. JAÏS, J. KALMAN, D. KEANE, Y.-H. KIM, P. KIRCHHOF, G. KLEIN, H. KOTTKAMP, K. KUMAGAI, B.D. LINDSAY, M. MANSOUR, F. E. MARCHLINSKI, P.M. MCCARTHY, J.L. MONT, F. MORADY, K. NADEMANEE, H. NAKAGAWA, A. NATALE, S. NATTEL, D.L. PACKER, C. PAPPONE, E. PRYSTOWSKY, A. RAVIELE, V. REDDY, J.N. RUSKIN, R.J. SHEMIN, H.-M. TSAO, AND D. WILBER, 2012 HRS/EHRA/ECAS Expert Consensus Statement on Catheter and Surgical Ablation of Atrial Fibrillation: recommendations for patient selection, procedural techniques, patient management and follow-up, definitions, endpoints, and research trial design., *Europace*, vol. 14, no. 4, pp. 528–606, Apr. 2012. [cited at p. 86, 87]
- [18] E.J. CANDÉS AND M.B. WAKIN, An Introduction To Compressive Sampling, *IEEE Signal Processing Magazine*., vol. 25, no. 2, pp. 21–30, 2008. [cited at p. 2]
- [19] S.S. CHEN, D.L. DONOHO, AND M.A. SAUNDERS, Atomic Decomposition by Basis Pursuit, *SIAM Review*, vol 43 (1), pp. 129–159, 2001. [cited at p. 100]
- [20] F. CHEN, H. GARNIER, AND M. GILSON, Robust identification of continuous-time models with arbitrary time-delay from irregularly sampled data, *Journal of Process Control*, pp. 19–27, 2015. [cited at p. 131]
- [21] M.B. CHOUDHARY, F. HOLMQVIST, J. CARLSON, H.-J. NILSSON, A. ROIJER, P.G. PLATONOV, Low atrial fibrillatory rate is associated with spontaneous conversion of recent-onset atrial fibrillation, *Europace* vol. 15, pp. 1445–1452, 2013. [cited at p. 121, 123]

- [22] VAŠEK CHVÁTAL, *Linear Programming*, W.H. Freeman and Company, New York, 1983. [cited at p. 10, 13]
- [23] W.J. CULVER, On the existence and uniqueness of the real logarithm of a matrix, *Proceedings of the American Mathematical Society* vol. 17 (5), pp. 1146–1151, 1966. [cited at p. 69]
- [24] J. DELFORGE, On local identifiability of linear systems, *Mathematical Biosciences*, vol. 70 (1), pp. 1–37, 1984. [cited at p. 51]
- [25] J.E. DENNIS AND R.B. SCHNABEL, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Prentice-Hall, New York, 1983. [cited at p. 50]
- [26] J.-M. DION, C. COMMAULT AND J. VAN DER WOUDE, Generic properties and control of linear structured systems: a survey, *Automatica* vol. 39, pp. 1125–1144, 2003. [cited at p. 51]
- [27] D.L. DONOHO, Compressed sensing, *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006. [cited at p. 2]
- [28] D.L. DONOHO, For most large underdetermined systems of linear equations the minimal l_1 -norm solution is also the sparsest solution, *Communications on Pure and Applied Mathematics*, vol. 59, no. 6, pp. 797–829, 2006. [cited at p. 36, 129]
- [29] D. L. DONOHO, M. ELAD, On the stability of the basis pursuit in the presence of noise, *Signal Processing*, vol 86 (3), pp. 511–532, 2006. [cited at p. 101]
- [30] J. ECKSTEIN, B. MAESEN, D. LINZ, S. ZEEMERING, A. VAN HUNNIK, S. VERHEULE, M. ALLESSIE, AND U. SCHOTTEN, Time course and mechanisms of endo-epicardial electrical dissociation during atrial fibrillation in the goat, *Cardiovascular Research*, vol. 89, no. 4, pp. 816–824, Feb. 2011. [cited at p. 132]
- [31] J. ECKSTEIN, S. ZEEMERING, D. LINZ, B. MAESEN, S. VERHEULE, A. VAN HUNNIK, H. CRIJNS, M.A. ALLESSIE, AND U. SCHOTTEN, Transmural Conduction Is the Predominant Mechanism of Breakthrough During Atrial Fibrillation: Evidence From Simultaneous Endo-Epicardial High-Density Activation Mapping, *Circulation: Arrhythmia and Electrophysiology*, vol. 6, no. 2, pp. 334–341, Apr. 2013. [cited at p. 132]
- [32] B. EFRON, T. HASTIE, I. JOHNSTONE AND R. TIBSHIRANI, Least angle regression, *The Annals of statistics*, vol. 32 (2), pp. 407–499, 2004. [cited at p. 43]
- [33] R. FLETCHER, *Practical Methods of Optimization*, John Wiley and Sons Ltd., Chichester, 1987. [cited at p. 18, 50]
- [34] J. FRIEDMAN, T. HASTIE AND R. TIBSHIRANI, Regularization Paths for Generalized Linear Models via Coordinate Descent, *Journal of Statistical Software*, vol. 33, no. 1, pp. 1–22, 2010. [cited at p. 124]
- [35] J.-J. FUCHS, More on sparse representations in arbitrary bases, *Proceedings of the 13th IFAC Symposium on System Identification*, Rotterdam, The Netherlands, pp. 1357–1362, 2003. [cited at p. 31]
- [36] J.-J. FUCHS, On sparse representations in arbitrary bases, *IEEE Transactions on Information Theory* vol. 50 (6), pp. 1341–1344, 2004. [cited at p. 7, 31, 101]

- [37] J.-J. FUCHS, Recovery conditions of sparse representations in the presence of noise, *ICASSP*, pp. 337–340, 2006. [cited at p. 24]
- [38] T.S. GARDNER, Inferring Genetic Networks and Identifying Compound Mode of Action via Expression Profiling, *Science*, vol. 301, no. 5629, pp. 102–105, Jul. 2003. [cited at p. 2]
- [39] H. GARNIER, M. MENSLER, AND A. RICHARD, Continuous-time model identification from sampled data: Implementation issues and performance evaluation, *International Journal of Control*, vol. 76, no. 13, pp. 1337–57, Jan. 2003. [cited at p. 131]
- [40] W.S. GRAY AND E.I. VERRIEST, Optimality properties of balanced realizations: Minimum sensitivity, *Proceedings of the 26th IEEE Conference on Decision and Control*, Los Angeles, CA, USA, pp. 124–128, 1987. [cited at p. 18, 130]
- [41] N.M.S. DE GROOT, R.P.M. HOUBEN, J.L. SMEETS, E. BOERSMA, U. SCHOTTEN, M.J. SCHALIJ, H. CRIJNS, AND M.A. ALLESSIE, Electropathological Substrate of Longstanding Persistent Atrial Fibrillation in Patients With Structural Heart Disease: Epicardial Break-through, *Circulation*, vol. 122, no. 17, pp. 1674–1682, Oct. 2010. [cited at p. 88]
- [42] N.K. GUPTA AND R.K. MEHRA, Computational aspects of maximum likelihood estimation and reduction in sensitivity function calculations, *IEEE Transactions on Automatic Control* vol. 19, pp. 774–783, 1974. [cited at p. 18]
- [43] E.J. HANNAN AND M. DEISTLER, *The Statistical Theory of Linear Systems*, John Wiley and Sons, New York, 1988. [cited at p. 49]
- [44] B. HANZON AND R.J. OBER, Overlapping block-balanced canonical forms for various classes of linear systems, *Linear Algebra and its Applications* vol. 281, pp. 171–225, 1998. [cited at p. 50]
- [45] A. E. HOERL AND R. W. KENNARD, Ridge regression: biased estimation for nonorthogonal problems, *Technometrics*, vol. 42 (1), Special 40th Anniversary Issue, pp. 80–86, 2000. [cited at p. 22]
- [46] N. HUSSAMI AND R. TIBSHIRANI, A Component Lasso, *arXiv.org*, 2013. [cited at p. 130]
- [47] T. KAILATH, *Linear Systems*, Prentice-Hall, Englewood Cliffs, NJ, 1980. [cited at p. 50]
- [48] W.B. KANNEL, P.A. WOLF, E.J. BENJAMIN, D. LEVY, Prevalence, incidence, prognosis, and predisposing conditions for atrial fibrillation: population-based estimates. *American Journal of Cardiology* vol. 82, no. 7, 1998. [cited at p. 123]
- [49] P. L. KEMPKER, Coordination Control of Linear Systems, *PhD Thesis*, Vrije Universiteit Amsterdam, 2012. [cited at p. 130]
- [50] J.M.H. KAREL, S.A.P. HADDAD, S. HISENI, R.L. WESTRA, W.A. SERDIJN, AND R.L.M. PEETERS, Implementing Wavelets in Continuous-Time Analog Circuits With Dynamic Range Optimization, *IEEE Transactions on Circuits and Systems*, vol. 59, no. 2, pp. 229–242, 2012. [cited at p. 2]
- [51] P. KIRCHHOF, G.Y.H. LIP, I.C. VAN GELDER, J. BAX, E. HYLEK, S. KAAB, U. SCHOTTEN, Comprehensive risk reduction in patients with atrial fibrillation: emerging diagnostic and therapeutic options—a report from the 3rd Atrial Fibrillation Competence NETWORK/European Heart Rhythm Association consensus conference, *Europace*, vol. 14, no. 1, pp. 8–27, Jan. 2012. [cited at p. 91]

- [52] K.T. KONINGS, C.J. KIRCHHOF, J.R. SMEETS, H.J. WELLENS, O.C. PENN, AND M.A. ALLESSIE, High-density mapping of electrically induced atrial fibrillation in humans, *Circulation*, vol. 89, no. 4, pp. 1665–1680, Apr. 1994. [cited at p. 88]
- [53] T.A.R. LANKVELD, S. ZEEMERING, H.J.G.M. CRIJNS, U. SCHOTTEN, The ECG as a tool to determine atrial fibrillation complexity, *Heart* vol.100, pp.1077–1084, 2014. [cited at p. 107]
- [54] W.E. LARIMORE, System identification, reduced order filters and modeling via canonical variate analysis, in: H.S. Rao and P. Dorato (eds.), *Proceedings of the 1983 American Control Conference 2*, Piscataway, NJ, pp. 445–451, 1983. [cited at p. 50]
- [55] C.L. LAWSON AND R.J. HANSON, *Solving Least Squares Problems*, Prentice-Hall, Englewood Cliffs, NJ, 1974. [cited at p. 50]
- [56] H.J.LIN, P.A. WOLF, M. KELLY-HAYES, A.S. BEISER, C.S. KASE, E.J. BENJAMIN, AND R.B. D’AGOSTINO, Stroke Severity in Atrial Fibrillation: The Framingham Study, *Stroke*, vol. 27, no. 10, pp. 1760–1764, Oct. 1996. [cited at p. 2]
- [57] L. LJUNG, *MATLAB System Identification Toolbox Users Guide*, Version 6, The Mathworks, 2004. [cited at p. 49]
- [58] L. LJUNG, *System Identification: Theory for the User* (2nd ed.), Prentice-Hall Inc., Englewood Cliffs, NJ, 1999. [cited at p. 6, 48, 49, 50]
- [59] L. LJUNG AND T. GLAD, On global identifiability for arbitrary model parametrizations, *Automatica* vol. 30 (2), pp. 265–276, 1994. [cited at p. 51]
- [60] B. MAESEN, S. ZEEMERING, C. AFONSO, J. ECKSTEIN, R.A.B. BURTON, A. VAN HUNNIK, D.J. STUCKEY, D. TYLER, J. MAESSEN, V. GRAU, S. VERHEULE, P. KOHL, AND U. SCHOTTEN, Rearrangement of atrial bundle architecture and consequent changes in anisotropy of conduction constitute the 3-dimensional substrate for atrial fibrillation, *Circulation: Arrhythmia and Electrophysiology*, vol. 6, no. 5, pp. 967–975, Oct. 2013. [cited at p. 132]
- [61] L.Y. DI MARCO, J.P. BOURKE, P. LANGLEY, Spatial complexity and spectral distribution variability of atrial activity in surface ECG recordings of atrial fibrillation, *Med. Biol. Eng. Comput.* vol 50, pp. 439–446, 2012. [cited at p. 110]
- [62] T. MCKELVEY, Fully parametrized state-space models in system identification, *Proceedings of the 10th IFAC Symposium on System Identification*, Copenhagen, Denmark, vol. 2, pp. 373–378, 1994. [cited at p. 50]
- [63] T. MCKELVEY AND A. HELMERSSON, System identification using an overparametrized model class - improving the optimization algorithm, *Proceedings of the 36th IEEE Conference on Decision and Control*, San Diego, California, USA, pp. 2984–2989, 1997. [cited at p. 50]
- [64] T. MCKELVEY, A. HELMERSSON AND T. RIBARITS, Data driven local coordinates for multi-variable linear systems and their application to system identification, *Automatica* vol. 40, pp. 1629–1635, 2004. [cited at p. 50]
- [65] M. MEO, V. ZARZOSO, O. MESTE, D.G. LATCU, N. SAOUDI, Spatial variability of the 12-lead surface ECG as a tool for noninvasive prediction of catheter ablation outcome in persistent atrial fibrillation, *IEEE Transactions on Biomedical Engineering*, vol. 60, pp. 20–27, 2013. [cited at p. 110]

- [66] A.J. MILLER, Selection of subsets of regression variables, *Journal of the Royal Statistical Society, Series A*, vol. 147, part 2, 398–425, 1984. [cited at p. 22]
- [67] A.J. MILLER, *Subset selection in regression*, Chapman and Hall, 1990. [cited at p. 22]
- [68] B.C. MOORE, Principal component analysis in linear systems: Controllability, observability and model reduction, *IEEE Transactions on Automatic Control* vol. 26, pp. 17–32, 1981. [cited at p. 130]
- [69] I. NAJFELD AND T.F. HAVEL, Derivatives of the Matrix Exponential and Their Computation, *Advances in Applied Mathematics*, vol. 16 (3), pp. 321–375, 1995. [cited at p. 55]
- [70] I. NAULT, N. LELLOUCHE, S. MATSUO, ET AL., Clinical value of fibrillatory wave amplitude on surface ECG in patients with persistent atrial fibrillation, *J Interv Card Electrophysiol* vol. 26, pp. 11–19, 2009. [cited at p. 109]
- [71] M.E.J. NEWMAN, The Structure and Function of Complex Networks, *SIAM Review*, 2003. [cited at p. 131]
- [72] B. NINNESS, A. WILLS AND S. GIBSON, The University of Newcastle Identification Toolbox (UNIT), *Proceedings of the 16th IFAC World Congress*, Prague, 2005. [cited at p. 49]
- [73] R.J. OBER, Balanced realizations: canonical form, parametrization, model reduction, *International Journal of Control* vol. 46, pp. 643–670, 1987. [cited at p. 130]
- [74] P. VAN OVERSCHEE AND B. DE MOOR, *Subspace Identification for Linear Systems*, Kluwer Academic Publishers, 1996. [cited at p. 50]
- [75] J. PAN AND W.J. TOMPKINS, A Real-Time QRS Detection Algorithm, *IEEE Transactions on Biomedical Engineering*, vol. 32, no. 3, pp. 230–236, Mar. 1985. [cited at p. 92]
- [76] R.L.M. PEETERS, System identification based on Riemannian geometry: theory and algorithms, *Tinbergen Institute Research Series* 64, Thesis Publishers, Amsterdam, 1994. [cited at p. 18, 50]
- [77] R.L.M. PEETERS AND R.L. WESTRA, On the identification of sparse gene regulatory networks, *Proceedings of the 16th International Symposium on the Mathematical Theory of Networks and Systems*, Leuven, Belgium, 2004. [cited at p. 12, 61]
- [78] L. PERNEBO AND L.M. SILVERMAN, Model reduction via balanced state space representations, *IEEE Transactions on Automatic Control* vol. 27, pp. 382–387, 1982. [cited at p. 130]
- [79] P.G. PLATONOV, V.D.A. CORINO, M. SEIFERT, F. HOLMQVIST, L. SÖRNMO, Atrial fibrillatory rate in the clinical context: natural course and prediction of intervention outcome, *Europace* vol. 16 (Suppl. 4), pp. iv110–iv119, 2014. [cited at p. 122]
- [80] J. QIAN, T. HASTIE, J. FRIEDMAN, R. TIBSHIRANI, AND N. SIMON, EDs., *Glmnet for Matlab* (2013). [Online]. Available: http://www.stanford.edu/hastie/glmnet_matlab/. [Accessed: 17-Mar-2015]. [cited at p. 112, 130]
- [81] T. RIBARITS, The role of parametrizations in identification of linear dynamic systems, *PhD Thesis*, TU Wien, 2002. [cited at p. 50, 56]

- [82] T. RIBARITS, M. DEISTLER AND B. HANZON, On new parametrization methods for the estimation of linear state-space models, *International Journal of Adaptive Control and Signal Processing* vol. 18, Special Issue on Subspace-based Identification, pp.17–743, 2004. [cited at p. 56]
- [83] T. RIBARITS, M. DEISTLER AND B. HANZON, An analysis of separable least squares data driven local coordinates for maximum likelihood estimation of linear systems, *Automatica* vol. 41 (3), Special Issue on Data-Based Modeling and System Identification, pp. 531–544, 2005. [cited at p. 50]
- [84] T. RIBARITS, M. DEISTLER AND T. MCKELVEY, An analysis of the parametrization by data driven local coordinates for multivariable linear systems, *Automatica* vol. 40 (5), pp. 789–803, 2004. [cited at p. 50]
- [85] U. RICHTER, L. FAES, A. CRISTOFORETTI, M. MASÈ, F. RAVELLI, M. STRIDH, AND L. SÖRNMO, A novel approach to propagation pattern analysis in intracardiac atrial fibrillation signals., *Annals of Biomedical Engineering*, vol. 39, no. 1, pp. 310–323, Jan. 2011. [cited at p. 132]
- [86] U. RICHTER, L. FAES, F. RAVELLI, AND L. SÖRNMO, Propagation Pattern Analysis During Atrial Fibrillation Based on Sparse Modeling, *IEEE Transactions on Biomedical Engineering*, vol. 59 (5), pp. 1319–1328, 2012. [cited at p. 99, 132]
- [87] U. SCHOTTEN, S. VERHEULE, P. KIRCHHOF, AND A. GOETTE, Pathophysiological Mechanisms of Atrial Fibrillation: A Translational Appraisal, *Physiological Reviews*, vol. 91, no. 1, pp. 265–325, Jan. 2011. [cited at p. 91]
- [88] U. SCHOTTEN, B. MAESEN, S. ZEEMERING, The need for standardization of time- and frequency-domain analysis of body surface electrocardiograms for assessment of the atrial fibrillation substrate, *Europace*, vol. 4 (8), pp. 1072–5, 2012. [cited at p. 108]
- [89] N. SIMON, J. FRIEDMAN, T. HASTIE, AND R. TIBSHIRANI, Regularization paths for Cox's proportional hazards model via coordinate descent, *Journal of statistical software*, vol. 39, no. 5, pp. 1–13, 2011. [cited at p. 112, 130]
- [90] N. SIMON, J. FRIEDMAN, T. HASTIE, AND R. TIBSHIRANI, A Sparse-Group Lasso, *Journal of Computational and Graphical Statistics*, vol. 22, no. 2, pp. 231–245, Apr. 2013. [cited at p. 130]
- [91] J. SJÖBERG, T. MCKELVEY AND L. LJUNG, On the use of regularization in system identification, *Proceedings of the 12th IFAC World Congress*, Sydney, Australia, vol. 7, pp. 381–386, 1993. [cited at p. 50]
- [92] T.S. SÖDERSTRÖM AND P. STOICA, *System Identification*, Prentice-Hall, New York, 1989. [cited at p. 49, 50, 101]
- [93] P. STOICA, Y. SELEN. Model-order selection: a review of information criterion rules, *Signal Processing Magazine, IEEE*, vol. 21 (4), pp. 36–47, 2004. [cited at p. 1]
- [94] D.S. STOFFER, D.E. TYLER, D.A. WENDT, The spectral envelope and its applications, *Statistical Science* pp. 224–253, 2000. [cited at p. 109]
- [95] R. TIBSHIRANI, Regression shrinkage and selection via the lasso, *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 58 (1), pp. 267–288, 1996. [cited at p. 22, 24]

- [96] L. ULDRY, J. VAN ZAEN, Y. PRUDAT, L. KAPPENBERGER, J-M. VESIN, Measures of spatiotemporal organization differentiate persistent from long-standing atrial fibrillation, *Europace*, vol. 14, pp. 1125–1131, 2012. [cited at p. 109, 114, 122]
- [97] C. VAN LOAN, Computing integrals involving the matrix exponential, *IEEE Transactions on Automatic Control* vol. 23 (3), pp. 395–404, 1978. [cited at p. 55]
- [98] M. VERHAEGEN, Identification of the deterministic part of MIMO state space models given in innovations form from input-output data, *Automatica* vol. 30, pp. 61–74, 1994. [cited at p. 50]
- [99] S. VERHEULE, E. TUYLS, A. VAN HUNNIK, M. KUIPER, U. SCHOTTEN, AND M. ALLESSIE, Fibrillatory Conduction in the Atrial Free Walls of Goats in Persistent and Permanent Atrial Fibrillation, *Circulation: Arrhythmia and Electrophysiology*, vol. 3 (6), pp. 590–599, 2010. [cited at p. 88, 103]
- [100] E.I. VERRIEST AND W.S. GRAY, A geometric approach to the minimum sensitivity design problem, *SIAM Journal on Control and Optimization* vol. 33 (3), pp. 863–881, 1995. [cited at p. 18, 130]
- [101] S. WALTER AND H. TIEMEIER, Variable selection: current practice in epidemiological studies, *European journal of epidemiology*, vol. 24, no. 12, pp. 733–736, 2009. [cited at p. 2]
- [102] A. WILLS, B. NINNESS AND S. GIBSON, On Gradient-Based Search for Multivariable System Estimates, *Proceedings of the 16th IFAC World Congress*, Prague, 2005. [cited at p. 17, 18, 51]
- [103] J.H. WILKINSON, The perfidious polynomial, *Studies in numerical analysis*, vol. 24, pp. 1–28, 1984. [cited at p. 63]
- [104] M. YUAN AND Y. LIN, Model selection and estimation in regression with grouped variables, *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, vol. 68, no. 1, pp. 49–67, 2006. [cited at p. 130]
- [105] H. ZOU AND T. HASTIE, Regularization and variable selection via the Elastic Net, *Journal of the Royal Statistical Society, Series B*, vol. 67, pp. 301–320, 2005 [cited at p. 22, 110]

Summary

In many practical settings it is desirable to describe the essence of a system by a model in a simplified, yet accurate way. Model complexity reduction can be achieved in many ways, for instance by determining if the value of a model parameter is significantly different from zero or not. In this thesis the notion of sparsity is employed as an underlying property of a system that is to be estimated, in the presence of insufficient amounts of available measurement data. Sparsity is defined here as the (relative) number of nonzero parameters within a system. Part I of this work is dedicated to the development of a sparse system identification framework and applications in two model classes, linear regression models and state-space models. In Chapter 2 a general framework for sparse identification of linear time-invariant (LTI) models is proposed that constitutes a hybrid approach to sparse parameter estimation. Minimizing a least-squares (ℓ_2) prediction error criterion ensures the quality of the overall fit of the model to the data. Minimizing the absolute value (ℓ_1 -norm) of the parameter vector within the model equivalence space aims to maximize parameter vector sparsity. This mixed ℓ_2/ℓ_1 optimization framework is subsequently applied to the class of linear regression models (Chapter 3) and the class of state-space models (Chapter 4), in an underdetermined setting where the number of parameters to be estimated is typically (much) larger than the available measurements. In the setting of linear regression models, a sparse solution can be found in the space of models with an equivalent minimal prediction error criterion value (Section 3.2). Experiments indicate that in a noiseless setting, the minimal amount of input-output (i/o) data needed to correctly estimate the data generating model parameters is much lower on average than stated by a theoretical (worst case) lower limit (Section 3.3).

An iterative version of the mixed ℓ_2/ℓ_1 optimization algorithm is applied to the general class of state-space models in innovations form (Section 4.5). Here, model equivalence is intrinsically present by the existence of a (non-singular) state-space transformation matrix that leads to equivalent i/o behaviour. After initial minimization of the prediction error criterion, the ℓ_1 -norm of the parameter vector is minimized in the space formed by the linear approximation of this equivalence space. The optimal solution is then taken as a direction in which a step is made to improve parameter vector sparsity. This ℓ_2/ℓ_1 minimization procedure is repeated to traverse along the manifold of increasingly sparse i/o equivalent state-space models. Experiments show

that this iterative version of mixed ℓ_2/ℓ_1 optimization is capable of increasing model sparsity while retaining an optimal fit to the data. The choice of the size of the step in the sparse search direction is however critical for the convergent properties of this algorithm. In a fully parameterized setting a connection with data driven local coordinates (DDL) can be made, achieving fast convergence to a sparse solution.

The special subclass of state-space models describing (sparse) network interactions is studied in more detail in Sections 4.6-4.8. These networks have a ring-like interaction structure, which makes large-scale generation of sparse, stable and minimal networks more feasible than for instance in the case of a random interaction structure. In sparse networks with interactions in discrete time (Section 4.6), the results obtained from the linear regression setting in Chapter 3 are directly transferrable, since a sparse linear regression problem can be solved for each state separately. Moving to networks with a sparse interaction matrix structure in continuous time introduces the complicating factor of the sampling period that determines the level of sparsity present in the discrete-time representation of the interaction matrix and the duration of the observed time interval. Estimating a sparse discrete-time interaction matrix, followed by transformation to continuous time is possible, but only for a narrow range of sampling times (Section 4.7). The continuous-time interactions can be estimated directly using the iterative ℓ_2/ℓ_1 optimization algorithm and some experiments show that it is indeed feasible to find a correct sparse solution in an underdetermined setting (Section 4.8), but the applicability of this technique is limited in this case due to the tendency of the algorithm to converge to a solution that corresponds to an unstable state-space model.

The notion of sparsity is investigated in the field of atrial fibrillation analysis. Atrial fibrillation (AF) is a common arrhythmia in which the normal, synchronized contraction of the atria is disturbed and multiple waves of electrical activations propagate over the atria. Invasive, high-density measurements of AF allow for a detailed description of the number and behaviour of fibrillation waves. In Chapter 6 an automated probabilistic approach to AF electrogram annotation and fibrillation wave construction is described that is shown to correspond well to manual pattern assessment in several recordings of human AF. An alternative approach to propagation pattern analysis is developed that is based on a sparse multivariate autoregressive model of electrogram interactions (Chapter 7), extending the work in Chapters 3 and 4 to time-delayed sparse linear regression. This technique enables identification of recurring propagation patterns in fibrillation waves, without manual or automated annotation. Example analyses in high-density mapping in a goat model of AF illustrate the ability of this application of sparse estimation to capture recurring patterns like wave trains, breakthrough waves and rotating wave fronts. Another application of sparse estimation is in the noninvasive analysis of AF (Chapter 8). Identification of the dominant predictors of successful pharmacological cardioversion (i.e. restoration of regular sinus rhythm) of AF is achieved using elastic net logistic regression, a shrinkage estimator employing a mixed regularization of the prediction error criterion by the ℓ_2 - and ℓ_1 -norm of the predictor coefficients. Here, sparse estimation is employed to select features (ECG-based complexity parameters and clinical features) in an overdetermined setting.

Samenvatting

In de praktijk is het vaak wenselijk om de essentie van een systeem te beschrijven met behulp van een model op een vereenvoudigde, maar toepasselijke manier. Het verminderen van de complexiteit van een model kan op aantal verschillende manieren, bijvoorbeeld door te bepalen of de waarde van een modelparameter significant verschilt van nul of niet. In dit proefschrift wordt het begrip sparsity (dungetheid of ijlheid zijn equivalente begrippen in de Nederlandse taal) gebruikt als een onderliggende eigenschap van een te schatten systeem, waarbij het aantal beschikbare waarnemingen onvoldoende is om de modelparameters uniek te kunnen schatten. Sparsity wordt hier gedefiniëerd als het (relatieve) aantal niet-nul parameters in een systeem. Deel I van dit proefschrift is gewijd aan de ontwikkeling van een raamwerk voor systeemidentificatie van dergelijke sparse systemen en de toepassing van dit raamwerk in een tweetal modelklassen, namelijk de klasse van lineaire regressiemodellen en de klasse van toestandsruimtemodellen. In Hoofdstuk 2 wordt dit algemene raamwerk voor sparse systeemidentificatie van lineaire, tijdsinvariante modellen geïntroduceerd, bestaande uit een hybride aanpak van sparse parameterschatting. Het minimaliseren van een kleinste kwadraten (ℓ_2) predictie-errorcriterium zorgt voor een goede overeenkomst tussen het model en de beschikbare meetgegevens. Het minimaliseren van de absolute waarde (ℓ_1 -norm) van de parameter vector in de ruimte van gelijkwaardige modellen, heeft als doel de sparsity van de parameter vector te maximaliseren. Dit gecombineerde ℓ_2/ℓ_1 optimalisatieraamwerk wordt vervolgens toegepast in de klasse van lineaire regressiemodellen (Hoofdstuk 3) en de klasse van toestandsruimtemodellen (Hoofdstuk 4), in beide gevallen in een onderbepaalde situatie waarin het aantal parameters dat geschat dient te worden gewoonlijk (veel) groter is dan het beschikbare aantal waarnemingen. In het geval van lineaire regressie modellen kan een sparse oplossing worden gevonden in de ruimte van modellen met een gelijke minimale waarde van het predictie-errorcriterium (Sectie 3.2). Experimenten laten zien dat in een situatie zonder ruis, de minimale hoeveelheid input-output (i/o) gegevens die nodig is om de parameters die de gegevens hebben gegenereerd, correct terug te schatten, gemiddeld veel lager is dan afgaande op een theoretische (worst case) ondergrens (Sectie 3.3).

Een iteratieve versie van het gecombineerde ℓ_2/ℓ_1 optimalisatiealgoritme wordt toegepast op de klasse van toestandsruimtemodellen in vernieuwingsvorm (Sectie 4.5).

Modelequivalentie is hier ingebakken in de beschrijving van het model, door het bestaan van een (niet-singuliere) toestandsruimte-transformatiematrix die leidt tot equivalent *i/o* gedrag. Na initiële minimalisatie van het predictie-errorcriterium wordt vervolgens de ℓ_1 -norm van de parameter vector geminimaliseerd in de ruimte beschreven door een lineaire benadering van deze equivalentieruimte. De optimale oplossing wordt dan gebruikt als een richting waarin een stap wordt genomen om de sparsity van de parameter vector te verbeteren. Deze ℓ_2/ℓ_1 minimalisatieprocedure wordt herhaald om een pad te volgen over de topologische ruimte van steeds sparsere *i/o*-equivalente toestandsruimtemodellen. Uit experimenten blijkt dat deze iteratieve versie van gecombineerde ℓ_2/ℓ_1 optimalisatie in staat is om de sparsity van een model te vergroten, terwijl de overeenkomst tussen het model en de waarnemingen gewaarborgd blijft. De grootte van de stap in de zoekrichting naar sparsity is echter cruciaal voor de convergentie van het algoritme. In een situatie waarin het toestandsruimte-model volledig geparameteriseerd is, kan een link worden gelegd naar de techniek *data driven local coordinates* (DDLC), waardoor een snellere convergentie naar een sparse oplossing kan worden bereikt.

De subklasse van toestandsruimtemodellen die een (sparse) interactienetwerk beschrijven, wordt verder onderzocht in Secties 4.6-4.8. De onderzochte netwerken bezitten een ringvormige interactiestructuur die het beter mogelijk maakt om op grote(re) schaal sparse, stabiele en minimale netwerken te creëren, vergeleken met bijvoorbeeld netwerken met een willekeurige interactiestructuur. In het geval van sparse netwerken met interacties in discrete tijd (Sectie 4.6) zijn de resultaten verkregen bij lineaire regressiemodellen direct overdraagbaar, aangezien een sparse lineair regressieprobleem kan worden opgelost voor elke afzonderlijke toestand. Het schatten van netwerken met een sparse interactiematrix in continue tijd brengt de complicerende factor met zich mee van de wijze van bemonstering. Het bemonsteringsinterval bepaalt namelijk zowel de mate van sparsity van de interactiematrix in discrete tijd, als de totale duur van het interval dat kan worden waargenomen. Het schatten van een sparse interactiematrix in discrete tijd, gevolgd door een transformatie naar continue tijd, is mogelijk, maar alleen voor een beperkt bereik van bemonsteringsintervallen (Sectie 4.7). De interacties in continue tijd kunnen direct worden geschat met behulp van het iteratieve ℓ_2/ℓ_1 optimalisatiealgoritme, en enkele experimenten bevestigen ook dat het mogelijk is om de correcte sparse oplossing te bepalen in een onderbepaalde situatie (Sectie 4.8), maar de praktische toepassing van deze techniek is beperkt in dit geval vanwege de neiging van het algoritme om naar een oplossing te convergeren die overeenkomt met een instabiel toestandsruimtemodel.

Het begrip sparsity wordt verder onderzocht in de analyse van atriumfibrilleren. Atriumfibrilleren (AF) is een veelvoorkomende hartritmestoornis waarbij de normale, gelijktijdige samentrekking van de atria is verstoord en in plaats daarvan meerdere golven de atria activeren. Invasieve metingen van AF met een hoge dichtheid van elektrodes maken een gedetailleerde beschrijving mogelijk van het aantal fibrillatiegolven en hun gedrag. In Hoofdstuk 6 wordt een geautomatiseerde, probabilistische aanpak van de annotatie van electrogrammen en de constructie van fibrillatiegolven beschreven, die goed blijkt overeen te komen met handmatige analyse van golfpatronen in ver-

scheidene patiënten. Een alternatieve manier om golfpatronen te analyseren is ontwikkeld, die gebaseerd is op een sparse multivariaat autoregressiemodel van interacties tussen electrogrammen (Hoofdstuk 7), waarbij de technieken uit Hoofdstuk 3 en 4 verder worden uitgebreid naar lineaire regressiemodellen met verschoven versies van de gegevens in de tijd. Deze techniek maakt het mogelijk om terugkerende patronen te identificeren in fibrillatiegolven, zonder handmatige of geautomatiseerde annotatie van electrogrammen. Voorbeeldanalyses van opnames in een geitenmodel van AF laten zien dat de toepassing van een sparse schattingsmethode het mogelijk maakt terugkerende patronen te onderscheiden, zoals golftreinen, doorbrekende golven en ronddraaiende golffronten. Het idee van een sparse schattingsmethode heeft ook een toepassing in de niet-invasieve analyse van AF (Hoofdstuk 8). De identificatie van belangrijke voorspellers van succesvolle farmacologische cardioversie (het herstellen van het normale patroon van sinusritme) van AF kan worden bereikt door elastic net logistische regressie toe te passen. De elastic net techniek is een shrinkage estimator die het predictie-errorcriterium regulariseert met een gecombineerde ℓ_2 - en ℓ_1 -norm van de coëfficiënten van de voorspellers. Hier wordt de sparse schattingsmethode ingezet om de belangrijkste kenmerken (complexiteitsparameters bepaald uit een ECG en klinische kenmerken) te selecteren in een overbepaalde situatie.

Valorization

Introduction

The notion of sparsity, as discussed in this thesis, is nowadays ubiquitous in science, for instance in feature selection in high-dimensional datasets (big data), or complex interaction network identification. The growing number of potential features that can be related to a certain trait or outcome is often not matched by the number of observations that can be recorded. This makes it more difficult to quantify the role of each feature in a unique and unambiguous way. Assuming sparsity in the relationship between all the candidate features or interactions and the observed behavior can overcome this limitation to some extent.

Sparse estimation

Sparse estimation of parameters in a linear regression problem (Chapter 3), given a limited set of available measurements, greatly improves the chances of correctly estimating the data generating parameter values, if the data generating parameter vector is indeed sparse. This is a technique that is applicable in many fields of research, in principle in every research question where the parameter estimation problem is underdetermined, i.e. when the number of parameters (greatly) exceeds the number of available measurements, or when one aims to determine the dominant regressors in a regression model. The application in this thesis of sparse (logistic) regression to determine the dominant predictors of successful pharmacological cardioversion, given a large set of candidate predictors, illustrates this application directly.

When considering state-space models (Chapter 4) there is an inherent model equivalence present (in the fully parameterized setting) that can be exploited to search for a maximally sparse solution within the set of models that exhibit equivalent input-output behaviour. In models where the interactions between states represent physical connections, finding a sparse solution can help to lower model complexity and reduce for instance the number of connections to be manufactured in a digital chip or the number of interactions between a large number of genes to take into account.

The case of sparse network interaction models (as a subclass of state-space models) has been studied in more detail in this thesis. In the case of a limited amount of available measurements, sparse estimation of network model parameters can also aid

in identifying the correct data generating parameters, using an adapted form of sparse linear regression (in discrete time) or an iterative version of the sparsity maximization algorithm (in continuous time). An application of sparse estimation of network interaction models in discrete time is presented in this thesis in Chapter 7. One limiting factor in this approach to sparse network interaction estimation, in discrete, but most notably in continuous time, is that appropriate sampling of the network inputs and outputs is critical to finding a meaningful solution. This issue needs to be investigated further, to determine necessary and sufficient conditions under which sparse estimation of network interactions is a viable approach.

Applications in atrial fibrillation

The developed software package to automatically process and annotate atrial electrograms (Chapter 6) is now used within our Physiology department, enabling fast and reproducible analysis of longer recordings of high-density contact mapping. This has substantially changed the way mapping data is being collected in our recent studies, moving from manually analyzing 4-second segments, which could take days, to analyzing several minutes of electrogram data recorded simultaneously at several atrial sites, in only a few minutes.

The algorithm developed for identification of recurring wave front propagation patterns (Chapter 7) has a potential application during (surgical) ablation of atrial fibrillation to identify regions with highly recurrent propagation patterns that are associated with the maintenance of AF. The readout of the algorithm is a directed graph of the mapping area that shows the dominant interactions between atrial locations, also indicating the strength of the interaction. This information will show a cardiologist an immediate impression of the prevailing conduction pattern in that area, which may be instrumental in guiding the ablation process.

The demonstrated predictive value of noninvasive AF complexity parameters (Chapter 8), when it comes to predicting successful pharmacological cardioversion of paroxysmal AF, has a clear potential to guide AF treatment. Using the software package that has been developed to extract the atrial signal from body surface ECGs and to compute complexity parameters in a standardized way, an initial noninvasive assessment of AF complexity can be made when a patient visits the clinic. Properly trained and validated prediction models can be developed and integrated in a knowledge support system that assists a physician in making an informed decision on patient treatment. This will eventually lead to a more patient-specific treatment, improving quality of care and reducing costs by abandoning likely unsuccessful treatment options.

Dankwoord / Acknowledgements

ἄνδρα μοι ἔννεπε, Μοῦσα, πολύτροπον, ὃς μάλα πολλὰ πλάγχθη,...
— Ὅμηρος, Ἰδύσσεια

Na vele jaren is er nu dan toch een eind gekomen aan een project dat ik in juni 2004 startte. Indertijd maakte ik vanuit mijn werk bij MaTeUM de overstap, letterlijk een trap omhoog in hetzelfde gebouw, naar een plek als promovendus bij de vakgroep Wiskunde, faculteit der Algemene Wetenschappen aan de Universiteit Maastricht. Officieel aangesteld op het project "Breedtestrategie", wat in de praktijk net zo vaag bleek als de naam doet vermoeden. Na aanvankelijk volstrekt geen idee te hebben gehad wat een promotieonderzoek zou moeten inhouden, lukte het uiteindelijk om een rudimentair promotieplan op te stellen. Sparse estimation moest het worden, toegepast op netwerken, met de beperkende factor dat er onvoldoende data beschikbaar was om een unieke oplossing te schatten. Er vanuit gaan dat veel van de modelparameters in feite niet meedoen, een nulwaarde hebben, kan in zo'n geval het schattingsprobleem vereenvoudigen. Helaas bleken vier jaren onvoldoende om het onderzoek volledig af te ronden. Tijdens mijn werk als software engineer bij Maastricht Instruments kwam ik in de jaren daaropvolgend veelvuldig in contact met de vakgroep Fysiologie, waar ik in januari 2011 aan de slag kon als onderzoeker op het gebied van signaalanalyse met betrekking tot atriumfibrilleren. In de loop van de tijd bleken er enkele raakvlakken te zijn tussen het onderwerp van het promotieonderzoek en de problemen die zich voordoen bij het analyseren van atriumfibrilleren. Het proefschrift dat er nu is, is daarom ook een combinatie van een licht theoretische verhandeling over systeem identificatie en praktische toepassingen op het gebied van atriumfibrilleren.

Can't bring back time. Like holding water in your hand.
— James Joyce, *Ulysses*

Het spreekt voor zich dat er in deze periode veel mensen zijn geweest die op een of andere manier bij mijn promotieonderzoek betrokken waren. Ik zal proberen om zoveel mogelijk van hen te bedanken, al vrees ik dat ik minstens een paar mensen zal vergeten.

Laat ik beginnen met Joris van de Klundert, toenmalig directeur van MaTeUM, en Koos Vrieze, ook betrokken bij MaTeUM, maar toentertijd natuurlijk voornamelijk hoogleraar bij de vakgroep Wiskunde. Joris, ik herinner me de tijd bij MaTeUM, toen nog een klein bedrijfje in een kamer bij de faculteit Algemene Wetenschappen, als een heel plezierige en inspirerende ervaring, waar altijd veel ruimte was voor mijn eigen inbreng en creativiteit. Ik denk ook dat hier mijn interesse voor een wetenschappelijke carrière is gewekt. Dat blijkt ook wel uit het feit dat een van mijn eerste wetenschappelijke publicaties uit een samenwerking met jou is voortgekomen. Koos, mijn promotor van het eerste uur, echter niet van het laatste, ik wil jou bedanken voor de kans die je me hebt gegeven om het eens als promovendus te proberen, en voor je bemoedigende woorden als we elkaar weer eens tegen kwamen na je pensionering. Dat brengt mij op mijn huidige promotor, Ralf Peeters. Ralf, ik zou je voor veel kunnen bedanken. Laat ik me beperken tot mijn oprechte dank voor je langdurige begeleiding, je nooit aflatende stroom van ideeën, oplossingen en achtergrondinformatie op zeer uiteenlopende gebieden binnen de wiskunde, en met name de systeemtheorie. Daarnaast wil ik ook mijn copromotor Ronald Westra bedanken voor zijn steun, zijn ongebreideld enthousiasme en zeker ook voor zijn hulp bij mijn eerste buitenlandse congres (in Kyoto nog wel). Ook mijn kamergenoten in de periode bij Wiskunde wil ik niet onvermeld laten. Joël, Ivo en Martijn, onze tijd samen op zolder, en later bij de bushalte (toen aangevuld met Jordi en Stefan), was meestal rustig en toch gezellig, met slechts af en toe een enkele wespenplaag.

Van mijn tijd bij Maastricht Instruments heb ik veel geleerd, en ik wil Frans Smeets, toenmalig directeur, bedanken voor de mogelijkheden die hij mij daar heeft geboden en de medewerking die ik heb gekregen toen ik besloot naar Fysiologie te vertrekken. Ook de mensen binnen het software team, Iwan, Ralf, Vincent, René en Bart, ben ik dankbaar voor de ervaring van het werken binnen een hecht en betrokken team.

“... it's a fool that looks for logic in the chambers of the human heart.”

— *O Brother, Where Art Thou*, dir. Joel Coen

Dat brengt mij tot Ulrich Schotten, hoogleraar bij de vakgroep Fysiologie en mijn tweede promotor. Uli, ook jou wil ik voor veel bedanken, maar misschien toch nog wel het meest voor je persoonlijke steun en motiverende gesprekken, niet alleen op het gebied van atriumfibrilleren. Zonder jouw kordate aanpak was dit proefschrift er waarschijnlijk nog lang niet geweest. Daarnaast heb ik de intensieve samenwerking met veel collega's binnen de groep als zeer verrijkend ervaren. Jens, ik kan me van jou vooral de lange sessies samen met Uli herinneren, waarin we tot in detail allerlei analysemethodes de revue lieten passeren, en je (geveinsde?) verbazing als er weer eens iets nieuws in de door mij ontwikkelde analysesoftware opdook. Bart, ook jou wil ik bedanken, niet alleen voor het aanleveren van mapping data van hoge kwaliteit, maar eveneens voor alle wetenschappelijke en vriendschappelijke discussies. Hetzelfde geldt zeker voor Arne, mijn huidige kamergenoot, alhoewel onze directe samenwerking wat trager verloopt, aangezien dit proefschrift toch een keer af moest. Theo, jouw inbreng

is uiteindelijk zeer belangrijk geweest om het verhaal mooi af te ronden, met de non-invasieve analyse van atriumfibrilleren, en ik hoop dat onze samenwerking nog veel moois zal opleveren in de toekomst.

I should not forget to mention my room mates from the last five years: Ali, even though our working hours were mostly out of phase, I enjoyed your company and I wish you all the best in Lugano. Pawel, our collaboration was also impaired by my work on my thesis, but I know we will collaborate again more efficiently in the future. Dennis, I enjoyed working with you, drinking wine with you and trying to keep a straight face during your constant innuendo.

I would like to extend a special thanks to my liaison at the department of Knowledge Engineering and one of my paranymphs. Pietro, the last couple of years we have been working more and more together on cross-departmental research in atrial fibrillation and this has proved to be more and more fruitful. In addition to that I appreciate your pleasant and warm personality, and I have fond memories of our yearly recurring CinC-trip.

Mijn andere paranimf, Jeroen. Als we de verhalen moeten geloven, zaten wij al als peuters op onze driewielers uitgebreid te discussiëren, in plaats van rond te fietsen. En het klopt, gedurende onze lange vriendschap zijn we een goede discussie inderdaad nooit uit de weg gegaan, maar gelukkig kunnen we in het algemeen toch best aardig met elkaar overweg. Fijn dat je de rol van paranimf graag op je wilt nemen, ik had niet anders verwacht.

Het is lastig om de precieze rol van familie in dit proefschrift te duiden, maar vast staat dat ik veel heb te danken aan mijn ouders, Paul en Herma. Jullie hebben me geleerd om altijd nieuwsgierig te zijn en me gesteund in de keuzes die ik heb gemaakt, en hebben slechts een enkele keer voorzichtig geïnformeerd hoe het er nou eigenlijk voorstond met mijn promotie. Nu kunnen we het eindelijk vieren.

Lieve Boukje, in de tijd dat ik met mijn promotie bezig ben geweest, is er heel wat veranderd in ons leven samen. Zo hebben we inmiddels twee geweldige dochters die ons ook goed bezig weten te houden. Soms verdween het onderwerp promotie zelfs zodanig naar de achtergrond, dat je je afvroeg of ik het ooit zou afronden. Dank je voor je begrip, al die keren dat je het even alleen moest zien te rooien met de kinderen, wanneer ik fysiek of geestelijk afwezig was. Nu deze horde genomen is, hoop ik dat we, nog meer dan voorheen, samen van het leven kunnen genieten en alles wat het te bieden heeft.

List of publications

- [1] S. ZEEMERING, T.A.R. LANKVELD, P. BONIZZI, H.J.G.M. CRIJNS, U. SCHOTTEN, Systematic comparison of techniques for noninvasive assessment of atrial fibrillation complexity to predict outcome of pharmacological cardioversion in patients with recent onset atrial fibrillation. *To be submitted.*
- [2] S. ZEEMERING, T.A.R. LANKVELD, H.J.G.M. CRIJNS, U. SCHOTTEN, Towards standardization of noninvasive atrial fibrillation substrate complexity quantification: Effect of choice of ECG-leads and complexity measure on prediction of pharmacological cardioversion. *Proceedings of the Computing in Cardiology Conference (CinC)*, pp. 879-882, 2013.
- [3] S. ZEEMERING, R.L.M. PEETERS, A. VAN HUNNIK, S. VERHEULE, U. SCHOTTEN, Identification of recurring wavefront propagation patterns in atrial fibrillation using basis pursuit. *Proceedings of the Engineering in Medicine and Biology Society (EMBC), 35th Annual International Conference of the IEEE*, pp. 2928-2931, 2013.
- [4] S. ZEEMERING, B. MAESEN, J. NIJS, D.H. LAU, M. GRANIER, S. VERHEULE, U. SCHOTTEN, Automated quantification of atrial fibrillation complexity by probabilistic electrogram analysis and fibrillation wave reconstruction. *Proceedings of the Engineering in Medicine and Biology Society (EMBC), 34th Annual International Conference of the IEEE*, pp. 6357-6360, 2012.
- [5] U. SCHOTTEN, B. MAESEN, S. ZEEMERING, The need for standardization of time- and frequency-domain analysis of body surface electrocardiograms for assessment of the atrial fibrillation substrate, *Europace*, vol. 4 (8), pp. 1072–5, 2012.
- [6] R.L.M. PEETERS, S. ZEEMERING, Sparse gene regulatory network identification. *Knowledge Discovery and Emergent Complexity in Bioinformatics*, pp. 171-182, 2007.
- [7] R.L.M. PEETERS, I.T.J. TOSSINGS, S. ZEEMERING, Sparse system identification by mixed ℓ_2/ℓ_1 -minimization. *Proceedings of the 17th International Symposium on Mathematical Theory of Networks And Systems*, p. 1466, 2006.
- [8] S. ZEEMERING, R.L.M. PEETERS, I.T.J. TOSSINGS, A prediction error method for sparse system identification, *Proceedings of the 25th Benelux Meeting on Systems and Control*, p. 21, 2006.

Curriculum vitae

Stef Zeemering was born on March 7, 1979 in Oldenzaal, the Netherlands. He studied Knowledge Engineering at Maastricht University from 1997 until 2002, completing a master in operations research with a thesis titled Gaussian mixture decomposition: clustering via maximum likelihood and the EM algorithm. In March 2002 he started working at MaTeUM b.v. (Mathematische Technologie Universiteit Maastricht) as a mathematical consultant. Then, in June 2004, he became involved in a Ph.D. project at the Department of Mathematics, (now part of the Department of Knowledge Engineering) at Maastricht University, on sparse estimation of (biological) model parameters in the presence of model equivalence. From June 2008 until the end of 2010 he worked at Maastricht Instruments b.v. as a software engineer and project leader. In the beginning of 2011 he joined the Department of Physiology as a researcher on atrial fibrillation complexity analysis. Part of his work there is presented in Part II of this thesis as an application of the theoretical framework developed in Part I. After completion of his Ph.D. thesis, Stef will continue working at the Department of Physiology as a post-doctoral fellow. Stef currently lives in Maastricht, the Netherlands, together with his partner Boukje and their two daughters, Elin en Tinde.

Acronyms

- AF** atrial fibrillation. 85
AFR atrial fibrillatory rate. 121
AUC area under the curve. 110
- BMI** body mass index. 117
- CV** cardioversion. 85, 108
- DDLC** data-driven local coordinates. 48
DF dominant frequency. 109
- ECG** electrocardiogram. 87
- FWA** fibrillation wave amplitude. 109
- GN** Gauss-Newton. 18
- i.i.d.** independent and identically distributed. 48
- LAD** left atrial diameter. 117
LP linear programming. 9
LTI linear time-invariant. 47, 48
LVEF left ventricular ejection fraction. 119
LVESD left ventricular end systolic diameter. 119
- MAW** main atrial wave. 110
MDF multidimensional dominant frequency. 109
MFWA multidimensional fibrillation wave amplitude. 109
MOI multidimensional organization index. 109
MSE multidimensional spectral entropy. 109
- OI** organization index. 109
- PCA** principal component analysis. 110
PEM prediction error method. 49
- RAV** right atrial volume. 117
RHE relative sub-band energy. 110
ROC receiver operating characteristics. 110
- SAE** sample entropy. 110
SE spectral entropy. 109
SNR signal to noise ratio. 37
SR sinus rhythm. 85
SVD singular value decomposition. 11
ZOH zero order hold. 65

Notation index

- β bound on the step in the sparse search direction s . 19
- \cong equivalence relation. 8
- $e(\theta)$ error vector. 5
- $H(\theta)$ Hessian matrix of the least squares criterion $V(\theta)$. 8
- $J(\theta)$ Jacobian matrix of the error vector $e(\theta)$. 8
- Ker** matrix kernel. 8
- \otimes Kronecker product. 56
- ℓ_p p-norm of a vector. 2
- \mathcal{M} model class. 5
- $M_{\min}^{(N,k)}$ minimum M to estimate a parameter vector of length N with k non-zero entries. 35
- \perp orthogonal complement. 17
- P_0 probability of a correct parameter vector estimate. 26
- P_0^{LS} probability of a correct parameter vector estimate, (least-squares solution). 32
- Φ regression matrix of regression vectors ϕ . 23
- ϕ regression vector. 22
- Φ^+ Moore-Penrose pseudo inverse. 23
- s local search direction for sparsity. 8
- S_c number of incorrectly estimated parameter values. 25
- θ model parameter vector. 5
- θ_{LS} least squares solution. 23
- θ_{MN} minimum norm solution. 23
- $\theta_{\min}^{\neq 0}$ minimum estimated non-zero parameter value. 39
- θ_{\max}^0 maximum estimated zero parameter value. 39
- $V(\theta)$ least squares criterion. 5
- Var** vector variance. 37
- vec**(\cdot) vectorization operator. 52

Index

- active set method, 11, 38
- AF cycle length distribution, 92
- atrial fibrillation, 85

- backward elimination, 22
- basis pursuit, 101
- bilinear transform, 55, 69
- breakthrough wave, 104

- cardioversion
 - electrical, 85
 - pharmacological, 85, 108
- conduction tortuosity, 96
- conduction velocity, 94
- convergence, 19
- Cox proportional hazards, 111

- data-driven local coordinates, 50
- dead time, 99
- decision variable, 10
- deflection detection, 92
- dependent variable, 21
- directional derivative, 55
- dominant frequency, 109
 - multidimensional, 109
- duality theorem, 13

- elastic net, 22, 25
- elastic net logistic regression, 110
- electrograms, 86
- error vector, 5

- f-wave power, 110
- far-field deflection, 92
- fibrillation wave amplitude, 109
 - multidimensional, 109
- fibrillation wave construction, 94
- forward selection, 22

- Gauss-Newton method, 18
- general identification process, 6

- Hessian matrix, 8, 53, 75
- high-density mapping, 91

- interaction matrix, 47
- interior point method, 11, 38
- intrinsic deflection, 92

- Jacobian matrix, 8, 52, 75

- Kalman gain matrix, 48
- kernel, 8
- Kronecker product, 56

- lasso, 22, 24
- least angle regression, 43
- least squares criterion, 5
 - linear regression, 23
 - nonlinear, 49
- Levenberg-Marquardt method, 18
- linear interaction network
 - continuous time, 65
 - discrete time, 60
- linear programming, 9
 - standard form, 10
- linear regression, 21
- logistic regression, 110
- long-standing AF, 85

- matrix coherence, 31

- maximum likelihood estimation, 50
- minimality, 48
- minimum norm solution, 23
- mixed optimization, 17
- model equivalence, 5
- model order, 1
- Moore-Penrose pseudo inverse, 23
- MVAR model, 100

- organization index, 109
- orthogonal complement, 17
- orthogonal projection, 60
- orthonormal matrix, 11
- overdetermined, 23

- paroxysmal AF, 85
- peripheral wave, 103
- permanent AF, 85
- persistent AF, 85
- prediction error method, 49
- prediction error process, 49
- principal component analysis, 110

- QRST-cancellation, 92

- rate control, 85
- recent onset AF, 85
- regressor, 21
- relative sub-band energy, 110
- rhythm control, 85
- ridge regression, 22
- rotating wave, 104

- sample entropy, 110
- search direction, 18
- Shannon's entropy, 109
- signal to noise ratio, 37
- simplex algorithm, 11
- simplex method, 38
- singular value decomposition, 11, 92
- sinus rhythm, 85
- sparsity, 1
 - partial, 12, 60
- spatial complexity, 110
- spatiotemporal stationarity, 110
- spectral concentration, 110
- spectral entropy, 109
 - multidimensional, 109
- spectral envelope, 109
- spectral organization index
 - multidimensional, 109
- spectral variability, 110
- state-space model
 - continuous time, 54
 - discrete time, 48
 - innovations form, 48
- subset selection, 21

- Taylor series, 8
- Tikhonov regularization, *see* ridge regression
- time-delay, 99
- tridiagonal matrix, 61

- underdetermined, 22
- unidentifiability, 50