
KERNEL ADAPTIVE FILTERING

A Comprehensive Introduction

Weifeng Liu, José C. Príncipe, and
Simon Haykin

 **WILEY**

JOHN WILEY & SONS, INC., PUBLICATION

KERNEL ADAPTIVE FILTERING

**Adaptive and Learning Systems for Signal Processing,
Communication, and Control**

Editor: Simon Haykin

A complete list of titles in this series appears at the end of this volume.

KERNEL ADAPTIVE FILTERING

A Comprehensive Introduction

Weifeng Liu, José C. Príncipe, and
Simon Haykin

 **WILEY**

JOHN WILEY & SONS, INC., PUBLICATION

Copyright © 2010 by John Wiley & Sons, Inc. All rights reserved

Published by John Wiley & Sons, Inc., Hoboken, New Jersey
Published simultaneously in Canada

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 750-4470, or on the web at www.copyright.com. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permission>.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services or for technical support, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic formats. For more information about Wiley products, visit our web site at www.wiley.com.

Library of Congress Cataloging-in-Publication Data:

Liu, Weifeng

Kernel adaptive filtering : a comprehensive introduction / Jose C. Principe, Weifeng Liu, Simon Haykin.

p. cm.

Includes bibliographical references and index.

ISBN 978-0-470-44753-6

1. Adaptive filters. 2. Kernel functions. I. Príncipe, José C. II. Haykin, Simon S., 1931– III. Title.

TK7872.F5P745 2010

621.382'23–dc22

2009042654

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

To our families

CONTENTS

PREFACE	xi
ACKNOWLEDGMENTS	xv
NOTATION	xvii
ABBREVIATIONS AND SYMBOLS	xix
1 BACKGROUND AND PREVIEW	1
1.1 Supervised, Sequential, and Active Learning / 1	
1.2 Linear Adaptive Filters / 3	
1.3 Nonlinear Adaptive Filters / 10	
1.4 Reproducing Kernel Hilbert Spaces / 12	
1.5 Kernel Adaptive Filters / 16	
1.6 Summarizing Remarks / 20	
Endnotes / 21	
2 KERNEL LEAST-MEAN-SQUARE ALGORITHM	27
2.1 Least-Mean-Square Algorithm / 28	
2.2 Kernel Least-Mean-Square Algorithm / 31	
2.3 Kernel and Parameter Selection / 34	
2.4 Step-Size Parameter / 37	
2.5 Novelty Criterion / 38	
2.6 Self-Regularization Property of KLMS / 40	
2.7 Leaky Kernel Least-Mean-Square Algorithm / 48	
2.8 Normalized Kernel Least-Mean-Square Algorithm / 48	

- 2.9 Kernel ADALINE / 49
- 2.10 Resource Allocating Networks / 53
- 2.11 Computer Experiments / 55
- 2.12 Conclusion / 63
- Endnotes / 65

3 KERNEL AFFINE PROJECTION ALGORITHMS 69

- 3.1 Affine Projection Algorithms / 70
- 3.2 Kernel Affine Projection Algorithms / 72
- 3.3 Error Reusing / 77
- 3.4 Sliding Window Gram Matrix Inversion / 78
- 3.5 Taxonomy for Related Algorithms / 78
- 3.6 Computer Experiments / 80
- 3.7 Conclusion / 89
- Endnotes / 91

4 KERNEL RECURSIVE LEAST-SQUARES ALGORITHM 94

- 4.1 Recursive Least-Squares Algorithm / 94
- 4.2 Exponentially Weighted Recursive Least-Squares Algorithm / 97
- 4.3 Kernel Recursive Least-Squares Algorithm / 98
- 4.4 Approximate Linear Dependency / 102
- 4.5 Exponentially Weighted Kernel Recursive Least-Squares Algorithm / 103
- 4.6 Gaussian Processes for Linear Regression / 105
- 4.7 Gaussian Processes for Nonlinear Regression / 108
- 4.8 Bayesian Model Selection / 111
- 4.9 Computer Experiments / 114
- 4.10 Conclusion / 119
- Endnotes / 120

5	EXTENDED KERNEL RECURSIVE LEAST-SQUARES ALGORITHM	124
5.1	Extended Recursive Least Squares Algorithm / 125	
5.2	Exponentially Weighted Extended Recursive Least Squares Algorithm / 128	
5.3	Extended Kernel Recursive Least Squares Algorithm / 129	
5.4	EX-KRLS for Tracking Models / 131	
5.5	EX-KRLS with Finite Rank Assumption / 137	
5.6	Computer Experiments / 141	
5.7	Conclusion / 150	
	Endnotes / 151	
6	DESIGNING SPARSE KERNEL ADAPTIVE FILTERS	152
6.1	Definition of Surprise / 152	
6.2	A Review of Gaussian Process Regression / 154	
6.3	Computing Surprise / 156	
6.4	Kernel Recursive Least Squares with Surprise Criterion / 159	
6.5	Kernel Least Mean Square with Surprise Criterion / 160	
6.6	Kernel Affine Projection Algorithms with Surprise Criterion / 161	
6.7	Computer Experiments / 162	
6.8	Conclusion / 173	
	Endnotes / 174	
	EPILOGUE	175
	APPENDIX	177
A	MATHEMATICAL BACKGROUND	177
A.1	Singular Value Decomposition / 177	
A.2	Positive-Definite Matrix / 179	
A.3	Eigenvalue Decomposition / 179	
A.4	Schur Complement / 181	

x CONTENTS

- A.5 Block Matrix Inverse / 181
- A.6 Matrix Inversion Lemma / 182
- A.7 Joint, Marginal, and Conditional Probability / 182
- A.8 Normal Distribution / 183
- A.9 Gradient Descent / 184
- A.10 Newton's Method / 184

B APPROXIMATE LINEAR DEPENDENCY AND SYSTEM STABILITY	186
REFERENCES	193
INDEX	204

PREFACE

For the first time, this book presents a comprehensive and unifying introduction to kernel adaptive filtering. Adaptive signal processing theory has been built on three pillars: the linear model, the mean square cost, and the adaptive least-square learning algorithm. When nonlinear models are required, the simplicity of linear adaptive filters evaporates and a designer has to deal with function approximation, neural networks, local minima, regularization, and so on. Is this the only way to go beyond the linear solution? Perhaps there is an alternative, which is the focus of this book. The basic concept is to perform adaptive filtering in a linear space that is related nonlinearly to the original input space. If this is possible, then all three pillars and our intuition about linear models can still be of use, and we end up implementing nonlinear filters in the input space.

This book will draw on the theory of reproducing kernel Hilbert spaces (RKHS) to implement the nonlinear transformation of the input to a high-dimensional feature space induced by a positive-definite function called *reproducing kernel*. If the filtering and adaptation operations to be performed in RKHS can be expressed by inner products of projected samples, then they can be directly calculated by kernel evaluations in the input space. We use this approach to introduce a family of adaptive filtering algorithms in RKHS:

- The kernel least-mean-square algorithm
- The kernel affine projection algorithms
- The kernel recursive least-squares algorithm
- The extended kernel recursive least-squares algorithm

These kernel-learning algorithms bridge closely two important areas of adaptive filtering and neural networks, and they embody beautifully two important methodologies of error-correction learning and memory-based learning. The bottlenecks of the RKHS approach to nonlinear filter design are the need for regularization, the need to select the kernel function, and the need to curtail the growth of the filter structure. This book will present in a mathematically rigorous manner the issues and the solutions to all these

problems, and it will illustrate with examples the performance gains of kernel adaptive filtering.

Chapter 1 starts with an introduction to general concepts in machine learning, linear adaptive filters, and conventional nonlinear methods. Then, the theory of reproducing kernel Hilbert spaces is presented as the mathematical foundation of kernel adaptive filters. We stress that kernel adaptive filters are universal function approximators, have no local minima during adaptation, and require reasonable computational resources.

Chapter 2 studies the kernel least-mean-square algorithm, which is the simplest among the family of kernel adaptive filters. We develop the algorithm in a step-by-step manner and delve into all the practical aspects of selecting the kernel function, picking the step-size parameter, sparsification, and regularization. Two computer experiments, one with Mackey–Glass chaotic time-series prediction and the other with nonlinear channel equalization, are presented.

Chapter 3 covers the kernel affine projection algorithms, which is a family of four similar algorithms. The mathematical equations of filtering and adaptation are thoroughly derived from first principles, and useful implementation techniques are discussed fully. Many well-known methods can be derived as special cases of the kernel affine projection algorithms. Three detailed applications are included to show their wide applicability and design flexibility.

Chapter 4 presents the kernel recursive least-squares algorithm and the theory of Gaussian process regression. A sparsification approach called approximate linear dependency is discussed. And with the aid of the Bayesian interpretation, we also present a powerful model selection method called “maximum marginal likelihood”. Two computer experiments are conducted to study the performance of different sparsification schemes and the effectiveness of maximum marginal likelihood to determine the kernel parameters.

Chapter 5 discusses the extended kernel recursive least-squares algorithm on the basis of the kernel recursive least-squares algorithm. We study systematically the problem of estimating the state of a linear dynamic system in RKHS from a sequence of noisy observations. Several important theorems are presented with proofs to outline the significance and basic approaches. This chapter contains two examples, Rayleigh channel tracking and Lorenz time-series modeling.

Chapter 6 is devoted to addressing the principal bottleneck of kernel adaptive filters, i.e., their growing structure. We introduce a subjective information measure called *surprise* and present a unifying sparsification scheme to curtail the growth effectively of kernel adaptive filters. Three interesting computer simulations are presented to illustrate the theories.

This book should appeal to engineers, computer scientists, and graduate students who are interested in adaptive filtering, neural networks, and kernel methods. A total of 12 computer-oriented experiments are distributed throughout the book that have been designed to reinforce the concepts discussed in the chapters. The computer experiments are listed in Table 1. Their MATLAB®

Table 1. A listing of all computer experiments in the book. MATLAB® programs that generate the results can be downloaded by all readers from the book's website <http://www.cnel.ufl.edu/~weifeng/publication.htm>.

Computer experiment	Topic
2.1	KLMS Applied to Mackey–Glass Time-Series Prediction
2.2	KLMS Applied to Nonlinear Channel Equalization
3.1	KAPA Applied to Mackey–Glass Time-Series Prediction
3.2	KAPA Applied to Noise Cancellation
3.3	KAPA Applied to Nonlinear Channel Equalization
4.1	KRLS Applied to Mackey–Glass Time-Series Prediction
4.2	Model Selection by Maximum Marginal Likelihood
5.1	EX-KRLS Applied to Rayleigh Channel Tracking
5.2	EX-KRLS Applied to Lorenz Time-Series Prediction
6.1	Surprise Criterion Applied to Nonlinear Regression
6.2	Surprise Criterion Applied to Mackey–Glass Time-Series Prediction
6.3	Surprise Criterion Applied to CO ₂ Concentration Forecasting

implementations can be downloaded directly from the website <http://www.cnel.ufl.edu/~weifeng/publication.htm>. To keep the codes readable, we placed simplicity over performance during design and implementation. These programs are provided without any additional guarantees.

We have strived to reflect fully the latest advances of this emerging area in the book. Each chapter concludes with a summary of the state of the art and potential future directions for research. This book should be a useful guide to both those who look for nonlinear adaptive filtering methodologies to solve practical problems and those who seek inspiring research ideas in related areas.

ACKNOWLEDGMENTS

We would like to start by thanking Dr. Puskal P. Pokharel, Seagate Technology; Dr. Murali Rao, University of Florida; Dr. Jay Gopalakrishnan, University of Florida; and Il Park, University of Florida, for their help in the development of the kernel adaptive filtering theory. We are most grateful to Dr. Aysegul Gunduz, Albany Medical College, Wadsworth Center; Dr. John Harris, University of Florida; and Dr. Steven Van Vaerenbergh, University of Cantabria, Spain, for providing many useful comments and constructive feedback on an early version of the manuscript of the book.

Many others have been kind enough to read critically selected chapters/sections of the book; in alphabetical order, they are as follows:

Erion Hasanbelliu, University of Florida, Gainesville, Florida
Dr. Kyu-hwa Jeong, Intel Corporation, Santa Clara, California
Dr. Ruijiang Li, University of California, San Diego, California
Dr. Antonio Paiva, University of Utah, Salt Lake City, Utah
Alexander Singh, University of Florida, Gainesville, Florida
Dr. Yiwen Wang, Hong Kong University of Science and Technology, Hong Kong
Dr. Jianwu Xu, University of Chicago, Chicago, Illinois

We also wish to thank (in alphabetical order): Dr. Peter Bartlett, University of California, Berkeley; Dr. Andrzej Cichocki, Riken, Brain Science Institute, Japan; Dr. Corinna Cortes, Google Lab; Dr. Graham C. Goodwin, University of Newcastle, UK; Dr. Michael Jordan, University of California, Berkeley; Dr. Thomas Kailath, Stanford University; Dr. Joel S. Kvitky, Rand Corporation; Dr. Yann LeCun, New York University; Dr. Derong Liu, University of Illinois at Chicago; Dr. David J. C. MacKay, University of Cambridge, UK; Dr. Tomaso Poggio, Massachusetts Institute of Technology; Dr. Ali H. Sayed, University of California, Los Angeles; Dr. Bernhard Schölkopf, Max Planck Institute for Biological Cybernetics, Germany; Dr. Sergios Theodoridis, University of Athens, Greece; Dr. Yoram Singer, Google Lab; Dr. Alexander J. Smola, Yahoo! research; Dr. Johan Suykens, Katholieke Universiteit

Leuven, Belgium; Dr. Paul Werbos, The National Science Foundation; Dr. Bernard Widrow, Stanford University; and Dr. Chris Williams, University of Edinburgh, UK.

We thank the staff at Wiley, publisher George Telecki, editorial assistant Lucy Hitz, and production editor Kris Parrish, as well as the project manager, Stephanie Sakson from Toppan Best-Set Premedia, for their full support and encouragement in preparing the manuscript of the book and for all their behind-the-scenes effort in the selection of the book cover and the production of the book.

Last, but by no means least, we are grateful to Lola Brooks, McMaster University, for typing many sections of the manuscript.

NOTATION

The book discusses many algorithms involving various mathematical equations. A convenient and uniform notation is a necessity to convey clearly the basic ideas of the kernel adaptive filtering theory. We think it is helpful to summarize and explain at the beginning of the text our notational guidelines for ease of reference.

There are mainly *three* types of variables we need to distinguish:

scalar, vector, and matrix variables

The following is a list of the notational conventions used in the book:

1. We use *small italic* letters to denote *scalar variables*. For example, the output of a filter is a scalar variable, which is denoted by y .
2. We use *CAPITAL ITALIC* letters to denote *SCALAR CONSTANTS*. For example, the order of a filter is a scalar constant, which is denoted by L .
3. We use **small bold** letters for **vectors**.
4. We use **CAPITAL BOLD** letters to denote **MATRICES**.
5. We use *parentheses* to denote the *time dependency* of any variables (either scalar, vector, or matrix). For example, $d(i)$ means the value of a scalar d at time (or iteration) i . $\mathbf{u}(i)$ means the value of a vector \mathbf{u} at time (or iteration) i . Similarly $\mathbf{G}(i)$ means the value of a matrix \mathbf{G} at time (or iteration) i . There is no rule without an exception. f_i is used to denote the estimate of an input–output mapping f at time (or iteration) i since parenthesis is preserved for input argument like $f_i(\mathbf{u})$.
6. We use the superscript T to denote *transposition*. For example, if

$$\mathbf{d} = \begin{bmatrix} d(1) \\ d(2) \\ \dots \\ d(N) \end{bmatrix}$$

then

$$\mathbf{d}^T = [d(1), d(2), \dots, d(N)].$$

7. All variables in our presentation are *real*. We do not discuss complex numbers in this book.
8. All vectors in our presentation are *column vectors* without exception.
9. We use subscript indices to denote 1) a component of a vector (or a matrix), 2) a general vector that the index is not related to time (or iteration). For example, \mathbf{c}_i could mean the i th vector in some set or the i th component of the vector \mathbf{c} according to the context.

We have made every effort to make the notation consistent and coherent for the benefit of the reader. The following Table 2 summarizes and lists some typical examples.

Table 2. Notation.

	Description	Examples
Scalars	Small <i>italic</i> letters	d
Vectors	Small bold letters	$\mathbf{w}, \boldsymbol{\omega}, \mathbf{c}_i$
Matrices	Capital BOLD letters	$\mathbf{U}, \boldsymbol{\Phi}$
Time or iteration	Indices in parentheses	$\mathbf{u}(i), d(i)$
Component of vectors/matrices	Subscript indices	$\mathbf{a}_j, \mathbf{G}_{i,j}$
Linear spaces	Capital mathbb letters	\mathbb{F}, \mathbb{H}
Scalar constants	Capital <i>ITALIC</i> letters	L, N

ABBREVIATIONS AND SYMBOLS

We collect here a list of the main abbreviations and symbols used throughout the text for ease of reference.

$(\cdot)^T$	vector or matrix transposition
\mathbf{A}^{-1}	inverse of matrix \mathbf{A}
$\mathbf{E}[\cdot]$	expected value of a random variable
$m(\cdot)$	the mean of a random variance
$\sigma^2(\cdot)$	the variance of a random variable
$\langle \cdot, \cdot \rangle$	inner product
$\ \cdot\ $	norm of a vector; square root of the inner product with itself
$ \cdot $	absolute value of a real number or determinant of a matrix
\propto	proportional to
\sim	distributed according to
∇	gradient
$\mathbf{0}$	zero vector or matrix
β	forgetting factor
$\mathcal{C}(i)$	dictionary or center set at iteration i
$d(i)$	desired output at time or iteration i (a real scalar)
$\text{diag}\{a, b\}$	a diagonal matrix with diagonal entries a and b
δ_1	distance threshold in novelty criterion
δ_2	prediction error threshold in novelty criterion
δ_3	threshold in approximate linear dependency test
δ_{ij}	Kronecker delta
$\Delta\mathbf{w}(i)$	weight adjustment at time or iteration i (a column vector in an Euclidean space)
\mathcal{D}	data set
$e(i)$	output estimation error at time or iteration i
\mathbb{F}	feature space induced by the kernel mapping
\mathbf{G}	Gram matrix of (transformed) input data
\mathbb{H}	reproducing kernel Hilbert space

I	identity matrix
$J(i)$	error cost at time or iteration i
$\mathbf{k}(i)$	Kalman gain (or gain vector) at time or iteration i
$K(\mathbf{A})$	condition number of a matrix \mathbf{A}
$\kappa(\mathbf{u}, \mathbf{u}')$	kernel (or covariance) function evaluated at \mathbf{u} and \mathbf{u}'
L	dimensionality of the input space
λ	regularization parameter
M	dimensionality of the feature space
\mathcal{M}	misadjustment of the least-mean-square algorithm
$\mathbf{n}(i)$	additive noise in the state space at time or iteration i
N	number of training data
η	step-size parameter
$O(\cdot)$	of the order of a number
P	state-error correlation matrix
$\varphi(\cdot)$	a mapping induced by a reproducing kernel
$\varphi(i)$	transformed filter input at time or iteration i (a column vector in a feature space)
Φ	transformed input data matrix
R	covariance matrix of (transformed) input data
\mathbb{R}	the set of real numbers
\mathbb{R}^L	L -dimensional real Euclidean space
ζ_{\max}	the maximum eigenvalue
$\text{tr}(\mathbf{A})$	trace of matrix \mathbf{A}
T_1	abnormality threshold in surprise criterion
T_2	redundancy threshold in surprise criterion
$\mathbf{u}(i)$	filter input at time or iteration i (a column vector in an Euclidean space)
\mathbb{U}	input domain
U	input data matrix
$v(i)$	additive noise in the output at time or iteration i
$\mathbf{w}(i)$	weight estimate at time or iteration i (a column vector in an Euclidean space)
$\boldsymbol{\omega}(i)$	weight estimate at time or iteration i (a column vector in a feature space)
z^{-1}	unit delay operator
AIC	Akaike information criterion
ALD	approximate linear dependency
APA	affine projection algorithm
BIC	Bayesian information criterion
CC	coherence criterion
CV	cross-validation
ENC	enhanced novelty criterion
EX-RLS	extended recursive least squares algorithm
EX-KRLS	extended kernel recursive least squares algorithm

GPR	Gaussian process regression
LMS	least-mean-square algorithm
LOOCV	leave-one-out cross-validation
LS	least squares
MAP	maximum a posterior
MDL	minimum description length
MSE	mean square error
MML	maximum marginal likelihood
NC	novelty criterion
NLMS	normalized least-mean-square algorithm
KA	kernel ADALINE
KAPA	kernel affine projection algorithm
KLMS	kernel least-mean-square algorithm
KRLS	kernel recursive least-squares algorithm
PCA	principal components analysis
PDF	probability density function
RAN	resource allocating network
RBF	radial-basis function
RKHS	reproducing kernel Hilbert space
RLS	recursive least-squares algorithm
RN	regularization network
RNN	recurrent neural network
SC	surprise criterion
SNR	signal-to-noise ratio
SVD	singular value decomposition
SVM	support vector machine
SW-KRLS	sliding window kernel recursive least-squares algorithm

