

A Study of Online Nonlinear Estimation Using Reproducing Kernels:  
 $\ell_p$ -regularized Least Squares and Kernel Weight Design

February 2021

Kwangjin Jeong

A Thesis for the Degree of Ph.D. in Engineering

A Study of Online Nonlinear Estimation Using Reproducing Kernels:  
 $\ell_p$ -regularized Least Squares and Kernel Weight Design

February 2021

Graduate School of Science and Technology  
Keio University

Kwangjin Jeong

# Acknowledgments

This thesis completes my studies in the doctoral course of Keio University since April 2018.

I would like to express my deepest gratitude to my advisor Associate Professor Masahiro Yukawa for his significant support and guidance to my entire study in the graduate school.

I would like also to express my deep gratitude to Professor Masaaki Ikehara of Keio university for his support in both aspects of my study and daily life, especially in my first year in the graduate school.

Furthermore, it is my pleasure to thank to all the current and former students of Yukawa laboratory for the friendship. In particular, I feel thankful deeply for having meaningful time including profound discussions with my roommates Mr. Masaaki Takizawa and Ms. Ran Iwamoto.

My special thanks go to Ms. Ishii for her careful support on my overall activities related to my research.

I would like also deeply thank to Professors Masaaki Ikehara and Yukihiro Sanada and Associate Professor Kei Kobayashi for serving as members of the examining committee of this thesis.

My works related to this thesis were supported by JSPS Grants-in-Aid (JP18J22032, JP18H01446).

# Contents

<b>1</b>	<b>General Introduction</b>	<b>7</b>
1.1	Sparse Optimization . . . . .	7
1.2	Kernel Adaptive Filtering . . . . .	10
1.3	This Study . . . . .	11
<b>2</b>	<b>Preliminaries</b>	<b>13</b>
2.1	Notations . . . . .	13
2.2	$\ell_p$ -regularized Least Squares Problem . . . . .	13
2.3	MKAF with Weighted Kernels . . . . .	15
2.3.1	Multikernel Adaptive Filtering . . . . .	15
2.3.2	Relation to MKL . . . . .	16
2.3.3	Imposing Kernel Weights . . . . .	17
<b>3</b>	<b>Critical-Point Paths</b>	<b>19</b>
3.1	Introduction . . . . .	19
3.2	A Path to the OLS Solution . . . . .	20
3.2.1	Existence of a Continuous Critical-point Path . . . . .	20
3.2.2	Characteristics of Critical-point Paths . . . . .	24
3.3	Paths to Sparsest OLS Solutions . . . . .	25
3.4	Paths from other OLS Solutions . . . . .	28
3.5	Concluding Remarks . . . . .	29
<b>4</b>	<b>Kernel Weight Design for MKAF</b>	<b>30</b>
4.1	Introduction . . . . .	30
4.2	Transient Behavior and Weight Derivation . . . . .	31
4.2.1	Preliminaries to Transient Analysis . . . . .	31
4.2.2	Derivation of Kernel Weights Equalizing Convergence Rates	33
4.2.3	Steady-state Error . . . . .	36
4.3	Implementation of the Kernel Weights . . . . .	37
4.3.1	Weight Estimation from Observations . . . . .	37
4.3.2	Implementation Based on Metric Design . . . . .	38
4.3.3	Use of Linear Kernel . . . . .	38
4.3.4	Calculation of Kernel Weights (4.22) under Gaussian Input	39
4.4	Numerical Experiments . . . . .	41
4.4.1	Experiments with Toy Data . . . . .	42
4.4.2	Experiments with Real Data - System Estimation . . . . .	45
4.4.3	Experiments with Real Data - Time Series Prediction . . . . .	48
4.4.4	Discussions on the Numerical Results . . . . .	54

<i>CONTENTS</i>	3
4.5 Concluding Remarks . . . . .	55
<b>5 General Conclusion</b>	<b>56</b>
<b>Bibliobraphy</b>	<b>64</b>
<b>Publications Related to the Dissertation</b>	<b>65</b>
<b>Other Publications</b>	<b>66</b>

# List of Figures

1.1	The solutions of ridge regression and LASSO. . . . .	8
1.2	The $\ell_0$ -norm, the $\ell_1$ -norm, and the $\ell_p$ -norm ( $0 < p < 1$ ). . . . .	9
1.3	The solution path and the critical-point path under an $\ell_p$ -regularization. . . . .	10
1.4	Effect of kernel weights on multikernel adaptive filtering. . . . .	12
2.1	A geometric interpretation of critical point. . . . .	14
2.2	Change of geometry by kernel weights. . . . .	18
3.1	Critical-point paths in a two-dimensional overdetermined case. . . . .	21
3.2	A case that Assumption 3.1 is violated but Theorem 3.1 still holds. . . . .	22
3.3	An illustration of Lemma 3.1(b). . . . .	23
3.4	The cases denied in Lemma 3.1 and Assumption 3.1. . . . .	24
3.5	An illustration of Remark 3.2. . . . .	25
3.6	Paths of critical points connect the origin and those OLS solutions satisfying (C2). . . . .	27
4.1	Results of Experiment 1a. . . . .	43
4.2	Weight evolution with Gaussian input(Experiment 1b.) . . . . .	44
4.3	Silverbox dataset for Experiment 2a. . . . .	45
4.4	Results of Experiment 2a. . . . .	46
4.5	Results of Experiment 2b. . . . .	47
4.6	Results of Experiment 3a. . . . .	49
4.7	Effect of the number of kernels ( $\max_q \sigma^{(q)} = 2$ ). . . . .	50
4.8	Effect of the number of kernels ( $\min_q \sigma^{(q)} = 2^{-4}$ ). . . . .	51
4.9	Results of Experiment 3b. . . . .	52
4.10	Results of Experiment 3c. . . . .	53

# List of Tables

4.1	Computational costs for Gaussian kernels . . . . .	42
4.2	Experiment 1a - Initial settings . . . . .	44
4.3	Experiment 2a - Initial settings . . . . .	45
4.4	Experiment 2b - Initial settings . . . . .	47
4.5	Experiment 3a - Initial settings . . . . .	48
4.6	Experiment 3b - Initial settings . . . . .	52
4.7	Experiment 3c - Initial settings . . . . .	53

# Abstract

Reproducing kernels have been extensively exploited for online nonlinear estimation tasks such as time-series data prediction which are of significant importance in many fields of science and engineering. However, two major challenges for online learning with kernels have also been recognized. First, the size of the dictionary grows linearly with the amount of observed data. Second, it is practically impossible to choose an appropriate kernel in an online situation where data are observed sequentially. Multikernel adaptive filtering has been proposed to solve the second issue, although a proper kernel weight design is essential to exploit its full potential.

In this thesis, the insights related to the above two issues are presented. The thesis consists of the following five chapters. Chapter 1 introduces the background and motivation of the thesis. Chapter 2 provides mathematical preliminaries on sparse optimization and online learning with kernels. The related works are introduced as well. In Chapter 3, sparse optimization by  $\ell_p$ -norm ( $0 < p < 1$ ) regularization is considered as a fundamental study for sparsification of the dictionary. Given an  $\ell_p$ -regularized least squares problem, the behavior of the critical points with a variation of the regularization parameter is investigated. In addition, it is proved that there exists a continuous, piecewise smooth path of critical points connecting the origin and the sparsest least square solution. Chapter 4 proposes a kernel weight design for multikernel adaptive filtering and shows that the proposed design is optimal in a certain sense. The proposed weight design is derived by equalizing the convergence speeds of kernel-wise coefficient errors. The derived weights can be calculated recursively using the observed data and the recursive kernel weighting is implemented to a multikernel adaptive filtering algorithm using metric projection. Numerical results show the superiority of the proposed algorithm over both of the carefully tuned preset weights and the machine learning algorithm based on multiple kernel learning, in the senses of the accuracy and the dictionary size. Chapter 5 summarizes the previous chapters and gives concluding remarks of this thesis as well as the future prospects for this study.



# Chapter 1

## General Introduction

The study of my Ph. D. course has been devoted for development of an online learning algorithm exploiting reproducing kernels, with a appropriate technique for dictionary sparsification. This thesis presents (i) a theoretical study on nonconvex regularization for sparse optimization and (ii) a multikernel adaptive filtering algorithm with an automatic tuning of kernel weights. This chapter briefly introduces the background of sparse optimization and kernel adaptive filtering.

### 1.1 Sparse Optimization

In signal processing, the way how to transform a signal, also called as *encoding-decoding* or *analysis-synthesis*, has always been a fundamental topic of study. Mallat and Zhang [1] have focused on the use of a large dictionary which consists of atoms with various waveforms including Fourier and wavelets. As an approach for flexible signal decomposition, a greedy algorithm called the matching pursuit (MP) has been proposed in [1]. Donoho and Johnstone [2] have proposed a wavelet selection technique for reconstruction of a noise-corrupted signal. The soft thresholding operator, which has been proposed in [2], has been also used for signal denoising in [3]. Both studies of [1] and [2] have consequently led to a signal representation by a small number of atoms from a redundant dictionary. Such tasks can be generalized to finding a *sparse* solution of an underdetermined least squares problem, which has more variables than equations, and hence there exist infinitely many solutions. From a practical viewpoint, sparse decomposition of a signal leads to small computational costs for processing tasks of the decomposed signal. Therefore, finding a sparse representation eventually improves efficiency of signal processing even though using a redundant dictionary may be memory-inefficient.

The  $\ell_1$ -norm is the absolute sum of the coefficients associated with the selected atoms. Chen *et al.* [4] have proposed linear programming approaches which solve a least squares problem regularized by the  $\ell_1$ -norm to obtain a sparse solution of signal decomposition. It has also been shown in [4] that the use of soft thresholding corresponds to solving a particular case of an  $\ell_1$ -regularized least squares problem. In the statistics community, a similar concept has been considered in an overdetermined system with a unique ordinary

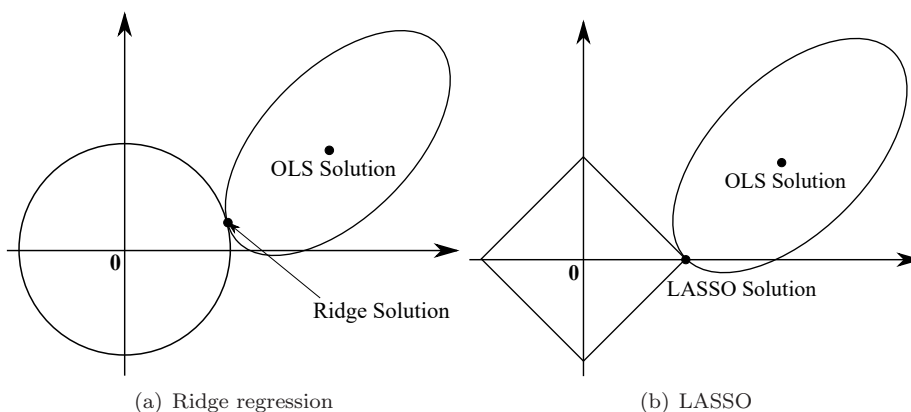


Figure 1.1: The solutions of ridge regression and LASSO.

least squares (OLS) solution. Tibshirani [5] has proposed a variable selection technique named least absolute shrinkage and selection operator (LASSO), to achieve a fine bias-variance tradeoff in regression. It is also shown that LASSO coincides with solving a least squares problem with an  $\ell_1$ -norm constraint, and LASSO achieves a sparser model than the ridge regression, which is constrained by the squared  $\ell_2$ -norm. See Figure 1.1 for comparison between solutions of ridge regression and LASSO.

Sparsity of a vector is measured by the  $\ell_0$ -norm, which counts the number of its nonzero entries. Since the  $\ell_0$ -norm is discontinuous, sparse optimization based on the  $\ell_0$ -norm yields a combinatorial optimization problem, which is NP-hard. The use of greedy algorithms to deal with the problem in affordable computational costs. The orthogonal matching pursuit (OMP) [6], an improved version of the MP, has been a celebrated example of such algorithms, since its potential ability to find a sparse solution has been reported in [7]. The OMP sets the initial coefficient vector at the origin, and then updates its step-by-step solution such that a residual error becomes zero at the solution. Thanks to simplicity and theoretical performance guarantee of the OMP, many improvements have been proposed. See [8, 9] for the ameliorated algorithms.

The  $\ell_1$ -norm can be regarded as a convex and continuous relaxation of an  $\ell_0$ -norm, which is nonconvex and discontinuous. See Figure 1.2 to compare the  $\ell_0$ -norm and the  $\ell_1$ -norm. Thanks to the convexity of the  $\ell_1$ -norm, an  $\ell_1$ -norm minimization problem with a zero-error constraint can be solved by linear programming. In addition, an  $\ell_1$ -norm minimization problem has been compared with an  $\ell_0$ -norm minimization problem profoundly. Donoho and Huo [10] have studied conditions for both problems to have unique solution. Candès and Tao [11] have proposed the restricted isometry property (RIP), which provides sufficient conditions that a solution of an  $\ell_1$ -norm minimization problem is the same with that of an  $\ell_0$ -norm minimization problem. It has been also shown that some random matrices satisfies a good RIP and this leads to the use of sparse optimization with using the  $\ell_1$ -norm not only in denoising but also in general signal recovery tasks and compressive sensing [12–14].

Sparse optimization using the  $\ell_1$ -norm has been also considered from a geometrical viewpoint. It has been reported in [5] that the LASSO solution moves

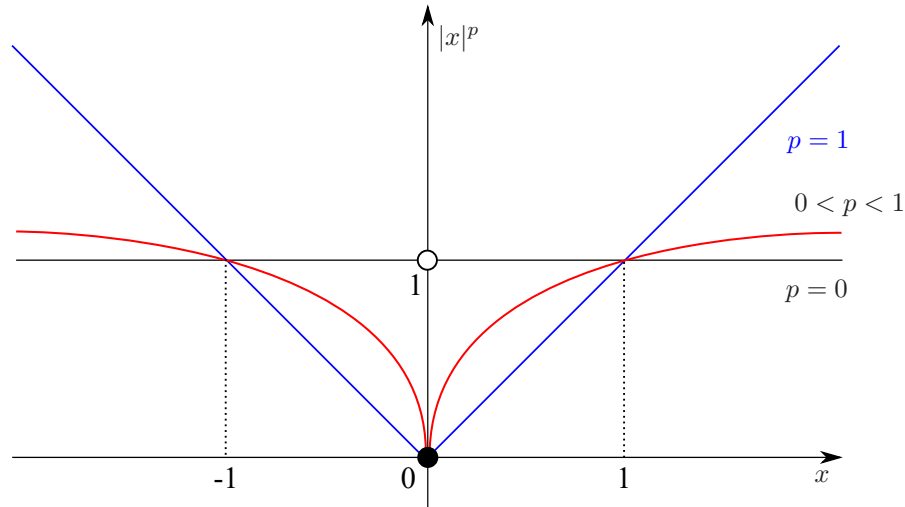


Figure 1.2: The  $\ell_0$ -norm, the  $\ell_1$ -norm, and the  $\ell_p$ -norm ( $0 < p < 1$ ).

when a hyperparameter for an  $\ell_1$ -norm constraint changes, and the solution draws a trajectory called a solution path. Osborne *et al.* [15] have shown that the solution path is piecewise straight and proposed the homotopy algorithm which obtains the entire path. Efron *et al.* [16] have proposed another algorithm called least angle regression (LARS). Although LARS is a greedy algorithm which starts from the origin and increases the number of nonzero coefficients, a modified version which allows nullifying a nonzero coefficient is equivalent to the homotopy algorithm. Donoho and Tsaig [17] have reported that the LARS and the homotopy algorithm, which consider an overdetermined system as LASSO does, can be used to solve an  $\ell_1$ -regularized least squares problem which considers an underdetermined system. In [17], the OMP has also been compared with the LARS, focused on the similarity in their greedy processes.

While the  $\ell_1$ -norm has been popular, the  $\ell_p$ -(quasi)-norm, where  $0 < p < 1$ , has also been considered as an attractive substitute of the  $\ell_0$ -norm for sparse optimization because it is a closer continuous approximation of the  $\ell_0$ -norm than the  $\ell_1$ -norm. Chartrand and Steneva [18] have presented an RIP for the  $\ell_p$ -norm and Foucart and Lai [19] have shown another sufficient condition for achieving the sparsest solution of an underdetermined system by  $\ell_p$ -regularization. Xu *et al.* [20] have focused on the case of  $p = 0.5$  and solved an  $\ell_{0.5}$ -regularized least squares problem using a thresholding operator derived from cubic equations. See the related algorithms in [21, 22]. Yukawa and Amari [23] have extended LARS to the  $\ell_p$  case. They have observed that the solution path under  $\ell_p$ -regularization is discontinuous, while the path consisting of the critical points is continuous and piecewise smooth. (See Figure 1.3.) It is shown that the breakpoints of the piecewise smooth path correspond to the step-by-step solution of the OMP, and this is a more direct link than that shown in [17].

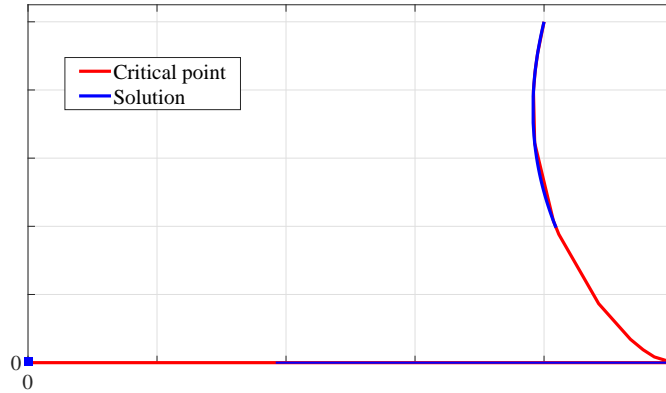


Figure 1.3: The solution path and the critical-point path under an  $\ell_p$ -regularization.

## 1.2 Kernel Adaptive Filtering

The core idea of kernel methods is to replace an inner product by a positive definite kernel function. This is eventually equivalent to map the data in a parametric space onto the higher-dimensional reproducing kernel Hilbert space (RKHS) associated with the positive definite kernel function. In addition, the representer theorem guarantees that the minimizer of a regularized loss function can be represented by a linear combination of a finite number of the basis functions of the RKHS. By the above properties, calculation and representation of an estimator in a possibly infinite-dimensional RKHS can be done by a finite number of calculations in a parametric space.

Enjoying such an advantage, kernel methods have been studied not only in batch settings, but also in online settings. Adaptive filtering [24, 25] is a representative online framework for linear estimation in the signal processing community. *Kernel Adaptive Filtering* [26] has emerged from direct application of an RKHS to the existing adaptive filtering algorithms. Engel *et al.* [27] have proposed an extension of the recursive least squares (RLS) algorithm. Kernelized versions of the affine projection algorithm (APA) and the least mean square (LMS) algorithm have been proposed by Liu *et al.* [28, 29]. Richard *et al.* [30] have also proposed the kernelized APA and LMS algorithms independently from [28, 29]. The main differences between the algorithms in [28, 29] and [30] are that (i) the formers explicitly consider calculations in an RKHS while the latter only use kernelized input in the parametric space and (ii) the latter are normalized algorithms. For simplicity, hereafter the two kernelized LMS algorithms in [29] and [30] are called the KLMS algorithm and the KNLMS algorithm, respectively.

Not only for kernel adaptive filtering, but most online kernel methods face two major challenges. One is undesirable growth of the dictionary, which contains the basis functions representing the estimator. Exploiting an infinite-dimensional RKHS inherently allows an excessively large basis. For online learning of large datasets, in particular, this implies the dictionary size and

the computational costs also increase linearly to the number of observed data. Wang *et al.* [31] have called such an undesired growth of dictionary as *the curse of kernelization*. The other is difficulty of kernel selection, which significantly affects on the performance of kernel methods including kernel adaptive filtering. It is practically impossible to choose a proper kernel before a sufficient number of observations. Kernel methods have been studied from various viewpoints to tackle the above challenges.

For kernel adaptive filtering, *sparsification* of the dictionary using certain criteria has been considered even in earlier works as [27]. The sparsification criteria are based on one or several measures such as a novelty [32], the approximate linear dependence (ALD) [27], the coherence [30], and the surprise [33]. The ALD measures a distance in an RKHS between the newly observed data and its projection onto the subspace spanned by all the basis functions in the dictionary. Van Vaerenbergh *et al.* [34] have found that, from the Bayesian viewpoint, the ALD can be regarded as the prediction variance which is used in the online Gaussian process regression [35]. The coherence has been already considered in greedy algorithms for sparse optimization [1, 4, 6, 36]. The use of the coherence for kernel adaptive filtering has been proposed in [30] and the quantized KLMS (QKLMS) [37] employs essentially the same criterion. The surprise is an information-theoretical measure which generalizes the novelty proposed by Platt [32] and the above other measures. Liu *et al.* [33] has shown that both of the ALD and the coherence are particular cases of the surprise which can be derived by assuming the Gaussian input distribution and unknown desired signal. Soft thresholding operator has also been used for sparse kernel adaptive filtering in [38].

To circumvent the difficulty of kernel selection, a simultaneous use of multiple kernels has been studied. Inspired by the multiple kernel learning (MKL) [39–44] which has emerged in the machine learning community, many kernel adaptive filtering using multiple kernels have been proposed [45, 46]. Yukawa has proposed another way of simultaneously using multiple kernels, called *multikernel adaptive filtering* [47, 48]. While the other methods consider a convex combination of kernels or an ensemble of a number of kernel adaptive filters, multikernel adaptive filtering takes different convex combinations so that different kinds of basis functions are included in the dictionary. Such a use of multiple kernel allows multikernel adaptive filtering a higher degree of freedom in estimation than the other methods [49]. Effectiveness of multikernel adaptive filtering has been shown in several applications including communications [50–52] and sensor network [53].

### 1.3 This Study

The main topics of this thesis are twofold.

- One is on critical-point paths of an  $\ell_p$ -regularized least squares problem in an underdetermined system. Although the study by Yukawa and Amari [23] have presented a LARS extension in an overdetermined system and observed a continuous path of critical points, the existence of a path connecting the origin and the ordinary least squares (OLS) solution has not been proved. In Chapter 3, it is proved that (i) there exists a continuous critical-point path connecting the origin and the (unique) OLS

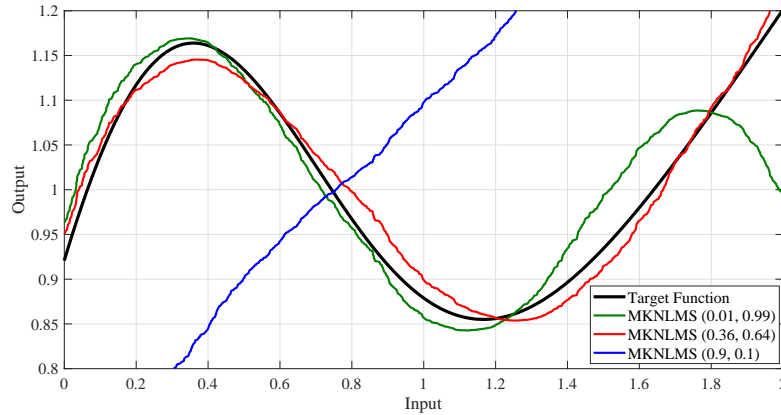


Figure 1.4: Effect of kernel weights on multikernel adaptive filtering.

solution in an overdetermined system, and (ii) there exists a continuous critical-point path connecting the origin and the sparsest OLS solution in an underdetermined system.

- The other is on a kernel weight design for multikernel adaptive filtering. The performance of multikernel adaptive filtering heavily depends on the choice of kernel weights, as shown by Figure 1.3, whereas the weights are learned in MKL. In Figure 1.3, the target function is estimated using two kernels: one is linear and the other is Gaussian. In the legend,  $\text{MKNLMS}(a, b)$  means that the kernel weights for the linear kernel and the Gaussian kernel are  $a$  and  $b = 1 - a$ , respectively. The proposed weight design in Chapter 4 is derived by equalizing the convergence speeds of kernel-wise coefficient errors. The derived weights can be calculated recursively using the observed data and the recursive kernel weighting is implemented to a multikernel adaptive filtering algorithm using metric projection. Numerical results show the superiority of the proposed algorithm over both of the carefully tuned preset weights and an online regression algorithm based on MKL, in the senses of the accuracy and the dictionary size.

The remainder of this thesis is organized as follows. Chapter 2 provides mathematical preliminaries on sparse optimization and online learning with kernels. The related works are introduced as well. In Chapter 3, sparse optimization by  $\ell_p$ -regularization is considered as a fundamental study for sparsification of the dictionary. Chapter 4 proposes a kernel weight design for multikernel adaptive filtering and shows that the proposed design is optimal in a certain sense. Chapter 5 summarizes the previous chapters and gives concluding remarks of this thesis as well as the future prospects for this study.

## Chapter 2

# Preliminaries

The main results of this thesis are derived from discussions dealing with various mathematical concepts. This chapter presents preliminary materials for the discussions in the main chapters.

### 2.1 Notations

Let  $\mathbb{R}$  and  $\mathbb{N}$  denote the sets of all real numbers and all nonnegative integers, respectively.  $\mathbb{R}^L$  and  $\mathbb{R}^{L_1 \times L_2}$  are the  $L$ -dimensional Euclidean space and the set of all  $L_1 \times L_2$  matrices, respectively. In addition, define  $\mathbb{R}_{++}^L := \{\mathbf{x} \in \mathbb{R}^L : x_i > 0, \forall i\}$ . The superscript  $\top$  stands for the transposition of a vector or a matrix.  $\mathbf{I}_r$  and  $\mathbf{O}$  are the  $r \times r$  identity matrix and the zero matrix of appropriate size, respectively.  $\mathbf{1}_r$  and  $\mathbf{0}$  are the vector of  $r$  ones and the zero vector of appropriate size, respectively. For a scalar  $a$ ,  $|a|$  stands for the absolute value of  $a$ . For a set  $\mathcal{I}$ ,  $|\mathcal{I}|$  is the cardinality, or size, of  $\mathcal{I}$ . For a vector  $\mathbf{x} \in \mathbb{L}$ , the  $\ell_p$ -norm of  $\mathbf{x}$  for  $0 < p < \infty$  is defined as follows:

$$\|\mathbf{x}\|_p = \left( \sum_{i=1}^L |x_i|^p \right)^{1/p}, \quad (2.1)$$

while the  $\ell_0$ -norm  $\|\mathbf{x}\|_0$  is the number of nonzero entries of  $\mathbf{x}$ . The range space and the null space of a matrix are denoted by  $\mathcal{R}(\cdot)$  and  $\mathcal{K}(\cdot)$ , respectively.

### 2.2 $\ell_p$ -regularized Least Squares Problem

Consider the linear system model

$$\mathbf{d} := \mathbf{U}^\top \mathbf{h}^* + \boldsymbol{\epsilon} \in \mathbb{R}^l, \quad (2.2)$$

where the sensing matrix  $\mathbf{U} := [\mathbf{u}_1 \cdots \mathbf{u}_l] \in \mathbb{R}^{m \times l}$  and the measurement vector  $\mathbf{d}$  are given, and  $\mathbf{h}^* \in \mathbb{R}^m$  and the noise vector  $\boldsymbol{\epsilon} \in \mathbb{R}^l$  are unknown. Assume that  $\mathbf{U}$  has full rank; *i.e.*,  $\text{Rank}(\mathbf{U}) = \min\{m, l\}$ . The linear system  $\mathbf{U}^\top \mathbf{h} = \mathbf{d}$ ,  $\mathbf{d} \in \mathbb{R}^l$ , has infinitely many OLS solutions if the system is underdetermined ( $m > l$ ), whereas it has a unique one if the system is overdetermined ( $m \leq l$ ).

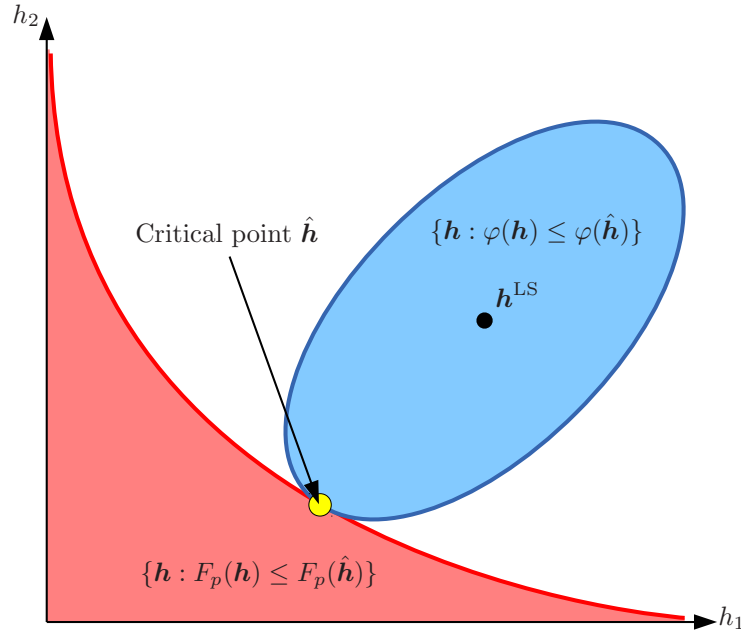


Figure 2.1: A geometric interpretation of critical point.

A problem of finding a sparsest OLS solution is formulated as follows:

$$\text{find } \mathbf{h}_{\text{sparsest}} \in \underset{\mathbf{h} \in \mathbb{R}^m}{\text{argmin}} \|\mathbf{h}\|_0 \text{ s.t. } \mathbf{d} = \mathbf{U}^T \mathbf{h}. \quad (2.3)$$

The  $\ell_0$ -norm involves combinatorial problem and hence solving (2.3) requires unaffordable computational costs for large systems. Consider the following  $\ell_p$ -regularized least squares problem for  $p > 0$ :

$$(\mathcal{L}_\lambda^p): \text{find } \hat{\mathbf{h}}_\lambda \in \underset{\mathbf{h} \in \mathbb{R}^m}{\text{argmin}} f_\lambda(\mathbf{h}) := \varphi(\mathbf{h}) + \lambda F_p(\mathbf{h}), \quad (2.4)$$

where  $\varphi(\mathbf{h}) := \frac{1}{2} \|\mathbf{U}^T \mathbf{h} - \mathbf{d}\|_2^2$ ,  $F_p(\mathbf{h}) := \frac{1}{p} \|\mathbf{h}\|_p^p = \frac{1}{p} \sum_{i=1}^m |h_i|^p$ , and  $\lambda \geq 0$  is the regularization parameter.

It has been shown in [23] that the set of global minima (or even local minima) for every possible  $\lambda \geq 0$  is discontinuous. In this thesis, the critical point is considered as a generalized concept of the stationary points. Indeed,  $F_p$  is differentiable only at those points which have no zero component. One cannot therefore define a stationary point of  $f_\lambda$  in the entire space, but instead can define a critical point by considering the partial derivatives of  $f_\lambda$  in terms of only nonzero components. Prior to the definition of a critical point, define the *support* of a vector as follows.

**Definition 1** (Support). *The support of a vector  $\mathbf{h}$  is defined as  $\text{supp}(\mathbf{h}) := \{i \in \{1, 2, \dots, n\} : h_i \neq 0\}$ , i.e., it is the set of the active indices of  $\mathbf{h}$ .*

For  $\mathcal{I} = \{i_1, i_2, \dots, i_m\} := \text{supp}(\mathbf{h})$  with its cardinality  $|\mathcal{I}| = m$ , denote by  $\nabla_{\mathcal{I}}$  the gradient with respect to the nonzero components; i.e., for any function



$f(\mathbf{h}) : \mathbb{R}^m \mapsto \mathbb{R}$ ,

$$\nabla_{\mathcal{I}} f(\mathbf{h}) := \left[ \frac{\partial f}{\partial h_{i_1}} \quad \frac{\partial f}{\partial h_{i_2}} \quad \cdots \quad \frac{\partial f}{\partial h_{i_m}} \right]^{\top}. \quad (2.5)$$

**Definition 2** (Critical point). *A vector  $\mathbf{h}$  is said to be a critical point of  $f_{\lambda}$  if it satisfies the following condition:*

$$\nabla_{\mathcal{I}} \varphi(\mathbf{h}) + \lambda \nabla_{\mathcal{I}} F_p(\mathbf{h}) = \mathbf{0}. \quad (2.6)$$

Definition 2 implies that the functions  $\varphi$  and  $F_p$  share the same tangent plane at a critical point. See Figure 2.1. Hereafter, the first order condition (2.6) is referred to as the critical point condition.

## 2.3 Multikernel Adaptive Filtering with Weighted Kernels

### 2.3.1 Multikernel Adaptive Filtering

Kernel adaptive filtering considers the following nonlinear system model:

$$d_n = \psi(\mathbf{u}_n) + \nu_n \in \mathbb{R}, \quad n \in \mathbb{N}, \quad (2.7)$$

where  $\mathbf{u}_n \in \mathbb{R}^L$  is the input vector,  $d_n$  is the corresponding output generated by the unknown function  $\psi$  of interest, and  $\nu_n$  is the zero-mean additive noise which is assumed independent of  $\mathbf{u}_n$ . Multikernel adaptive filtering simultaneously exploits multiple kernel functions to estimate the target function  $\psi$  given the observation  $\mathbf{u}_n$  and  $d_n$  up to each time instant  $n$ . Suppose that  $Q$  kernel functions  $\kappa_q(\cdot, \cdot)$ ,  $q = 1, 2, \dots, Q$ , are used. The dictionary for the  $q$ th kernel is defined as  $\mathcal{D}_{q,n} := \{\kappa_q(\mathbf{x}_1^{(q)}, \cdot), \dots, \kappa_q(\mathbf{x}_{r_{q,n}}^{(q)}, \cdot)\}$ , where  $r_{q,n}$  is the cardinality of  $\mathcal{D}_{q,n}$  and the vectors  $\mathbf{x}_j^{(q)} \in \mathbb{R}^L$ ,  $j = 1, 2, \dots, r_{q,n}$ , determine the dictionary elements. The instantaneous estimate is then obtained as

$$\hat{\varphi}_n(\mathbf{u}_n) = \sum_{q=1}^Q \sum_{j=1}^{r_{q,n}} h_{n,j}^{(q)} \kappa_q(\mathbf{x}_j^{(q)}, \mathbf{u}_n) = \sum_{q=1}^Q \mathbf{k}_n^{(q)\top} \mathbf{h}_n^{(q)}, \quad (2.8)$$

where

$$\begin{aligned} \mathbf{k}_n^{(q)} &= [\kappa_q(\mathbf{x}_1^{(q)}, \mathbf{u}_n), \dots, \kappa_q(\mathbf{x}_{r_{q,n}}^{(q)}, \mathbf{u}_n)]^{\top} \in \mathbb{R}^{r_{q,n}}, \\ \mathbf{h}_n^{(q)} &= [h_{n,1}^{(q)}, \dots, h_{n,r_{q,n}}^{(q)}]^{\top} \in \mathbb{R}^{r_{q,n}}, \end{aligned}$$

for all  $q = 1, 2, \dots, Q$ . Concatenating the vectors as

$$\begin{aligned} \mathbf{k}_n &= [\mathbf{k}_n^{(1)\top}, \dots, \mathbf{k}_n^{(Q)\top}]^{\top} \in \mathbb{R}^{r_n}, \\ \mathbf{h}_n &= [\mathbf{h}_n^{(1)\top}, \dots, \mathbf{h}_n^{(Q)\top}]^{\top} \in \mathbb{R}^{r_n}, \end{aligned}$$

where  $r_n = \sum_{q=1}^Q r_{q,n}$ , reduces (2.8) to

$$\hat{\varphi}_n(\mathbf{u}_n) = \mathbf{k}_n^{\top} \mathbf{h}_n. \quad (2.9)$$

The nonlinear estimation of  $\psi$  then reduces to linear estimation of the Euclidean vector  $\mathbf{h}^*$ , which is a minimizer of some error function, given the *kernelized* input vector  $\mathbf{k}_n$ ,  $n \in \mathbb{N}$ . For  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^L$ , the celebrated examples of positive definite reproducing kernel are given as below.

**Normalized Gaussian kernel**

$$\kappa_G(\mathbf{u}, \mathbf{v}) = \frac{1}{(\sqrt{2\pi}\sigma)^L} \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|_2^2}{2\sigma^2}\right), \quad \sigma > 0. \quad (2.10)$$

**Linear kernel**

$$\kappa_L(\mathbf{u}, \mathbf{v}) = \mathbf{u}^\top \mathbf{v} + c, \quad c \in \mathbb{R}. \quad (2.11)$$

**Polynomial kernel**

$$\kappa_P(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^\top \mathbf{v} + c)^a, \quad c \in \mathbb{R}, a \in \mathbb{N}. \quad (2.12)$$

In the multikernel adaptive filtering algorithms, the dictionaries  $\mathcal{D}_{q,n}$  for the Gaussian kernel are constructed in online fashion, as in the case of the other kernel adaptive filtering algorithms, whereas the dictionary for the linear kernel is fixed. This is because the RKHS for the linear kernel has finite dimension and its basis is known a priori whereas that for the Gaussian kernel has infinite dimension and an appropriate dictionary is unknown in practice. Some criteria for dictionary construction can be found in the literature [27, 30, 32, 33, 35]. A popular criterion for the Gaussian kernel is the coherence criterion [30], which tests whether the maximal entry of the kernelized input vector  $\mathbf{k}_n^{(q)}$  is below a prespecified threshold, say  $\delta > 0$ . If the maximal entry is below  $\delta$ , then the input is considered to be novel (since it is incoherent to any existing elements of the dictionary) and is thus added to the dictionary. The resulting dictionary consists of incoherent elements, meaning that the centers of the Gaussian functions in the dictionary are apart from each other at the prespecified level.

The MKNLMS algorithm [47, 48] is a multikernel adaptive filtering algorithm of which the error function is the mean squared error  $\text{MSE}(\mathbf{h}) := E[(\mathbf{k}_n^\top \mathbf{h} - d_n)^2]$ , where  $\mathbf{h} \in \mathbb{R}^{r_n}$  and  $E[\cdot]$  stands for the expectation with respect to the input  $\mathbf{u}_n$  and the noise  $\nu_n$ . Without a manipulation of dictionaries, each update of the MKNLMS algorithm can be written as follows:

$$\mathbf{h}_{n+1} = \mathbf{h}_n - \mu \frac{\mathbf{k}_n^\top \mathbf{h}_n - d_n}{\|\mathbf{k}_n\|_2^2} \mathbf{k}_n, \quad (2.13)$$

where  $\mu \in (0, 2)$  is the step size. When  $Q = 1$ , the MKNLMS algorithm reduces to the KNLMS algorithm [30].

### 2.3.2 Relation to MKL

MKL has extensively been studied in machine learning community. In a similar form to (2.8), the instantaneous estimate by MKL can be written as follows:

$$\begin{aligned} \hat{\varphi}_n(\mathbf{u}_n) &= \sum_{j=1}^{r_n} h_{n,j} \left( \sum_{q=1}^Q \theta_{q,n} \kappa_q(\mathbf{x}_j, \mathbf{u}_n) \right) \\ &= \sum_{j=1}^{r_n} h_{n,j} \kappa_n(\mathbf{x}_j, \mathbf{u}_n), \end{aligned} \quad (2.14)$$

**Algorithm 1** OMKR (hedge)

- 
- 0)  $\mathcal{D}_0 = \emptyset$ ,  $\mathbf{h}_0^{(q)} = [\ ]$ ,  $\theta_{q,1} = \frac{1}{Q}$ ,  $q = 1, 2, \dots, Q$ ,  $\mu > 0$ : step size,  
 $\beta \in (0, 1)$ : discounting parameter.
  - 1) Prediction:  $\hat{\varphi}_n(\mathbf{u}_n) = \sum_{q=1}^Q \theta_{q,n} \mathbf{k}_n^{(q)\top} \mathbf{h}_n^{(q)}$ .
  - 2)  $\epsilon_{q,n} = d_n - \mathbf{k}_n^{(q)\top} \mathbf{h}_n^{(q)}$ .
  - 3)  $\mathbf{h}_{n+1}^{(q)} = \begin{bmatrix} \mathbf{h}_n^{(q)} \\ \mu \epsilon_{q,n} \end{bmatrix}$ ,  $\tilde{\theta}_{q,n+1} = \theta_{q,n} \beta^{\epsilon_{q,n}^2}$ .
  - 4)  $\theta_{q,n+1} = \frac{\tilde{\theta}_{q,n+1}}{\sum_{p=1}^Q \tilde{\theta}_{p,n+1}}$ .
  - 5)  $\mathcal{D}_{q,n+1} = \mathcal{D}_{q,n} \cup \{\kappa_q(\mathbf{u}_n, \cdot)\}$ .
- 

where  $\theta_{q,n} \geq 0$ ,  $q = 1, 2, \dots, Q$ , are the kernel weights such that  $\sum_{q=1}^Q \theta_{q,n} = 1$ . Most MKL algorithms, such as those in [42–44], first update  $\theta_{q,n}$  to construct a convex combination of multiple kernels, and then update the coefficients  $h_{n,j}$ . One can see from (2.14) that MKL uses the appropriately designed single kernel  $\kappa_n$ , and the degree of freedom in MKL is  $Q + r_n$ . On the other hand, as seen from (2.8), the degree of freedom in multikernel adaptive filtering is  $Q \times r_n$ , which is much larger than  $Q + r_n$  in practical situations. Thanks to the higher degree of freedom, multikernel adaptive filtering enjoys a more compact representation than MKL by excluding those kernels which are unnecessary for estimation. Moreover, multikernel adaptive filtering can potentially extract the local structures of the target containing multiple components, since different kernels can be used to express different components. In Section 5, the proposed method is experimentally compared with the OMKR algorithm [43, 44]. The OMKR algorithm has been developed as an MKL algorithm for online regression, while most MKL algorithms deal with classification tasks under a batch setting. It is therefore a particular example of MKL that can be directly compared with multikernel adaptive filtering in numerical examples. Algorithm 1 shows the OMKR algorithm with a discounting parameter, which controls the kernel weights according to the corresponding kernel-wise prediction errors.

### 2.3.3 Imposing Kernel Weights

A severe degradation in performance of multikernel adaptive filtering may occur when there exists some imbalance among the kernels. Such imbalance can be alleviated by introducing kernel weights. Given positive definite kernels  $\kappa_1, \kappa_2, \dots, \kappa_Q$ , imposing weights  $w_q > 0$ ,  $q = 1, 2, \dots, Q$ , on the kernels changes the geometry of the original RKHS [48]. Geometry of the space is important for multikernel adaptive filtering, of which the update is based on the projection onto a zero-instantaneous-error hyperplane. The use of unweighted kernels, as in Figure 2.2(a), leaves the possible imbalance among the kernels and may result in slow convergence of the errors which are associated with some of the kernels. Figure 2.2(b) shows the change of geometry by the appropriate kernel weights such that the imbalance is alleviated. Applying the kernel weights to the update

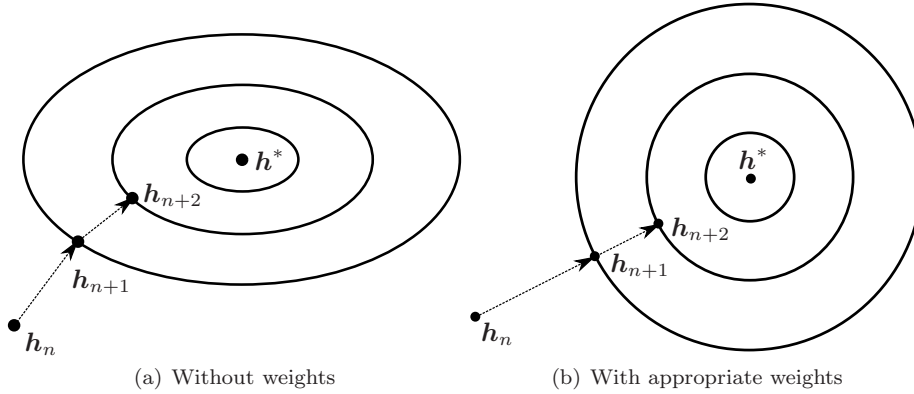


Figure 2.2: Change of geometry by kernel weights.

(2.13) yields

$$\mathbf{g}_{n+1} = \mathbf{g}_n - \mu \frac{\mathbf{k}_n^\top \mathbf{W} \mathbf{g}_n - d_n}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \mathbf{W} \mathbf{k}_n, \quad (2.15)$$

where  $\mathbf{g}_n := \mathbf{W}^{-1} \mathbf{h}_n$  and the weight matrix  $\mathbf{W}$  is given as

$$\mathbf{W} := \begin{bmatrix} w_1 \mathbf{I}_{r_{1,n}} & \mathbf{O} & \cdots & \mathbf{O} \\ \mathbf{O} & w_2 \mathbf{I}_{r_{2,n}} & \cdots & \mathbf{O} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{O} & \mathbf{O} & \cdots & w_Q \mathbf{I}_{r_{Q,n}} \end{bmatrix}.$$

The choice of  $\mathbf{W}$  governs the transient behavior of the coefficient error  $\mathbf{z}_n := \mathbf{g}_n - \mathbf{g}^*$ , where  $\mathbf{g}^* := \mathbf{W}^{-1} \mathbf{h}^*$ . Specifically, relative convergence rates among the kernel-wise coefficient errors, which coincide with the subvectors of  $\mathbf{z}_n$ , are determined by the ratios among the weights  $w_1, w_2, \dots, w_Q$ , as well as the input autocorrelation matrix. An appropriate choice of kernel weights leads to uniform convergence rates of the kernel-wise coefficient errors. The kernel weights, however, need to be chosen carefully because a wrong choice may result in a worse performance than the use of unweighted kernels and even than the use of a single kernel.

## Chapter 3

# Critical-Point Paths under $\ell_p$ -regularization ( $0 < p < 1$ )

### 3.1 Introduction

Regularization by the  $\ell_p$ -norm ( $0 < p < 1$ ) has been regarded as an attractive approach of sparse optimization, since (i) the  $\ell_p$ -norm is a better approximation of the  $\ell_0$ -norm (the number of nonzero entries); and (ii) it has been reported that solving an  $\ell_p$ -norm minimization problem yields a sparser solution than minimizing the  $\ell_1$ -norm [18, 19, 22, 54, 55]. By extending the idea of the LARS [16], the solution path of a least squares problem regularized by the  $\ell_p$ -norm has been investigated [23, 56, 57]. The solution path consists of all the solutions, each of which corresponds to a value of the regularization parameter. In the  $\ell_1$  case, the solution moves along a gradual change of the regularization parameter, and hence the solution path is continuous. On the other hand, in the  $\ell_p$  case, it has been reported that the solution suddenly jumps, whereas the corresponding regularization parameter changes gradually [23, 57]. Because of nonconvexity of the  $\ell_p$ -norm, the  $\ell_p$ -regularized least squares problem can have multiple critical points, where the gradient is zero, including the solutions. (See also Definition 2.) By tracking the critical points with changing the regularization parameter, one can obtain a continuous and piecewise smooth paths.

In this chapter, it is shown that every sparsest OLS solution is connected to the origin by a continuous path of critical points (Theorem 3.2). The path of critical points is shortly referred to as a *critical-point path*. The critical-point path is investigated first in an overdetermined system, where the OLS solution is unique. It is proved that (i) there exists a unique path emerging from the OLS solution and (ii) it does not stop (the velocity vector of the path does not become zero) as long as the support does not change (see Lemma 3.1). This implies that the path eventually falls out of the region corresponding to the current support unless the path ends up with some loop(s). It is ensured automatically that a path does not bifurcate excluding those critical points at which the Hessian matrix of the  $\ell_p$ -regularized squared error function is singular. Such a singular point always exists on a path indeed (see Proposition 3.2), and therefore it is assumed that the path does not bifurcate at any singular points in order to make the analysis feasible. Although the assumption excluding bifurcation is

required for the main theorem, the theorem holds even when the assumption is violated (see Remark 3.1).

The above argument tells us that the path goes continuously in the region corresponding to the current support and finally reaches a point out of the region. The support of this end point is a proper subset of that of the OLS solution (which means that the end point has at least one zero component), and it is the OLS solution of a modified squared error function forcing those components associated with the off-support to be zero. (See Proposition 3.1.) Repetitions of this lead the path to the origin, which verifies the existence of a piecewise smooth path (with each end point being its breakpoint) from the OLS solution to the origin in the overdetermined case. (See Theorem 3.1.) The basic idea of the proof for an underdetermined case is reducing the problem to an overdetermined one in terms of a small number of components associated with the support of an OLS solution. The reduction is only possible for such a solution whose support has its associated submatrix of the fat sensing matrix be full column rank. (See Theorem 3.3.) It is proved that being full-rank is a necessary condition for a reduced overdetermined system to include the sparsest OLS solution of the entire underdetermined system. (See Lemma 3.2) There is also a possibility that a path exists even though the problem cannot be reduced to an overdetermined one. For such an issue, a necessary condition for an OLS solution to be connected with the origin by a continuous path is presented. (See Proposition 3.3.)

## 3.2 A Path to the OLS Solution

An overdetermined system, where the OLS solution is unique, is considered in this section. Figure 3.1 illustrates critical-point paths, which are obtained by gradually increasing the  $\ell_p$ -norm of critical points. One can see that the blue path connects the origin and the OLS solution  $\mathbf{h}^{\text{LS}}$ , while the red and green paths reach some other point. The existence of a critical-point path connecting the origin and the OLS solution is proved under a mild condition. It is also shown that the path passes through a point at which the Hessian matrix of  $f_\lambda := \varphi(\mathbf{h}) + \lambda F_p(\mathbf{h})$  is singular.

### 3.2.1 Existence of a Continuous Critical-point Path

Let us arbitrarily choose a vector  $\hat{\mathbf{h}} \in \mathbb{R}^m$  with its support  $\mathcal{I} := \text{supp}(\hat{\mathbf{h}})$  satisfying  $\nabla_{\mathcal{I}}\varphi(\hat{\mathbf{h}}) = \mathbf{0}$ , and define

$$S_{\mathcal{I}} := \{\mathbf{h} : \text{supp}(\mathbf{h}) = \mathcal{I}\} (\ni \hat{\mathbf{h}}). \quad (3.1)$$

The matrix  $\mathbf{G} := \mathbf{U}\mathbf{U}^{\text{T}} \in \mathbb{R}^{m \times m}$  is positive definite in an overdetermined case, from which it can be verified that  $\hat{\mathbf{h}}$  is a unique critical point in  $S_{\mathcal{I}}$  for  $\lambda = 0$ .

The critical point condition (2.6) can be rewritten as

$$\mathbf{U}_{\mathcal{I}}\mathbf{U}^{\text{T}}(\mathbf{h} - \mathbf{h}^{\text{LS}}) + \lambda \nabla_{\mathcal{I}}F_p(\mathbf{h}) = \mathbf{0}, \quad (3.2)$$

where  $\mathbf{U}_{\mathcal{I}} \in \mathbb{R}^{|\mathcal{I}| \times l}$  is a submatrix of  $\mathbf{U}$  consisting of those row vectors associated with the support  $\mathcal{I}$ . Since

$$\nabla_{\mathcal{I}}\varphi(\hat{\mathbf{h}}) = \mathbf{U}_{\mathcal{I}}\mathbf{U}^{\text{T}}(\hat{\mathbf{h}} - \mathbf{h}^{\text{LS}}) = \mathbf{0}, \quad (3.3)$$

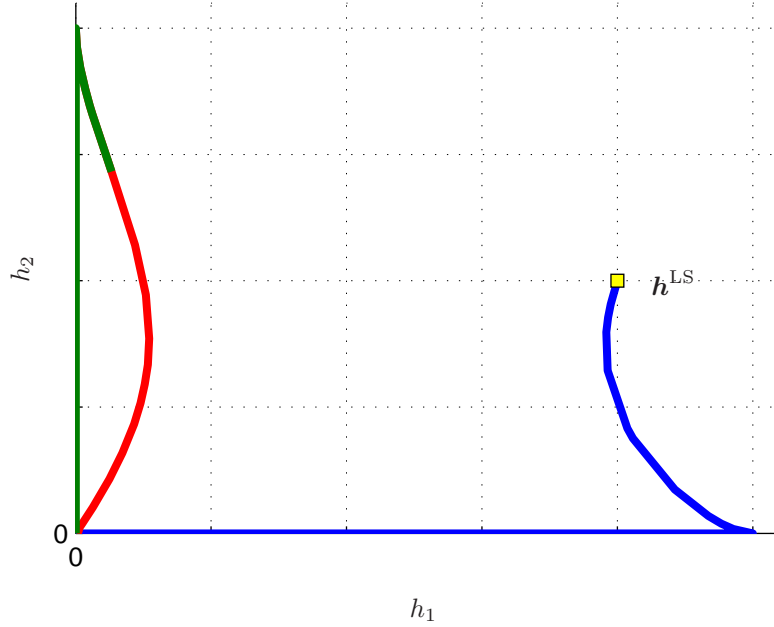


Figure 3.1: Critical-point paths in a two-dimensional overdetermined case.

one can obtain

$$\begin{aligned}
 \mathbf{U}_{\mathcal{I}}\mathbf{U}^{\top}(\mathbf{h} - \mathbf{h}^{\text{LS}}) - \nabla_{\mathcal{I}}\varphi(\hat{\mathbf{h}}) + \lambda\nabla_{\mathcal{I}}F_p(\mathbf{h}) &= \mathbf{U}_{\mathcal{I}}\mathbf{U}^{\top}(\mathbf{h} - \hat{\mathbf{h}}) + \lambda\nabla_{\mathcal{I}}F_p(\mathbf{h}) \\
 &= \mathbf{G}_{\mathcal{I}}(\mathbf{h}_{\mathcal{I}} - \hat{\mathbf{h}}_{\mathcal{I}}) + \lambda\nabla_{\mathcal{I}}F_p(\mathbf{h}) \\
 &= \mathbf{0}
 \end{aligned} \tag{3.4}$$

for any critical point  $\mathbf{h} \in S_{\mathcal{I}}$ , where  $\mathbf{G}_{\mathcal{I}} := \mathbf{U}_{\mathcal{I}}\mathbf{U}_{\mathcal{I}}^{\top}$  and  $\mathbf{h}_{\mathcal{I}} \in \mathbb{R}^{|\mathcal{I}|}$  is a subvector of  $\mathbf{h}$  consisting of those components associated with  $\mathcal{I}$ . The second equality holds because  $\mathbf{U}^{\top}\mathbf{h} = \mathbf{U}_{\mathcal{I}}^{\top}\mathbf{h}_{\mathcal{I}}$ . A total differential of (3.4) leads to

$$\mathbf{K}_{\mathcal{I}}d\mathbf{h}_{\mathcal{I}} = -d\lambda\nabla_{\mathcal{I}}F_p(\mathbf{h}). \tag{3.5}$$

where  $\mathbf{K}_{\mathcal{I}} := \nabla_{\mathcal{I}}\nabla_{\mathcal{I}}f_{\lambda}(\mathbf{h}) = \mathbf{G}_{\mathcal{I}} + \lambda\nabla_{\mathcal{I}}\nabla_{\mathcal{I}}F_p(\mathbf{h})$  is the Hessian matrix of  $f_{\lambda}$ .

As shown below, the matrix  $\mathbf{K}_{\mathcal{I}}$  becomes singular at some point on a critical-point path. Assume that the following statements hold to exclude the case that a path bifurcates at singular points, guaranteeing that the path has no loop.

**Assumption 3.1.** For any critical point  $\tilde{\mathbf{h}}$  such that the Hessian matrix  $\mathbf{K}_{\mathcal{I}}$  is singular, the following hold:

- (i) the zero eigenvalue of  $\mathbf{K}_{\mathcal{I}}$  is simple; and
- (ii)  $\nabla_{\mathcal{I}}F_p(\tilde{\mathbf{h}}) \notin \mathcal{R}(\mathbf{K}_{\mathcal{I}})$ .

The main theorem of this section is stated below.

**Theorem 3.1.** In an overdetermined case, there exists a continuous path of critical points connecting the origin and the OLS solution.

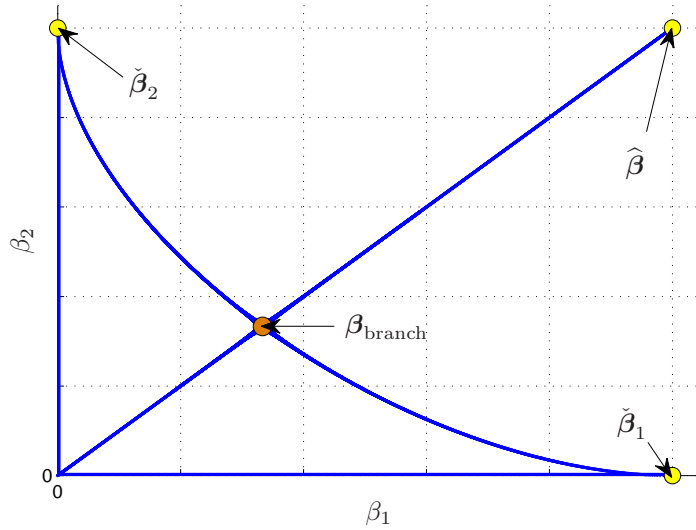


Figure 3.2: A case that Assumption 3.1 is violated but Theorem 3.1 still holds.

**Remark 3.1** (On Assumption 3.1).

(a) Assumption 3.1 is violated in the following cases:

- (i) two eigenvalues become zero coincidentally:
- (ii) the vector  $\nabla_{\mathcal{I}} F_p(\tilde{\mathbf{h}})$  happens to lie in the subspace  $\mathcal{R}(\mathbf{K}_{\mathcal{I}})$  which has no volume in  $\mathbb{R}^{|\mathcal{I}|}$ .

(b) There exists a case where Assumption 3.1 is violated but Theorem 3.1 still holds.

Remark 3.1(b) implies that Assumption 3.1 is a sufficient, but not necessary, condition for Theorem 3.1. Figure 3.2 illustrates the case of  $\mathbf{G} = \mathbf{I}$ ,  $\hat{\boldsymbol{\beta}} = [1 \ 1]^T$  and  $p = 0.5$ ; the Hessian matrix  $\mathbf{K}_{\mathcal{I}}$  becomes a zero matrix at  $\boldsymbol{\beta}_{\text{branch}}$  and has the zero eigenvalues with multiplicity two. As Assumption 3.1 is violated, bifurcation happens on the path emerging from  $\hat{\boldsymbol{\beta}}$  at the branch point  $\boldsymbol{\beta}_{\text{branch}}$  into three branches. However, one can find that all the branches reach the origin.

Theorem 3.1 can be verified by applying the following proposition recursively, as shown later on.

**Proposition 3.1.** *In the region  $S_{\mathcal{I}}$ , there exists a path connecting the  $\hat{\mathbf{h}}$  and some  $\tilde{\mathbf{h}}$  such that (i)  $\text{supp}(\tilde{\mathbf{h}}) = \tilde{\mathcal{I}} \subsetneq \mathcal{I}$  and (ii)  $\nabla_{\tilde{\mathcal{I}}} \varphi(\tilde{\mathbf{h}}) = \mathbf{0}$  ( $\lambda = 0$ ). The edge point  $\tilde{\mathbf{h}}$  lies on the boundary of the open set  $S_{\mathcal{I}}$ .*

Proposition 3.1 will be proved with the following lemma.

**Lemma 3.1.**

- (a) The  $\hat{\mathbf{h}}$  has a neighborhood in which there is a unique path that emerges from  $\hat{\mathbf{h}}$  and lies in  $S_{\mathcal{I}}$ .



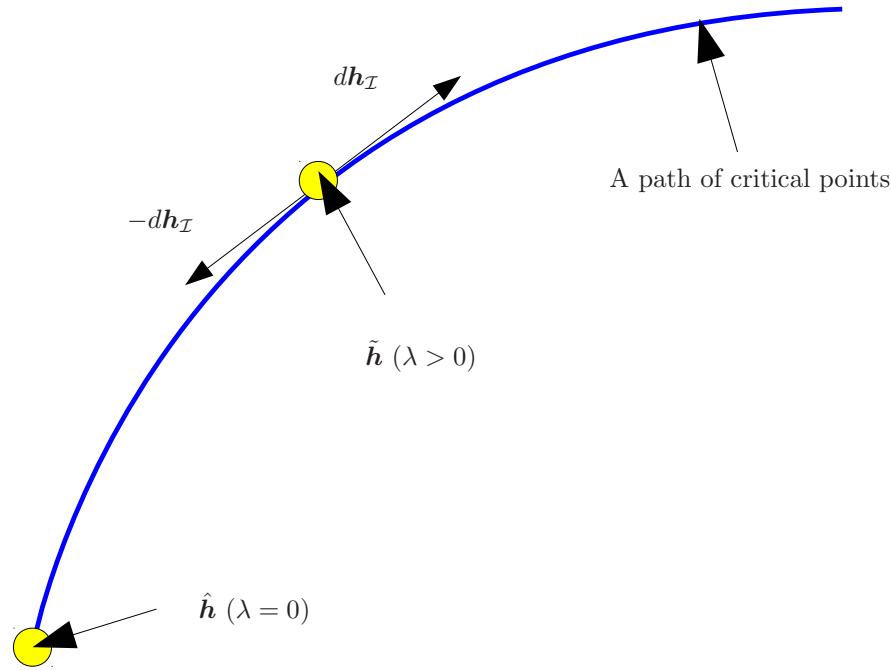


Figure 3.3: An illustration of Lemma 3.1(b).

(b) For any critical point  $\tilde{\mathbf{h}} \in S_{\mathcal{I}}$  for  $\lambda > 0$ , there exists a unique path in  $S_{\mathcal{I}}$  that passes through  $\tilde{\mathbf{h}}$ ; i.e., the path never stops at any critical point for  $\lambda > 0$ .

(c) The path emerging from  $\hat{\mathbf{h}}$  does not form any loop on  $S_{\mathcal{I}}$ .

**Proof of Lemma 3.1(a):** In a sufficiently small vicinity of  $\hat{\mathbf{h}}$ ,  $\mathbf{K}_{\mathcal{I}}$  is nonsingular (because  $\mathbf{G}_{\mathcal{I}}$  is). Left-multiplying  $\mathbf{K}_{\mathcal{I}}^{-1}$  to (3.5) yields the following differential.

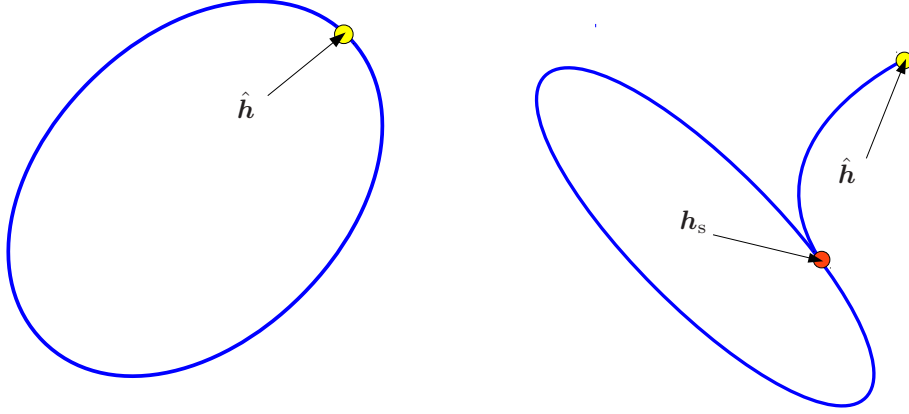
$$\frac{d\mathbf{h}_{\mathcal{I}}}{d\lambda} = -\mathbf{K}_{\mathcal{I}}^{-1} \nabla_{\mathcal{I}} F_p(\hat{\mathbf{h}}). \quad (3.6)$$

This verifies the claim.

**Proof of Lemma 3.1(b):** The case that  $\mathbf{K}_{\mathcal{I}}$  is nonsingular has already been verified in the proof of Lemma 3.1(a). Suppose now that  $\mathbf{K}_{\mathcal{I}}$  is singular. Then, an eigenvector  $\boldsymbol{\varepsilon}_0$  of  $\mathbf{K}_{\mathcal{I}}$  corresponding to the zero eigenvalue satisfies  $\mathbf{K}_{\mathcal{I}} \boldsymbol{\varepsilon}_0 = \mathbf{0}$ , meaning that the pair  $(d\mathbf{h}_{\mathcal{I}}, d\lambda) = (\boldsymbol{\varepsilon}_0, 0)$  satisfies (3.5). With Assumption 3.1, this verifies the existence and uniqueness of a path passing through  $\tilde{\mathbf{h}}$ .

**Proof of Lemma 3.1(c):** Lemma 3.1(a) ensures that the path never returns to the starting point  $\hat{\mathbf{h}}$ . Hence, a loop can exist only if the path bifurcates. However, bifurcation happens neither at nonsingular points (Lemma 3.1(b)) nor at singular points (Assumption 3.1). Figure 3.4 illustrates the cases denied in this work. The cases in Figure 3.4(a) and Figure 3.4(b) are denied by Lemma 3.1(a) and Assumption 3.1, respectively.  $\square$

**Proof of Proposition 3.1:** Lemma 3.1(c) verifies that the path emerging from  $\hat{\mathbf{h}}$  never stops at any critical point  $\tilde{\mathbf{h}} \in S_{\mathcal{I}}$  for  $\lambda > 0$ . By Lemma 3.1(b), the path does not form a loop and finally falls out of  $S_{\mathcal{I}}$ . The support of  $\mathbf{h}$  diminishes



(a) A case denied by Lemma 3.1.

(b) A case denied by Assumption 3.1.

Figure 3.4: The cases denied in Lemma 3.1 and Assumption 3.1.

and this verifies  $\check{\mathcal{I}} \subsetneq \mathcal{I}$ . Let  $h_{i^*}$  be a vanishing component, where  $i^* \in \mathcal{I} \setminus \check{\mathcal{I}}$ . By the critical point condition (2.6), at any point  $\tilde{\mathbf{h}} \in S_{\mathcal{I}}$  on the path excluding the edge point  $\hat{\mathbf{h}} \notin S_{\mathcal{I}}$ , it holds that  $\frac{\partial \varphi}{\partial h_{i^*}}(\tilde{\mathbf{h}}) = -\lambda \frac{\partial F_p}{\partial h_{i^*}}(\tilde{\mathbf{h}})$ . As  $\tilde{\mathbf{h}} \rightarrow \hat{\mathbf{h}}$ ,  $\lambda \rightarrow 0$  since  $\left| \frac{\partial F_p}{\partial h_{i^*}}(\tilde{\mathbf{h}}) \right| \rightarrow \infty$ , while  $\left| \frac{\partial \varphi}{\partial h_{i^*}}(\tilde{\mathbf{h}}) \right| < \infty$ . It is hence verified by (2.6) again that  $\nabla_{\check{\mathcal{I}}} \varphi(\hat{\mathbf{h}}) = \mathbf{0}$ .  $\square$

**Proof of Theorem 3.1:** Let  $\hat{\mathbf{h}} = \mathbf{h}^{\text{LS}}$ . Then, by Proposition 3.1, the  $\hat{\mathbf{h}}$  is connected with another critical point  $\check{\mathbf{h}} \in S_{\check{\mathcal{I}}}$  for some  $S_{\check{\mathcal{I}}} \subsetneq S_{\mathcal{I}}$  and  $\lambda = 0$ . Since  $\check{\mathbf{h}}$  is also a critical point for  $\lambda = 0$ , it is also possible to consider a path in  $S_{\check{\mathcal{I}}}$  emerging from the  $\check{\mathbf{h}}$  by just the same way with the first stage except replacing  $\hat{\mathbf{h}}$  by  $\check{\mathbf{h}}$ . Repeating such stages recursively, the path eventually reaches to the origin corresponding to the empty support, and the existence of a path connecting the origin and  $\mathbf{h}^{\text{LS}}$  is verified.  $\square$

### 3.2.2 Characteristics of Critical-point Paths

Letting both of  $\hat{\mathbf{h}} \in S_{\mathcal{I}}$  and  $\check{\mathbf{h}} \in S_{\check{\mathcal{I}}}$  be critical points for  $\lambda = 0$ , the following proposition on a critical point at which  $\mathbf{K}_{\mathcal{I}}$  is singular can be proved.

**Proposition 3.2.** *A path connecting the  $\hat{\mathbf{h}}$  and  $\check{\mathbf{h}}$  passes through a point at which  $d\lambda = 0$  (see (3.5),) and the Hessian matrix  $\mathbf{K}_{\mathcal{I}}$  of  $f_{\lambda}$  is singular.*

*Proof:* Both  $\hat{\mathbf{h}}$  and  $\check{\mathbf{h}}$  are critical points for  $\lambda = 0$  and the other points on the path are for  $\lambda > 0$ . Therefore,  $\lambda$  along the path emerging from  $\hat{\mathbf{h}}$  increases up to some point and then starts to decrease. At the changing point, it holds that  $d\lambda = 0$ . Substituting  $d\lambda = 0$  into (3.5) yields  $\mathbf{K}_{\mathcal{I}} d\mathbf{h}_{\mathcal{I}} = \mathbf{0}$ . Since  $d\mathbf{h}_{\mathcal{I}} \neq \mathbf{0}$  on the path, it follows that  $\mathbf{K}_{\mathcal{I}}$  is singular at the changing point and  $d\mathbf{h}_{\mathcal{I}}$  is an eigenvector corresponding to the zero eigenvalue.  $\square$

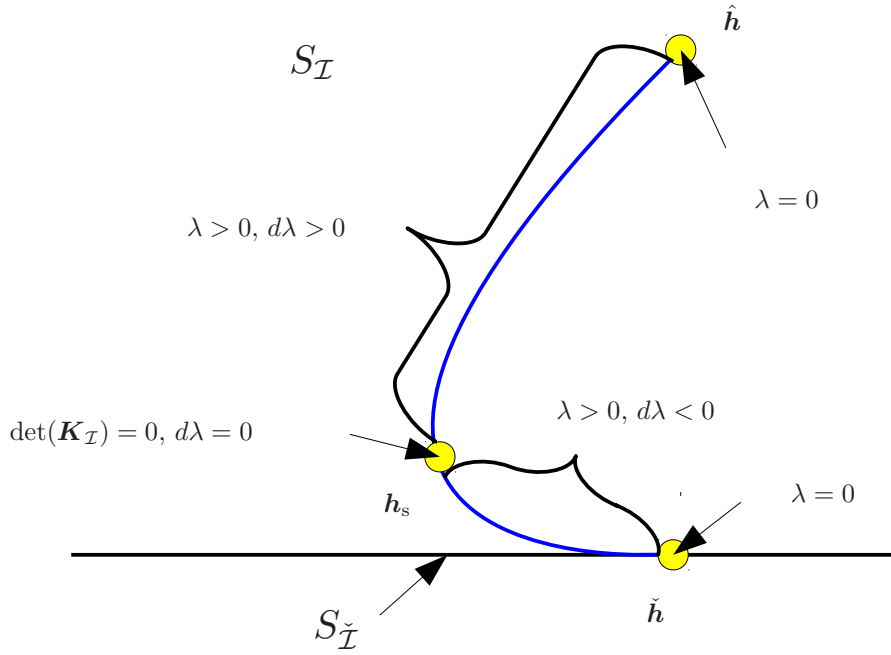


Figure 3.5: An illustration of Remark 3.2.

Below are remarks on the behavior of  $\lambda$  under the light of Proposition 3.2. (See Figure 3.5.)

**Remark 3.2** (Behavior of  $\lambda$  along a path emerging from  $\hat{h}$ ).

- (a) At  $\hat{h}$  ( $\lambda = 0$ ) : The nonsingularity of  $\mathbf{G}_{\mathcal{I}}$  ensures the uniqueness and existence of  $d\mathbf{h}_{\mathcal{I}}$ .
- (b) From  $\hat{h}$  to the singular point  $h_s$  ( $d\lambda = 0$ ) : By Lemma 3.1(c) and Assumption 3.1, if  $\mathbf{K}_{\mathcal{I}}$  is nonsingular at a critical point, a single path passes the point.
- (c) At  $h_s$  : By Proposition 3.2, there is a point at which  $d\lambda = 0$  and  $\mathbf{K}_{\mathcal{I}}$  is singular.
- (d) From  $h_s$  to  $\check{h}$  : The path goes inside  $S_{\mathcal{I}}$  with decreasing  $\lambda$  and falls out of  $S_{\mathcal{I}}$  to reach  $\check{h}$ .

### 3.3 Paths to Sparsest OLS Solutions

The primal interest of this chapter concerns the paths that connect the origin and the set  $V^*$  of OLS solutions in the underdetermined case. More specifically, only the paths containing a single OLS solution are considered. Although many OLS solutions exist, only some of them are connected with the origin by paths

of critical points. The existence of the paths, each of which leads to a sparsest OLS solution, are proved.

The statement below is the main result of this chapter.

**Theorem 3.2.** *Every sparsest OLS solution is connected with the origin by a continuous path of critical points.*

The following lemma is useful to prove Theorem 3.2.

**Lemma 3.2.** *Given a fixed OLS solution  $\hat{\mathbf{h}} \in V^*$ , let  $\mathcal{I} = \text{supp}(\hat{\mathbf{h}})$ . Consider the following three conditions.*

(C1)  $\mathbf{G}_{\mathcal{I}}$  is nonsingular.

(C2)  $\{\mathbf{h} \in V^* : \text{supp}(\mathbf{h}) \subseteq \mathcal{I}\} = \{\hat{\mathbf{h}}\}$ ; i.e.,  $\hat{\mathbf{h}}$  is the unique OLS solution over the subspace in which  $h_i = 0, \forall i \in \bar{\mathcal{I}} := \{1, 2, \dots, n\} \setminus \mathcal{I}$ .

(C3)  $\hat{\mathbf{h}} \in \text{argmin}_{\mathbf{h} \in V^*} \|\mathbf{h}\|_0$ ; i.e.,  $\hat{\mathbf{h}}$  is sparsest among the OLS solutions.

Then, the following statements hold.

(a) (C1) and (C2) are equivalent.

(b) (C3) implies (C1) and (C2).

(c) Assume that

$$\min_{\mathbf{h} \in V^*} \|\mathbf{h}\|_0 = r := \text{Rank}(\mathbf{U}). \quad (3.7)$$

In this case, all the conditions (C1), (C2) and (C3) are equivalent.

**Proof of Lemma 3.2(a):**

(C2)  $\Rightarrow$  (C1): Assume that  $\mathbf{G}_{\mathcal{I}}$  is singular, i.e.,  $\mathcal{K}(\mathbf{U}_{\mathcal{I}}^T) \setminus \{\mathbf{0}\}$  is nonempty. It is shown that  $\hat{\mathbf{h}}$  is only an element in  $\{\mathbf{h} \in V^* : \text{supp}(\mathbf{h}) \subseteq \mathcal{I}\}$ , or equivalently that there exists some  $\tilde{\mathbf{h}} \neq \hat{\mathbf{h}}$  such that  $\tilde{\mathbf{h}} \in V^*$  and  $\text{supp}(\tilde{\mathbf{h}}) \subseteq \mathcal{I}$ . Then,

$$\begin{aligned} \tilde{\mathbf{h}} \in V^* &\Leftrightarrow \nabla \varphi(\tilde{\mathbf{h}}) = \mathbf{0} \\ &\Leftrightarrow \mathbf{U}(\mathbf{U}_{\mathcal{I}}^T \tilde{\mathbf{h}}_{\mathcal{I}} - \mathbf{d}) = \mathbf{0} \\ &\Leftrightarrow \mathbf{U}\mathbf{U}_{\mathcal{I}}^T \tilde{\mathbf{h}}_{\mathcal{I}} = \mathbf{U}\mathbf{d}. \end{aligned} \quad (3.8)$$

The vector  $\hat{\mathbf{h}}_{\mathcal{I}}$  is clearly a solution of (3.8) and the singularity of  $\mathbf{G}_{\mathcal{I}}$  suggests that  $\tilde{\mathbf{h}}_{\mathcal{I}} := \hat{\mathbf{h}}_{\mathcal{I}} + \mathbf{a}$  for any  $\mathbf{a} \in \mathcal{K}(\mathbf{U}_{\mathcal{I}}^T) \neq \{\mathbf{0}\}$  satisfies (3.8).

(C1)  $\Rightarrow$  (C2): Assume that  $\mathbf{G}_{\mathcal{I}}$  is nonsingular. It is clear that  $\{\hat{\mathbf{h}}\} \subset \{\mathbf{h} \in V^* : \text{supp}(\mathbf{h}) \subseteq \mathcal{I}\}$ . The inclusion  $\{\hat{\mathbf{h}}\} \supset \{\mathbf{h} \in V^* : \text{supp}(\mathbf{h}) \subseteq \mathcal{I}\}$  can be verified by showing that  $\tilde{\mathbf{h}} := \hat{\mathbf{h}} + \mathbf{a}$  satisfies (3.8) only for  $\mathbf{a} = \mathbf{0}$ . Since  $\hat{\mathbf{h}}_{\mathcal{I}}$  is a solution of (3.8), substituting  $\tilde{\mathbf{h}} := \hat{\mathbf{h}} + \mathbf{a}$  into (3.8) yields

$$\begin{aligned} \mathbf{U}\mathbf{U}_{\mathcal{I}}^T \mathbf{a} = \mathbf{0} &\Leftrightarrow \mathbf{U}_{\mathcal{I}}^T \mathbf{a} \in \mathcal{K}(\mathbf{U}) \subset \mathcal{K}(\mathbf{U}_{\mathcal{I}}) \\ &\Leftrightarrow \mathbf{U}_{\mathcal{I}}^T \mathbf{a} = \mathbf{0}. \end{aligned} \quad (3.9)$$

The nonsingularity of  $\mathbf{G}_{\mathcal{I}}$  verifies that (3.8) holds for  $\mathbf{a} = \mathbf{0}$ .

**Proof of Lemma 3.2(b):** It is shown that (C3)  $\Rightarrow$  (C2). Assume that  $\hat{\mathbf{h}}$  is sparsest over  $V^*$ . Suppose that there exists another vector  $(\hat{\mathbf{h}} \neq) \tilde{\mathbf{h}} \in V^*$

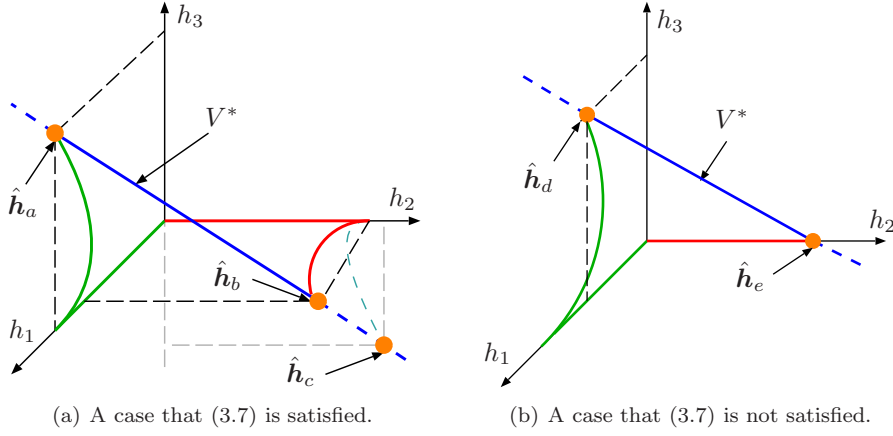


Figure 3.6: Paths of critical points connect the origin and those OLS solutions satisfying (C2).

satisfying  $\text{supp}(\tilde{\mathbf{h}}) = \text{supp}(\hat{\mathbf{h}})$ . Then, there exists a sparser vector  $\alpha\hat{\mathbf{h}} + (1 - \alpha)\tilde{\mathbf{h}} \in V^*$  for some  $\alpha \in \mathbb{R}$  than  $\hat{\mathbf{h}}$ . This contradicts the assumption and hence (C2) follows.

**Proof of Lemma 3.2(c):** (C3)  $\Rightarrow$  (C2) has already been proved in the proof of Lemma 3.2(b) and the equivalence between (C1) and (C2) has also verified by Lemma 3.2(a). It is thus sufficient to show here that (C1)  $\Rightarrow$  (C3) for completion of the proof. Assume that  $\hat{\mathbf{h}} \in V^*$  does not satisfy (C3); i.e.,  $\hat{\mathbf{h}} \notin \text{argmin}_{\mathbf{h} \in V^*} \|\mathbf{h}\|_0$ . Then,  $\|\hat{\mathbf{h}}\|_0 > r = \min_{\mathbf{h} \in V^*} \|\mathbf{h}\|_0$  by the assumption (3.7). This implies that  $\mathbf{U}_{\mathcal{I}}$  is row rank deficient and hence  $\mathbf{G}_{\mathcal{I}} = \mathbf{U}_{\mathcal{I}}\mathbf{U}_{\mathcal{I}}^{\top}$  is singular. This completes the proof.  $\square$

It must be pointed out that the assumption (3.7) is sufficient (but not necessary) to have the equivalence among (C1) – (C3). To elaborate the role of the assumption better, let us consider the three-dimensional examples described in Figure 3.6. Figure 3.6(a) illustrates the case that  $V^*$  is a one-dimensional linear manifold which intersects each of the coordinate planes  $h_1$ - $h_2$ ,  $h_2$ - $h_3$ , and  $h_3$ - $h_1$  at a single point out of the  $h_1$ ,  $h_2$ , and  $h_3$  coordinates. In this case, one can readily verify that (3.7) is satisfied; note that  $2 = \min_{\mathbf{h} \in V^*} \|\mathbf{h}\|_0 = \text{Rank}(\mathbf{U}) = 2$ . Looking at the points  $\hat{\mathbf{h}}_a$ ,  $\hat{\mathbf{h}}_b$ , and  $\hat{\mathbf{h}}_c$  in the figure, one can find that the equivalence between (C2) and (C3) holds in this case. Now, turn to Figure 3.6(b), which illustrates the case that  $V^*$  intersects each of the coordinate planes  $h_1$ - $h_2$ ,  $h_2$ - $h_3$ , and  $h_3$ - $h_1$  at a single point as the previous case but the intersection point  $\hat{\mathbf{h}}_e$  is on the  $h_2$  coordinate. In this case, (3.7) does not hold since  $1 = \min_{\mathbf{h} \in V^*} \|\mathbf{h}\|_0 \neq \text{Rank}(\mathbf{U}) = 2$ . The point  $\hat{\mathbf{h}}_e$  indicates that (C3) implies (C2), while the point  $\hat{\mathbf{h}}_d$  indicates that (C2) does not implies (C3).

The following theorem is more general than Theorem 3.2 in the sense that it guarantees the existence of a critical-point path not only for sparsest OLS solutions but also for some other ones.

**Theorem 3.3.** *For any  $\hat{\mathbf{h}} \in V^*$  satisfying (C1), there exists a continuous path of critical points connecting the origin and  $\hat{\mathbf{h}}$ .*

Proof: Consider critical points in the subspace in which  $h_i = 0, \forall i \in \bar{\mathcal{I}}$ . Define

$\varphi_{\mathcal{I}}(\mathbf{h}_{\mathcal{I}}) := \frac{1}{2}\|\mathbf{U}_{\mathcal{I}}^{\top}\mathbf{h}_{\mathcal{I}} - \mathbf{d}\|_2^2$  and  $F_{p,\mathcal{I}}(\mathbf{h}_{\mathcal{I}}) := \frac{1}{p}\|\mathbf{h}_{\mathcal{I}}\|_p^p$ . Then, it holds that

$$\begin{aligned}\varphi(\mathbf{h}) &= \varphi_{\mathcal{I}}(\mathbf{h}_{\mathcal{I}}), & F_p(\mathbf{h}) &= F_{p,\mathcal{I}}(\mathbf{h}_{\mathcal{I}}) \\ \nabla_{\mathcal{I}}\varphi(\mathbf{h}) &= \nabla\varphi_{\mathcal{I}}(\mathbf{h}_{\mathcal{I}}), & \nabla_{\mathcal{I}}F_p(\mathbf{h}) &= \nabla F_{p,\mathcal{I}}(\mathbf{h}_{\mathcal{I}})\end{aligned}$$

in the subspace. The critical point condition is therefore reduced to

$$\nabla\varphi_{\mathcal{I}}(\mathbf{h}_{\mathcal{I}}) + \lambda\nabla F_{p,\mathcal{I}}(\mathbf{h}_{\mathcal{I}}) = \mathbf{0}. \quad (3.10)$$

The critical points in the subspace can be expressed as those for the following problem:

$$\text{find } \mathbf{h}_{\lambda,\mathcal{I}}^* \in \underset{\mathbf{h}_{\mathcal{I}} \in \mathbb{R}^{|\mathcal{I}|}}{\text{argmin}} f_{\lambda,\mathcal{I}}(\mathbf{h}) := \varphi_{\mathcal{I}}(\mathbf{h}_{\mathcal{I}}) + \lambda F_{p,\mathcal{I}}(\mathbf{h}_{\mathcal{I}}). \quad (3.11)$$

Due to the nonsingularity of  $\mathbf{G}_{\mathcal{I}}$ , Theorem 3.1 can be applied to the problem in (3.11) to verify the existence of a continuous path of critical points connecting  $\hat{\mathbf{h}}$  and the origin.  $\square$

Lemma 3.2(b) ensures that every sparsest solution satisfies (C2), and hence Theorem 3.3 directly implies Theorem 3.2. In Figure 3.6, one can see that all the points  $\hat{\mathbf{h}}_a, \hat{\mathbf{h}}_b, \hat{\mathbf{h}}_c, \hat{\mathbf{h}}_d$ , and  $\hat{\mathbf{h}}_e$  that satisfy (C2) are connected with the origin by paths of critical points. This illustrates the validity of Theorem 3.3 under Lemma 3.2.

### 3.4 Paths from other OLS Solutions

In this section, the possibilities of critical-point paths emerging from some of the OLS solutions, which have not been considered in the previous section, are discussed. In fact, being sparsest (or having a sufficiently small support) is a sufficient, but not necessary, condition for an OLS solution to be connected with the origin by a path. The following proposition gives a necessary condition for an OLS solution to have a path emerging from the solution itself.

**Proposition 3.3.** *Let  $\hat{\mathbf{h}} \in V^*$  with its support  $\mathcal{I}$ . Then, in a vicinity of  $\hat{\mathbf{h}}$ , there exists a path of critical points with support  $\mathcal{I}$  emerging from  $\hat{\mathbf{h}}$  only if*

$$\nabla_{\mathcal{I}}F_p(\hat{\mathbf{h}}) \in \mathcal{R}(\mathbf{U}_{\mathcal{I}}). \quad (3.12)$$

Here, a path such that  $\lambda$  is continuously zero (in part) is neglected.

Proof: Suppose that there exists a path  $\mathcal{C}$  of critical points with support  $\mathcal{I}$  leading to  $\hat{\mathbf{h}}$ . Then, define a sequence  $(\mathbf{h}_n)_{n \in \mathbb{N}} \subset \mathcal{C}$  converging to  $\hat{\mathbf{h}}$ . Since  $\lambda = 0$  at  $\hat{\mathbf{h}} \in V^*$ ,  $\lambda \neq 0$  in the vicinity of  $\hat{\mathbf{h}}$  and hence there exists a sufficiently large  $N_0 \in \mathbb{N}$  such that  $\nabla_{\mathcal{I}}F_p(\mathbf{h}_n)$  is aligned to  $\nabla_{\mathcal{I}}\varphi(\mathbf{h}_n)$  for all  $n \geq N_0$ . This implies that  $(\nabla_{\mathcal{I}}F_p(\mathbf{h}_n))_{n \geq N_0} \subset \mathcal{R}(\mathbf{U}_{\mathcal{I}})$ . The closedness of the subspace  $\mathcal{R}(\mathbf{U}_{\mathcal{I}})$  and the continuity of the operator  $\nabla_{\mathcal{I}}F_p$  ensure  $\nabla_{\mathcal{I}}F_p(\hat{\mathbf{h}}) \in \mathcal{R}(\mathbf{U}_{\mathcal{I}})$ .  $\square$

The existence of an OLS solution  $\hat{\mathbf{h}} \in V^*$  with full support which satisfies the condition in Proposition 3.3 can be considered as follows. Without loss of generality, suppose that  $V^* \cap \mathcal{R}(\mathbf{U}) \subset \mathbb{R}_{++}^n := \{\mathbf{h} \in \mathbb{R}^m : h_i > 0, \forall i\}$ ;  $V^* \cap \mathcal{R}(\mathbf{U})$  contains a sole vector which is the minimum Euclidean-norm OLS solution. Define dual coordinates by  $\boldsymbol{\eta} := \nabla F_p(\mathbf{h})$  for  $\mathbf{h} \in \mathbb{R}_{++}^n$ . It is then not difficult to see that the correspondence between  $\mathbf{h}$  and  $\boldsymbol{\eta}$  is one to one over  $\mathbb{R}_{++}^n$ . Consider

the special case of  $\text{Rank}(\mathbf{U}) = 1$ . In this case,  $\mathcal{R}(\mathbf{U})$  is a one-dimensional subspace and  $V^*$  is a hyperplane in  $\mathbb{R}^m$ . One can see that the hyperplane  $V^*$  separates the two convex sets in  $\mathbb{R}_{++}^n$ :  $S_1 := (0, \epsilon)^n$  and  $S_2 := (\delta, \infty)^n$  for some small  $\epsilon > 0$  and large  $\delta > 0$ . The subspace  $\mathcal{R}(\mathbf{U})$  contains some  $\boldsymbol{\eta}_1 \in (1/\sqrt{\epsilon}, \infty)^n$  and  $\boldsymbol{\eta}_2 \in (0, 1/\sqrt{\delta})^n$ . Then the inverse images of  $\boldsymbol{\eta}_1$  and  $\boldsymbol{\eta}_2$  under the nonlinear mapping  $\nabla F_p$  satisfy  $\mathbf{h}_1 \in (0, \epsilon)^n$  and  $\mathbf{h}_2 \in (\delta, \infty)^n$ , respectively. This implies that the inverse image of the ray  $\mathcal{R}(\mathbf{U}) \cap \mathbb{R}_{++}^n$  contains a continuous curve connecting the two points  $\mathbf{h}_1$  and  $\mathbf{h}_2$  which are separated by the hyperplane  $V^*$ . Hence, the curve crosses  $V^*$ , meaning that there exists an OLS solution with full support satisfying (3.12).

The same discussion can be also applied to the case that  $\text{Rank}(\mathbf{U}) = n - 1$ . However, the above discussion gives us doubts that, in the case that  $1 < \text{Rank}(\mathbf{U}) < n - 1$ , there exists in general no full-support OLS solution satisfying (3.12). In this case, one can extend the discussion for  $\text{Rank}(\mathbf{U}) = n - 1$  by restricting oneself to an  $r + 1$  subspace

$$M_{\mathcal{I}} := \{\mathbf{h} : \text{supp}(\mathbf{h}) \subseteq \mathcal{I}\} \quad (3.13)$$

for any  $\mathcal{I} \subset \{1, 2, \dots, n\}$  with  $|\mathcal{I}| = r + 1$  so that  $V_{\mathcal{I}}^* := \mathbf{h}_{\mathcal{I}} \in \mathbb{R}^{r+1} : [\mathbf{h}_{\mathcal{I}}^T \mathbf{0}^T] \in V^*$  is a one-dimensional line and  $\mathcal{R}(\mathbf{U}_{\mathcal{I}}) \subset \mathbb{R}^{r+1}$  is a hyperplane in  $\mathbb{R}^{r+1}$ .

### 3.5 Concluding Remarks

In this chapter, critical points of a least squares problem under  $\ell_p$ -regularization ( $0 < p < 1$ ) have been studied. It has been proved under a mild condition that there exists a critical-point path for every sparsest OLS solution which connects it with the origin. The major contributions of this work are summarized as follows.

- In an overdetermined case, there is a path connecting the unique OLS solution and the origin (Theorem 3.1).
- The path connecting the OLS solution and the origin consists of curves each of which connects critical points for  $\lambda = 0$ , hence includes a point where the Hessian matrix  $\mathbf{K}_{\mathcal{I}}$  is singular (Proposition 3.2, Remark 3.2).
- Every sparsest OLS solution in an underdetermined case is connected with the origin by a continuous path of the critical points (Theorem 3.2).
- In an underdetermined case, there is a path connecting an OLS solution and the origin if the solution can be regarded as that of an overdetermined case (Theorem 3.3).
- An OLS solution has a path of critical points emerging from itself only if  $\nabla_{\mathcal{I}} F_p(\mathbf{h})$  lies in the range of  $\mathbf{G}_{\mathcal{I}}$  (Proposition 3.3).

It has turned out that  $\ell_p$ -regularization ( $0 < p < 1$ ) does yield sparsest solutions of a least squares problem. The existence theorem presented in this paper will become a driving force for developing an efficient method to compute a sparsest solution by using  $\ell_p$ -regularization.

## Chapter 4

# Kernel Weight Design for Multikernel Adaptive Filtering

### 4.1 Introduction

How can multiple kernel functions cooperate well for estimating an unknown nonlinear function of interest in online fashion? This is a highly important question for the signal processing community because many tasks in signal processing and machine learning can be cast as online nonlinear estimation problems. Reproducing kernel has widely been studied due to a number of reasons including the so-called kernel trick to reduce an inner product computation in a possibly infinite dimensional space to finite-dimensional arithmetics, as well as the solid theoretical basis. Kernel adaptive filtering [26–28, 30, 33, 37, 38, 58–63] has become a popular approach to online nonlinear estimation during the last decade. The success of kernel adaptive filtering, however, relies heavily on the choice of kernel, and it is difficult in practice to design a kernel fitting a given specific task well.

In the machine learning community, MKL [39–44] has extensively been studied mainly under batch settings motivated by the difficulty of kernel design, and online algorithms have also been proposed recently [43, 44]. Most MKL algorithms consist of two steps. In the first step, an appropriate kernel which fits the observed data well is constructed by a convex combination of multiple kernels, which is then used to estimate the unknown nonlinear function. The second step of MKL is essentially a coefficient update with a single kernel, and hence MKL has a severe limitation in estimating a target function which consists of multiple components such as linear and nonlinear functions. In the signal processing community, on the other hand, multikernel adaptive filtering [47–49] has been studied as an effective method utilizing multiple reproducing kernels under online settings. See [45, 46, 64–69] for other approaches using multiple kernels. Multikernel adaptive filtering has a higher degree of freedom in estimation than MKL. Thanks to the higher degree of freedom, multikernel adaptive filtering enjoys (a) a more compact representation than MKL, and (b) a potential ability



to extract the local structures of the target containing multiple components.

Appropriate treatments of kernel weights are necessary for multikernel adaptive filtering since a naive use of unweighted kernels never guarantees desirable performance when there is some imbalance among the kernels. To use an analogy, suppose the following situation in linear adaptive filtering. When the power of the input vector is concentrated on a small subset of entries, the input autocorrelation matrix has large eigenvalue spread. It has been known that large eigenvalue spread of the input autocorrelation matrix causes unbalanced convergence rates among the eigenvectors, and hence deteriorates the performance of adaptive filtering [24, 25]. The eigenvalue spread can be reduced by imposing weights on the entries so that the powers of weighted entries become equal to each other. The only difference in the case of multikernel adaptive filtering from the linear case is that each kernel weight is not imposed on a single entry of the input vector, but on all the entries associated with the corresponding kernel. The power balance and convergence rates are then considered in a *kernel-wise* sense. A principled way of tuning the kernel weights, which lead to uniform convergence rates of the kernel-wise coefficient errors, for multikernel adaptive filtering has not been investigated so far. Such hyper-parameter tuning problems have been empirically dealt with under some pretraining or validation process, which is undesirable in online settings. To exploit the full potential of multikernel adaptive filtering, it is therefore essential to develop a technique which yields the appropriate kernel weights for an online manner.

The headline question presented at the beginning can be subdivided as follows: (i) Is there a set of kernel weights which equalizes the convergence rates of the kernel-wise coefficient errors? (ii) If such a set of weights exists, how can it be implemented in online fashion? In this chapter, an answer to the above questions is presented by deriving a recursive kernel weighting technique which updates the kernel weights at each iteration, instead of presetting the weights.

## 4.2 Transient Behavior and Weight Derivation

### 4.2.1 Preliminaries to Transient Analysis

In advance of discussions on transient analysis, some preliminary informations related to this chapter must be reminded. The nonlinear adaptive filtering model is given as follows.

$$d_n = \psi(\mathbf{u}_n) + \nu_n \in \mathbb{R}, \quad n \in \mathbb{N}, \quad (4.1)$$

where  $\mathbf{u}_n \in \mathbb{R}^L$  is the input vector,  $d_n$  is the corresponding output generated by the unknown function  $\psi$  of interest, and  $\nu_n$  is the zero-mean additive noise which is assumed independent of  $\mathbf{u}_n$ . Given  $Q$  kernel functions  $\kappa_q(\cdot, \cdot)$ ,  $q = 1, 2, \dots, Q$ , and the corresponding dictionaries  $\mathcal{D}_{q,n} := \{\kappa_q(\mathbf{x}_1^{(q)}, \cdot), \dots, \kappa_q(\mathbf{x}_{r_{q,n}}^{(q)}, \cdot)\}$ , whose cardinality is  $r_{q,n}$ , multikernel adaptive filtering achieves the following estimate of the unknown function  $\psi$ .

$$\hat{\varphi}_n(\mathbf{u}_n) = \sum_{q=1}^Q \sum_{j=1}^{r_{q,n}} h_{n,j}^{(q)} \kappa_q(\mathbf{x}_j^{(q)}, \mathbf{u}_n) = \sum_{q=1}^Q \mathbf{k}_n^{(q)\top} \mathbf{h}_n^{(q)} = \mathbf{k}_n^\top \mathbf{h}_n, \quad (4.2)$$

where

$$\begin{aligned}\mathbf{k}_n^{(q)} &= [\kappa_q(\mathbf{x}_1^{(q)}, \mathbf{u}_n), \dots, \kappa_q(\mathbf{x}_{r_{q,n}}^{(q)}, \mathbf{u}_n)]^\top \in \mathbb{R}^{r_{q,n}}, \\ \mathbf{h}_n^{(q)} &= [h_{n,1}^{(q)}, \dots, h_{n,r_{q,n}}^{(q)}]^\top \in \mathbb{R}^{r_{q,n}},\end{aligned}$$

for all  $q = 1, 2, \dots, Q$ , and

$$\begin{aligned}\mathbf{k}_n &= [\mathbf{k}_n^{(1)\top}, \dots, \mathbf{k}_n^{(Q)\top}]^\top \in \mathbb{R}^{r_n}, \\ \mathbf{h}_n &= [\mathbf{h}_n^{(1)\top}, \dots, \mathbf{h}_n^{(Q)\top}]^\top \in \mathbb{R}^{r_n}\end{aligned}$$

are the concatenated vectors, where  $r_n = \sum_{q=1}^Q r_{q,n}$ . The nonlinear estimation of  $\psi$  then reduces to linear estimation of the Euclidean vector  $\mathbf{h}^*$ , which is a minimizer of some error function. The MKNLMS algorithm [47, 48] is a multikernel adaptive filtering algorithm of which the error function is the mean squared error  $\text{MSE}(\mathbf{h}) := E[(\mathbf{k}_n^\top \mathbf{h} - d_n)^2]$ , where  $\mathbf{h} \in \mathbb{R}^{r_n}$  and  $E[\cdot]$  stands for the expectation with respect to the input  $\mathbf{u}_n$  and the noise  $\nu_n$ . Without a manipulation of dictionaries, each update of the MKNLMS algorithm can be written as follows:

$$\mathbf{h}_{n+1} = \mathbf{h}_n - \mu \frac{\mathbf{k}_n^\top \mathbf{h}_n - d_n}{\|\mathbf{k}_n\|_2^2} \mathbf{k}_n, \quad (4.3)$$

Applying the kernel weights to the update (4.3) yields

$$\mathbf{g}_{n+1} = \mathbf{g}_n - \mu \frac{\mathbf{k}_n^\top \mathbf{W} \mathbf{g}_n - d_n}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \mathbf{W} \mathbf{k}_n, \quad (4.4)$$

where  $\mathbf{g}_n := \mathbf{W}^{-1} \mathbf{h}_n$  and the weight matrix  $\mathbf{W}$  is given as

$$\mathbf{W} := \begin{bmatrix} w_1 \mathbf{I}_{r_{1,n}} & \mathbf{O} & \dots & \mathbf{O} \\ \mathbf{O} & w_2 \mathbf{I}_{r_{2,n}} & \dots & \mathbf{O} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{O} & \mathbf{O} & \dots & w_Q \mathbf{I}_{r_{Q,n}} \end{bmatrix}. \quad (4.5)$$

On the other hand, the coefficient error  $\mathbf{z}_n := \mathbf{g}_n - \mathbf{g}^*$ , where  $\mathbf{g}^* := \mathbf{W}^{-1} \mathbf{h}^*$ , is updated by equation (4.4) as follows:

$$\mathbf{z}_{n+1} = \mathbf{z}_n - \mu \frac{\mathbf{k}_n^\top \mathbf{W} \mathbf{z}_n - \nu_n}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \mathbf{W} \mathbf{k}_n. \quad (4.6)$$

For the  $q$ th block, the blockwise update is given as

$$\mathbf{z}_{n+1}^{(q)} = \mathbf{z}_n^{(q)} - \mu w_q \frac{\mathbf{k}_n^\top \mathbf{W} \mathbf{z}_n - \nu_n}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \mathbf{k}_n^{(q)}, \quad (4.7)$$

$$\begin{aligned}\|\mathbf{z}_{n+1}^{(q)}\|_2^2 &= \|\mathbf{z}_n^{(q)}\|_2^2 - 2\mu w_q \frac{\mathbf{k}_n^\top \mathbf{W} \mathbf{z}_n - \nu_n}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \mathbf{k}_n^{(q)\top} \mathbf{z}_n^{(q)} \\ &\quad + (\mu w_q)^2 \frac{\|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \frac{(\mathbf{k}_n^\top \mathbf{W} \mathbf{z}_n - \nu_n)^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2}.\end{aligned} \quad (4.8)$$

It is clear that the transient behavior of the squared norm (4.8) depends on the choice of the kernel weights. It is now to focus on the convergence rates of  $\|\mathbf{z}_{n+1}^{(q)}\|_2^2$ ,  $q = 1, 2, \dots, Q$ , and find the kernel weights which lead to the uniform convergence rate, *i.e.*,  $E[\|\mathbf{z}_{n+1}^{(q)}\|_2^2] = \rho E[\|\mathbf{z}_n^{(q)}\|_2^2]$  with the same  $\rho \in (0, 1)$  for all  $q$ . For simplicity, the transient analysis of the coefficient errors is considered in the case that the dictionaries are fixed ( $r_{q,n}$  is fixed as well for all  $q$ ). In addition to fixing the dictionaries, the following properties are also assumed.

**Assumption 4.1.**

- (a) At each time  $n$ , the kernelized input  $\mathbf{k}_n$  and the additive noise  $\nu_n$  are independent of each other.
- (b) The fluctuations in the squared norms  $\|\mathbf{k}_n^{(q)}\|_2^2$ ,  $q = 1, 2, \dots, Q$ , and  $\|\mathbf{W}\mathbf{k}_n\|_2^2$  are sufficiently small to justify the approximations

$$E\left[\frac{\|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2}\right] \approx \frac{E[\|\mathbf{k}_n^{(q)}\|_2^2]}{E[\|\mathbf{W}\mathbf{k}_n\|_2^2]},$$

where the matrix  $\mathbf{W}$  is given as (4.5).

The above assumptions are commonly used for performance analysis of the NLMS algorithm [24, 25].

The goal of this chapter is to derive a set of weights that equalizes the expected convergence rates of the kernel-wise subvectors  $\mathbf{z}_n^{(q)}$ ,  $q = 1, 2, \dots, Q$ , of the coefficient error vector  $\mathbf{z}_n$  in the sense of the squared norm. The coefficient error vector  $\mathbf{z}_n$  itself, however, is unknown due to its dependency on the unknown target  $\mathbf{g}^*$ . To tackle with this issue, we assume  $\mathbf{z}_n$  as a random variable which has the following properties:

**Assumption 4.2** (Assumption on  $\mathbf{z}_n$ ).

- (a) The coefficient error vector  $\mathbf{z}_n$  is distributed spherically with zero mean.
- (b) The distribution of  $\mathbf{z}_n$  is nearly even in the larger area than the area where the weighted input  $\mathbf{W}\mathbf{k}_n$  can exist, so that  $\mathbf{z}_n$  is approximately independent of  $\mathbf{W}\mathbf{k}_n$ .

The above assumption implies the following properties:

- The subvectors  $\mathbf{z}_n^{(q)}$ ,  $q = 1, 2, \dots, Q$  are independent of each other.
- The probability of  $\mathbf{z}_n$  depends only on its squared norm  $\|\mathbf{z}_n\|_2^2$ .

### 4.2.2 Derivation of Kernel Weights Equalizing Convergence Rates

Using the expectation  $E_{\mathbf{z}_n}[\|\mathbf{z}_n^{(q)}\|_2^2]$  with respect to  $\mathbf{z}_n$ , the transient analysis on (4.8) becomes tractable even if  $\mathbf{z}_n^{(q)}$  itself is unknown. The expectation of

(4.8) with respect to  $\mathbf{z}_n$  is

$$\begin{aligned}
E_{\mathbf{z}_n} [\|\mathbf{z}_{n+1}^{(q)}\|_2^2] &= E_{\mathbf{z}_n} [\|\mathbf{z}_n^{(q)}\|_2^2] - 2\mu w_q E_{\mathbf{z}_n} \left[ \frac{\mathbf{z}_n^{(q)\top} \mathbf{k}_n^{(q)} \mathbf{k}_n^\top \mathbf{W} \mathbf{z}_n}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \right] \\
&\quad + (\mu w_q)^2 \frac{\|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2} E_{\mathbf{z}_n} \left[ \frac{(\mathbf{k}_n^\top \mathbf{W} \mathbf{z}_n)^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \right] \\
&\quad + (\mu w_q)^2 \frac{\|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \frac{\nu_n^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2}. \tag{4.9}
\end{aligned}$$

The quadratic terms in (4.9) are calculated with the following lemma.

**Lemma 4.1.** *Let  $\mathbf{A} \in \mathbb{R}^{r_n \times r_n}$  be an arbitrary matrix. If  $\mathbf{z}_n$  is spherically distributed with zero mean, it holds that*

$$E_{\mathbf{z}_n} [\mathbf{z}_n^\top \mathbf{A} \mathbf{z}_n] = \frac{E_{\mathbf{z}_n} [\|\mathbf{z}_n\|_2^2]}{r_n} \text{Trace}(\mathbf{A}). \tag{4.10}$$

Proof: It is trivial that

$$\begin{aligned}
E_{\mathbf{z}_n} [\mathbf{z}_n^\top \mathbf{A} \mathbf{z}_n] &= \text{Trace}(E_{\mathbf{z}_n} [\mathbf{z}_n \mathbf{z}_n^\top \mathbf{A}]) \\
&= \text{Trace}(E_{\mathbf{z}_n} [\mathbf{z}_n \mathbf{z}_n^\top] \mathbf{A}). \tag{4.11}
\end{aligned}$$

Since  $\mathbf{z}_n$  is spherically distributed with zero mean, it holds that

$$E_{\mathbf{z}_n} [\mathbf{z}_n \mathbf{z}_n^\top] = \lambda_n \mathbf{I}_{r_n}, \tag{4.12}$$

for some  $\lambda_n \geq 0$ . From (4.12), the variance  $\lambda_n$  can be calculated as

$$\lambda_n = \frac{\text{Trace}(E_{\mathbf{z}_n} [\mathbf{z}_n \mathbf{z}_n^\top])}{\text{Trace}(\mathbf{I}_{r_n})} = \frac{E_{\mathbf{z}_n} [\|\mathbf{z}_n\|_2^2]}{r_n}. \tag{4.13}$$

Substituting both of (4.12) and (4.13) into (4.11) leads to

$$\begin{aligned}
E_{\mathbf{z}_n} [\mathbf{z}_n^\top \mathbf{A} \mathbf{z}_n] &= \text{Trace}(E_{\mathbf{z}_n} [\mathbf{z}_n \mathbf{z}_n^\top] \mathbf{A}) \\
&= \text{Trace}(\lambda_n \mathbf{I}_{r_n} \mathbf{A}) \\
&= \lambda_n \text{Trace}(\mathbf{A}) \\
&= \frac{E_{\mathbf{z}_n} [\|\mathbf{z}_n\|_2^2]}{r_n} \text{Trace}(\mathbf{A}), \tag{4.14}
\end{aligned}$$

which completes the proof.  $\square$

By Lemma 4.1 and Assumption 4.2(b), (4.9) leads to

$$\begin{aligned}
E_{\mathbf{z}_n} [\|\mathbf{z}_{n+1}^{(q)}\|_2^2] &\approx E_{\mathbf{z}_n} [\|\mathbf{z}_n^{(q)}\|_2^2] - 2 \frac{\mu w_q^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \|\mathbf{k}_n^{(q)}\|_2^2 \frac{E_{\mathbf{z}_n} [\|\mathbf{z}_n^{(q)}\|_2^2]}{r_{q,n}} \\
&\quad + \frac{(\mu w_q)^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \|\mathbf{k}_n^{(q)}\|_2^2 \left( \frac{E_{\mathbf{z}_n} [\|\mathbf{z}_n\|_2^2]}{r_n} + \frac{\nu_n^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \right) \\
&= \left( 1 - \frac{\mu(2-\mu) w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2 r_{q,n}} \right) E_{\mathbf{z}_n} [\|\mathbf{z}_n^{(q)}\|_2^2] \\
&\quad + (\mu w_q)^2 \frac{\|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2} \frac{\nu_n^2}{\|\mathbf{W} \mathbf{k}_n\|_2^2}, \tag{4.15}
\end{aligned}$$

where the second equality holds by

$$E_{\mathbf{z}_n} [\|\mathbf{z}_n^{(q)}\|_2^2] = \frac{r_{q,n}}{r_n} E_{\mathbf{z}_n} [\|\mathbf{z}_n\|_2^2], \quad (4.16)$$

which is readily verified by Lemma 4.1. The expectation of (4.15) with respect to  $\mathbf{k}_n$  and  $\nu_n$  is

$$\begin{aligned} E[E_{\mathbf{z}_n} [\|\mathbf{z}_{n+1}^{(q)}\|_2^2]] &= E \left[ 1 - \frac{\mu(2-\mu)}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{r_{q,n}} \right] E_{\mathbf{z}_n} [\|\mathbf{z}_n^{(q)}\|_2^2] \\ &\quad + (\mu w_q)^2 E \left[ \frac{\|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \frac{\nu_n^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right]. \end{aligned} \quad (4.17)$$

Since  $\mu \in (0, 2)$ , it holds that

$$0 < \frac{\mu(2-\mu)}{r_{q,n}} \leq \frac{1}{r_{q,n}}. \quad (4.18)$$

In addition, by  $w_q > 0$ ,  $q = 1, 2, \dots, Q$ , it holds that

$$0 < \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} < 1. \quad (4.19)$$

From (4.18) and (4.19), it follows that

$$E \left[ 1 - \frac{\mu(2-\mu)}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{r_{q,n}} \right] = 1 - \frac{\mu(2-\mu)}{r_{q,n}} E \left[ \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right] \in (0, 1). \quad (4.20)$$

This implies that the expected squared norms  $E_{\mathbf{z}_n} [\|\mathbf{z}_n^{(q)}\|_2^2]$  of the kernel-wise coefficient errors decrease monotonically after each update (4.4). Then the expected squared norms converge to zero and (4.20) can be regarded as the convergence rate. If a uniform convergence rate, say  $\rho > 0$ , is achieved by some kernel weight setting, then it holds that

$$E \left[ 1 - \frac{\mu(2-\mu)}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{r_{q,n}} \right] = \rho, \quad (4.21)$$

for all  $q = 1, 2, \dots, Q$ , as shown in the following proposition.

**Proposition 4.1.** *Setting the kernel weights as*

$$w_q = \gamma \left( \frac{E[\|\mathbf{k}_n^{(q)}\|_2^2]}{r_{q,n}} \right)^{-\frac{1}{2}}, \quad \gamma > 0, \quad q = 1, 2, \dots, Q, \quad (4.22)$$

*achieves the uniform convergence rate*

$$\rho = 1 - \frac{\mu(2-\mu)}{r_n}, \quad (4.23)$$

*for all  $q = 1, 2, \dots, Q$ .*

Proof: From Assumption 4.1(b) and (4.22),

$$\begin{aligned}
E \left[ \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right] &\approx \frac{w_q^2 E[\|\mathbf{k}_n^{(q)}\|_2^2]}{E[\|\mathbf{W}\mathbf{k}_n\|_2^2]} \\
&= \frac{w_q^2 E[\|\mathbf{k}_n^{(q)}\|_2^2]}{\sum_{p=1}^Q w_p^2 E[\|\mathbf{k}_n^{(p)}\|_2^2]} \\
&= \frac{r_{q,n} \gamma^2}{\sum_{p=1}^Q r_{p,n} \gamma^2} = \frac{r_{q,n}}{r_n}, \tag{4.24}
\end{aligned}$$

for all  $q$ . This leads to

$$\begin{aligned}
E \left[ 1 - \frac{\mu(2-\mu)}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{r_{q,n}} \right] &= 1 - \frac{\mu(2-\mu)}{r_{q,n}} E \left[ \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right] \\
&= 1 - \frac{\mu(2-\mu)}{r_{q,n}} \frac{r_{q,n}}{r_n} \\
&= 1 - \frac{\mu(2-\mu)}{r_n}, \tag{4.25}
\end{aligned}$$

which does not depend on  $q$ .  $\square$

### 4.2.3 Steady-state Error

The summation (4.17) with respect to  $q$  yields the following transient behavior of the entire error.

$$\begin{aligned}
E[E_{z_n}[\|z_{n+1}\|_2^2]] &= \sum_{q=1}^Q E[E_{z_n}[\|z_{n+1}^{(q)}\|_2^2]] \\
&= E_{z_n}[\|z_n\|_2^2] - \frac{\mu(2-\mu)}{r_n} E_{z_n}[\|z_n\|_2^2] \sum_{q=1}^Q E \left[ \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right] \\
&\quad + \mu^2 \sum_{q=1}^Q E \left[ \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \frac{\nu_n^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right] \\
&= \left( 1 - \frac{\mu(2-\mu)}{r_n} \right) E_{z_n}[\|z_n\|_2^2] + \mu^2 E \left[ \frac{\nu_n^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right], \tag{4.26}
\end{aligned}$$

where the second equality is due to (4.16), and the last equality holds because

$$\sum_{q=1}^Q E \left[ \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right] = E \left[ \sum_{q=1}^Q \frac{w_q^2 \|\mathbf{k}_n^{(q)}\|_2^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right] = 1. \tag{4.27}$$

From (4.26), it holds that

$$\begin{aligned}
E_{z_n}[\|z_n\|_2^2] &= \frac{r_n}{\mu(2-\mu)} \mu^2 E \left[ \frac{\nu_n^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right] \\
&= \frac{\mu r_n}{2-\mu} E \left[ \frac{\nu_n^2}{\|\mathbf{W}\mathbf{k}_n\|_2^2} \right] \tag{4.28}
\end{aligned}$$

at the steady state. The above steady-state error depends on the choice of  $\mathbf{W}$ . In particular, if the weights  $w_q$ ,  $q = 1, 2, \dots, Q$ , become larger, the steady-state error becomes smaller even if the ratios among the weights are fixed. By considering the normalized error

$$\frac{E_{\mathbf{z}_n} [\|\mathbf{z}_n\|_2^2]}{E_{\mathbf{z}_n} [\|\mathbf{g}^*\|_2^2]} = \frac{E_{\mathbf{z}_n} [\|\mathbf{z}_n\|_2^2]}{E_{\mathbf{z}_n} [\|\mathbf{z}_n\|_2^2] + \|\mathbf{g}_n\|_2^2}, \quad (4.29)$$

such dependency on the amplitudes is eliminated. In addition, both of (4.23) and (4.29) are independent of the value of  $\gamma$  if the kernel weights are chosen by (4.22). Hereafter,  $\gamma = 1$  without loss of generality.

**Remark 4.1** (Relation to Batch Normalization). *The case where  $r_{q,n} = 1$  of the kernel weight (4.22) corresponds to the use of entrywise kernel weights. Imposing the entrywise weights is similar to the batch normalization, which has been popular in machine learning community [70–72]. The batch normalization aims to balance the scales of the entries of the observed input vector. In the kernel weights (4.22),  $\frac{E[\|\mathbf{k}_n^{(q)}\|_2^2]}{r_{q,n}}$  can be seen as a scale of the  $q$ th kernel because it holds that*

$$\frac{E[\|\mathbf{k}_n^{(q)}\|_2^2]}{r_{q,n}} = E\left[\frac{\|\mathbf{k}_n^{(q)}\|_2^2}{r_{q,n}}\right] \approx \int_{\mathbb{R}^L} E[(\kappa_q(\mathbf{u}_n, \mathbf{x}))^2] d\mathbf{x}, \quad (4.30)$$

for sufficiently large  $r_{q,n}$ . The essential difference between imposing the kernel weight (4.22) and the batch normalization are twofold.

- The kernel weight does not consider the mean of the kernelized input, where the batch normalization includes the subtraction of the mean. Imposing the kernel weights is therefore equivalent to the batch normalization when the mean is zero.
- The kernel weight does not use the mini-batch in order to approximate the expectation. (In practice, the mini-batch can be used to implement the kernel weights.)

## 4.3 Implementation of the Kernel Weights

### 4.3.1 Weight Estimation from Observations

In kernel adaptive filtering, the length  $r_{q,n}$  of the  $q$ th kernelized input may change in time. In addition,  $E[\|\mathbf{k}_n^{(q)}\|_2^2]$  is unavailable because the statistical properties of  $\mathbf{k}_n^{(q)}$  are unknown. The expectation in the kernel weight (4.22) can be approximated by the sample mean

$$E\left[\frac{\|\mathbf{k}_n^{(q)}\|_2^2}{r_{q,n}}\right] \approx \frac{1}{n} \sum_{l=1}^n \frac{\|\mathbf{k}_l^{(q)}\|_2^2}{r_{q,l}}. \quad (4.31)$$

By the above approximation, the kernel weights can be updated recursively. Although an exponentially weighted average with a forgetting factor can be

used instead, the simple sample mean is solely focused in the present study. Under (4.31) with  $\gamma = 1$ , (4.22) is approximated by

$$w_{q,n} = \left( \frac{1}{n} \sum_{l=1}^n \frac{\|\mathbf{k}_l^{(q)}\|_2^2}{r_{q,l}} \right)^{-\frac{1}{2}}. \quad (4.32)$$

### 4.3.2 Implementation Based on Metric Design

Replacing the fixed weight matrix  $\mathbf{W}$  in (4.4) to  $\mathbf{W}_n$  which is time-varying, the following update is obtained.

$$\mathbf{g}_{n+1} = \mathbf{g}_n - \mu \frac{\mathbf{k}_n^\top \mathbf{W}_n \mathbf{g}_n - d_n}{\mathbf{k}_n^\top \mathbf{W}_n^2 \mathbf{k}_n} \mathbf{W}_n \mathbf{k}_n, \quad (4.33)$$

where

$$\mathbf{W}_n := \begin{bmatrix} w_{1,n} \mathbf{I}_{r_{1,n}} & \mathbf{O} & \cdots & \mathbf{O} \\ \mathbf{O} & w_{2,n} \mathbf{I}_{r_{2,n}} & \cdots & \mathbf{O} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{O} & \mathbf{O} & \cdots & w_{Q,n} \mathbf{I}_{r_{Q,n}} \end{bmatrix}. \quad (4.34)$$

Since (4.33) can be regarded as the update equation of the NLMS algorithm for the modified input vector  $\mathbf{W}_n \mathbf{k}_n$ , it is clear that the algorithm tracks the vector  $\mathbf{g}_n^* = \mathbf{W}_n^{-1} \mathbf{h}^*$ , so that the output of the filter to this modified input vector satisfies  $(\mathbf{W}_n \mathbf{k}_n)^\top \mathbf{g}_n^* = \mathbf{k}_n^\top \mathbf{h}^*$ . This, however, means that the target vector  $\mathbf{g}_n^*$  changes as  $\mathbf{W}_n$  changes, even if  $\mathbf{h}^*$  is fixed. This may cause slow convergence of the coefficients which comes from strong time dependency of the kernel weights during the initial phase of adaptation. This dependency issue can be fixed by left-multiplying the both sides of (4.33) by  $\mathbf{W}_n$ , which leads to

$$\mathbf{h}_{n+1} = \mathbf{h}_n - \mu \frac{\mathbf{k}_n^\top \mathbf{h}_n - d_n}{\mathbf{k}_n^\top \mathbf{W}_n^2 \mathbf{k}_n} \mathbf{W}_n^2 \mathbf{k}_n, \quad (4.35)$$

with  $\mathbf{h}_n = \mathbf{W}_n \mathbf{g}_n$ . The above update tracks the original target vector  $\mathbf{h}^*$ , and hence is free from the dependency of the target vector on  $\mathbf{W}_n$ . Provided that  $\mathbf{h}^*$  stays constant, it can then be guaranteed that the coefficient vector  $\mathbf{h}_n$  approaches monotonically to  $\mathbf{h}^*$  after each update in terms of the metric  $\mathbf{W}_n^{-2}$  since the update equation in (4.35) can be interpreted as a relaxed projection in terms of the metric  $\mathbf{W}_n^{-2}$ . See [73–77] for the metric projection-based adaptive filtering algorithms. The entire MKNLMS algorithm equipped with the proposed metric design is presented in Algorithm 2. As a measure for the sparsification criterion, the novelty [32], the coherence [30], or the surprise [33] can be used.

### 4.3.3 Use of Linear Kernel

For the linear kernel, a fixed dictionary  $\mathcal{D}_L := \{\mathbf{e}_1^\top(\cdot), \mathbf{e}_2^\top(\cdot), \dots, \mathbf{e}_L^\top(\cdot), \mathbf{1}\}$ , where  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_L\}$  is the standard basis of the Euclidean space  $\mathbb{R}^L$  and  $\mathbf{1} : \mathbb{R}^L \mapsto \mathbb{R} : \mathbf{x} \mapsto 1$ , can be considered since the dictionary  $\mathcal{D}_L$  completes a basis of an  $L + 1$  dimensional Euclidean space. From (2.11), such a use of the linear kernel is actually the simultaneous use of two linear kernel functions



**Algorithm 2** MKNLMS with Adaptive Kernel Weighting

- 
- 0)  $\mathcal{D}_{q,n} = \emptyset$ ,  $\mathbf{h}_0^{(q)} = [\ ]$ ,  $q = 1, 2, \dots, Q$ ,  $\mu \in (0, 2)$ : step size.
- 1) Based on a sparsification criterion, decide whether to add the input  $\kappa_q(\mathbf{u}_n, \cdot)$  to the dictionary  $\mathcal{D}_{q,n}$  or not.
- 2) Vector augmentation:
- $$\tilde{\mathbf{k}}_n^{(q)} = \begin{cases} \begin{bmatrix} \mathbf{k}_n^{(q)} \\ \kappa_q(\mathbf{u}_n, \mathbf{u}_n) \end{bmatrix} & \text{if } \mathcal{D}_q \text{ grows,} \\ \mathbf{k}_n^{(q)} & \text{otherwise,} \end{cases}$$
- $$\tilde{\mathbf{h}}_n^{(q)} = \begin{cases} \begin{bmatrix} \mathbf{h}_n^{(q)} \\ 0 \end{bmatrix} & \text{if } \mathcal{D}_q \text{ grows,} \\ \mathbf{h}_n^{(q)} & \text{otherwise.} \end{cases}$$
- 3)  $\tilde{\mathbf{k}}_n = \begin{bmatrix} \tilde{\mathbf{k}}_n^{(1)} \\ \tilde{\mathbf{k}}_n^{(2)} \\ \vdots \\ \tilde{\mathbf{k}}_n^{(Q)} \end{bmatrix}$ ,  $\tilde{\mathbf{h}}_n = \begin{bmatrix} \tilde{\mathbf{h}}_n^{(1)} \\ \tilde{\mathbf{h}}_n^{(2)} \\ \vdots \\ \tilde{\mathbf{h}}_n^{(Q)} \end{bmatrix}$ .
- 4)  $w_{q,n}^2 = \left( \frac{1}{n} \sum_{l=1}^n \frac{\|\mathbf{k}_l^{(q)}\|_2^2}{r_l^{(q)}} \right)^{-1}$ ,  $\forall q = 1, 2, \dots, Q$ .
- 5)  $\mathbf{h}_{n+1} = \tilde{\mathbf{h}}_n - \mu \frac{\tilde{\mathbf{k}}_n^\top \tilde{\mathbf{h}}_n - d_n}{\tilde{\mathbf{k}}_n^\top \mathbf{W}_n^2 \tilde{\mathbf{k}}_n} \mathbf{W}_n^2 \tilde{\mathbf{k}}_n$ .
- 

$\kappa_L(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}$  and  $\kappa_C(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y} + 1$ , with the corresponding dictionaries  $\mathcal{D}_L := \{\kappa_L(\mathbf{e}_1, \cdot), \kappa_L(\mathbf{e}_2, \cdot), \dots, \kappa_L(\mathbf{e}_L, \cdot)\}$  and  $\mathcal{D}_C := \{\kappa_C(\mathbf{0}, \cdot)\}$ , respectively. An application of the above kernel functions and dictionaries to the weight (4.22) yields the following weights:

$$w_L = \left( E \left[ \frac{\|\mathbf{u}_n\|_2^2}{L} \right] \right)^{-\frac{1}{2}}, \quad w_C = 1. \quad (4.36)$$

The above weights can be approximated at each iteration by

$$w_{L,n} = \left( \frac{1}{n} \sum_{l=1}^n \frac{\|\mathbf{u}_l\|_2^2}{L} \right)^{-\frac{1}{2}}, \quad (4.37)$$

and

$$w_{C,n} = 1. \quad (4.38)$$

#### 4.3.4 Calculation of Kernel Weights (4.22) under Gaussian Input

Here the kernel weights for the linear kernel and a Gaussian kernel at the steady state are calculated in the case that the input vector is under the i.i.d. Gaussian

distribution  $\mathcal{N}(0, \sigma_u^2)$ . The dictionary size  $r$  for the Gaussian kernel is constant at the steady state and then the kernelized input is given as follows:

$$\mathbf{k}_n = \begin{bmatrix} \mathbf{u}_n \\ \mathbf{k}_n^G \end{bmatrix} \in \mathbb{R}^{L+r}. \quad (4.39)$$

**Case 1:** Linear kernel (Sec. 4.3.3)

Since  $w_C = 1$  is fixed, it is sufficient to calculate  $w_L$ .

$$E \left[ \frac{\|\mathbf{u}_n\|_2^2}{L} \right] = \frac{1}{L} \sum_{i=1}^L (2\pi\sigma_u^2)^{-\frac{L}{2}} \int_{-\infty}^{\infty} u_i^2 \exp\left(-\frac{\|\mathbf{u}\|_2^2}{2\sigma_u^2}\right) d\mathbf{u}. \quad (4.40)$$

For all  $i$ , the integral part can be rewritten as follows:

$$\begin{aligned} & \int_{-\infty}^{\infty} u_i^2 \exp\left(-\frac{\|\mathbf{u}\|_2^2}{2\sigma_u^2}\right) d\mathbf{u} \\ &= \left[ \int_{-\infty}^{\infty} u_i^2 \exp\left(-\frac{u_i^2}{2\sigma_u^2}\right) du_i \right] \prod_{j \neq i} \int_{-\infty}^{\infty} \exp\left(-\frac{u_j^2}{2\sigma_u^2}\right) du_j. \end{aligned} \quad (4.41)$$

The equality holds because the input follows the i.i.d. Gaussian distribution. Since

$$\int_{-\infty}^{\infty} \exp\left(-\frac{u_j^2}{2\sigma_u^2}\right) du_j = \sqrt{2\pi\sigma_u^2} \quad (4.42)$$

for all  $j$ , the following holds.

$$\begin{aligned} E \left[ \frac{\|\mathbf{u}_n\|_2^2}{L} \right] &= \frac{1}{L} \sum_{i=1}^L (2\pi\sigma_u^2)^{-\frac{L}{2}} \int_{-\infty}^{\infty} u_i^2 \exp\left(-\frac{u_i^2}{2\sigma_u^2}\right) du_i \\ &= \frac{1}{L} \sum_{i=1}^L \sigma_u^2 = \sigma_u^2. \end{aligned} \quad (4.43)$$

Then, by (4.36),

$$w_L = \frac{1}{\sigma_u}. \quad (4.44)$$

**Case 2:** Gaussian kernel defined in (2.10)

$$\begin{aligned} E \left[ \frac{\|\mathbf{k}_n^G\|_2^2}{r} \right] &= \frac{1}{r} (2\pi\sigma_u^2)^{-\frac{L}{2}} (2\pi\sigma^2)^{-L} \\ &\quad \times \sum_{i=1}^r \int_{-\infty}^{\infty} \exp\left(-\frac{\|\mathbf{u} - \mathbf{x}_i\|_2^2}{\sigma^2}\right) \exp\left(-\frac{\|\mathbf{u}\|_2^2}{2\sigma_u^2}\right) d\mathbf{u}, \end{aligned} \quad (4.45)$$

where  $\mathcal{D}_G = \{\kappa_G(\mathbf{x}_1, \cdot), \kappa_G(\mathbf{x}_2, \cdot), \dots, \kappa_G(\mathbf{x}_r, \cdot)\}$  is the dictionary. As in Case 1,

the integral part can be calculated as follows:

$$\begin{aligned}
& \int_{-\infty}^{\infty} \exp\left(-\frac{\|\mathbf{u} - \mathbf{x}_i\|_2^2}{\sigma^2}\right) \exp\left(-\frac{\|\mathbf{u}\|_2^2}{2\sigma_u^2}\right) d\mathbf{u} \\
&= \int_{-\infty}^{\infty} \exp\left(-\left(\frac{1}{\sigma^2} + \frac{1}{2\sigma_u^2}\right)\|\mathbf{u}\|_2^2 + \frac{2\mathbf{u}^\top \mathbf{x}_i}{\sigma^2} - \frac{\|\mathbf{x}_i\|_2^2}{\sigma^2}\right) d\mathbf{u} \\
&= \exp\left(-\frac{\|\mathbf{x}_i\|_2^2}{2\sigma_u^2 + \sigma^2}\right) \\
&\quad \times \int_{-\infty}^{\infty} \exp\left(-\left(\frac{1}{\sigma^2} + \frac{1}{2\sigma_u^2}\right)\left\|\mathbf{u} - \frac{2\sigma_u^2}{2\sigma_u^2 + \sigma^2}\mathbf{x}_i\right\|_2^2\right) d\mathbf{u} \\
&= \left(\frac{2\pi\sigma_u^2\sigma^2}{2\sigma_u^2 + \sigma^2}\right)^{\frac{L}{2}} \exp\left(-\frac{\|\mathbf{x}_i\|_2^2}{2\sigma_u^2 + \sigma^2}\right). \tag{4.46}
\end{aligned}$$

Then it holds that

$$\begin{aligned}
E\left[\frac{\|\mathbf{k}_n^G\|_2^2}{r}\right] &= \frac{1}{r}(2\pi\sigma^2)^{-L} \left(\frac{\sigma^2}{2\sigma_u^2 + \sigma^2}\right)^{\frac{L}{2}} \\
&\quad \times \sum_{i=1}^r \exp\left(-\frac{\|\mathbf{x}_i\|_2^2}{2\sigma_u^2 + \sigma^2}\right), \tag{4.47}
\end{aligned}$$

and the kernel weight is calculated as

$$w_G = \sqrt{\frac{r(2\pi\sigma^2)^L}{\left(\frac{\sigma^2}{2\sigma_u^2 + \sigma^2}\right)^{\frac{L}{2}} \sum_{i=1}^r \exp\left(-\frac{\|\mathbf{x}_i\|_2^2}{2\sigma_u^2 + \sigma^2}\right)}}. \tag{4.48}$$

For a sufficiently large  $r$ , the following approximation can be used.

$$\begin{aligned}
\frac{1}{r} \sum_{i=1}^r \exp\left(-\frac{\|\mathbf{x}_i\|_2^2}{2\sigma_u^2 + \sigma^2}\right) &\approx \int_{\mathbb{R}^L} \exp\left(-\frac{\|\mathbf{x}\|_2^2}{2\sigma_u^2 + \sigma^2}\right) d\mathbf{x} \\
&= (2\pi(2\sigma_u^2 + \sigma^2))^{\frac{L}{2}}. \tag{4.49}
\end{aligned}$$

Substituting (4.49) into (4.48) leads to

$$w_G = (2\pi\sigma^2)^{\frac{L}{4}}. \tag{4.50}$$

## 4.4 Numerical Experiments

In this section, numerical experiments related to system estimation and time series prediction are presented. The following algorithms are studied in those experiments:

1. The OMKR algorithm [43, 44] (Algorithm 1),
2. The MKNLMS algorithm [47, 48] with unweighted kernels, *i.e.*,  $w_q = 1/Q$  for all  $q$ ,

Table 4.1: Computational costs for Gaussian kernels

	# Additions	# Multiplications
OMKR (hedge)	$(2Lr_n + 1)Q - 1$	$((L + 4)r_n + 3)Q + 1$
MKNLMS (manual)	$(2L + 2)r_n - 1$	$(L + 3)r_n + 1$
MKNLMS (proposed)	$(2L + 3)r_n - 1$	$(L + 4)r_n + 3Q + 1$

3. The MKNLMS algorithm with weighted kernels, with manually tuned kernel weights which achieve the smallest MSE, and
4. The MKNLMS algorithm with weighted kernels, whose weights are tuned by the proposed automatic weighting.

As a measure for the sparsification criterion of the MKNLMS algorithm, the coherence [30] is employed. The best preset kernel weights for the MKNLMS algorithm have been found by grid search. In advance of the numerical results, see the computational costs of the above algorithms in Table 4.1, where  $L$ ,  $r_n$ , and  $Q$  are the input length, the dictionary size, and the number of kernels, respectively.

#### 4.4.1 Experiments with Toy Data

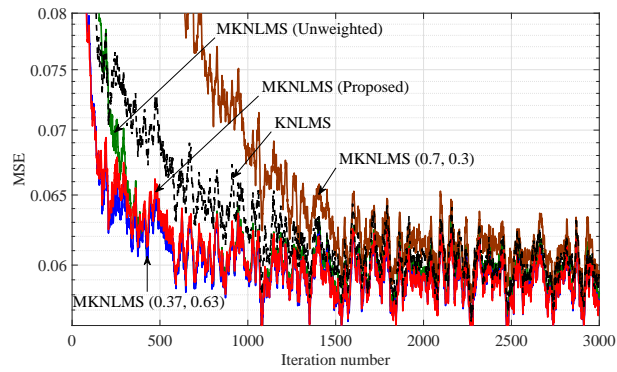
**Experiment 1a:** The efficacy of kernel weights is investigated by a simple regression task. Three preset weights for the MKNLMS algorithm are tested as well as the proposed automatic weighting. In addition, the NLMS algorithm [78] and the KNLMS algorithm [30] are also compared with the MKNLMS algorithms. The NLMS and KNLMS algorithms correspond to the cases using the linear kernel and the Gaussian kernel, respectively. The target system is given as follows:

$$\psi(u) = 0.6u + \exp\left(-\frac{(u - 0.2)^2}{2 \times 0.5^2}\right). \quad (4.51)$$

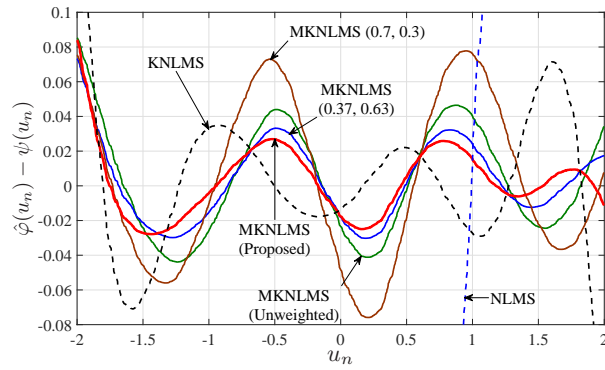
The regression task is to estimate  $\psi(u)$  by observing the current input  $u_n \in \mathbb{R}$  ( $L = 1$ ) and the noisy output  $d_n = \psi(u_n) + \nu_n$ , where  $u_n$  is randomly generated under the uniform distribution in the interval  $(-2, 2)$  and  $\nu_n$  is the additive Gaussian noise satisfying  $\text{SNR} = 10$  dB. The initial settings for this experiment are given in Table 4.2. The averaged results of 200 trials are presented in Figure 4.1. Note that the NLMS algorithm is not shown in Figure 4.1(a) because its MSE is excessively large. Figure 4.1(b) shows the differences  $\hat{\varphi}(u_n) - \psi(u_n)$ , where each  $\hat{\varphi}(u_n)$  is the estimated function by an algorithm. The kernel weights in Figure 4.1(d) are normalized so that the sum of the weights for the linear and Gaussian kernels becomes one.

**Experiment 1b:** This experiment shows that the proposed technique generates the weight sequence convergent to the weights derived in Section 4.2.2. The target system is given as

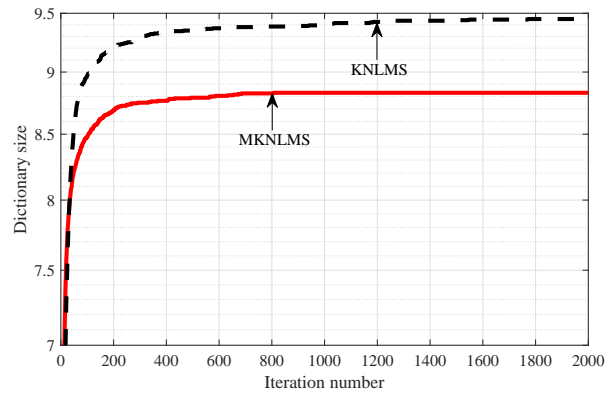
$$d_n = \beta_1^\top \mathbf{u}_n + 30 \exp\left(-\frac{\|\mathbf{u}_n - \beta_2\|_2^2}{2 \times 0.5^2}\right) - 10 \exp\left(-\frac{\|\mathbf{u}_n - \beta_3\|_2^2}{2 \times 0.5^2}\right), \quad (4.52)$$



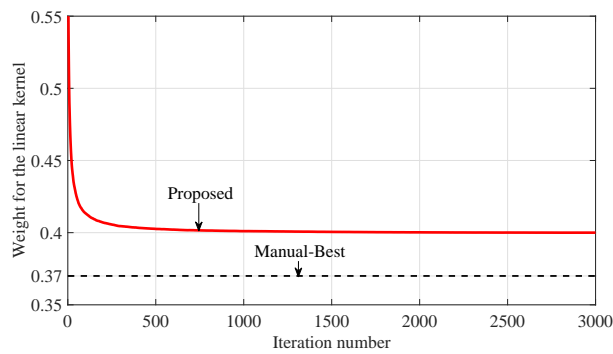
(a)



(b)



(c)

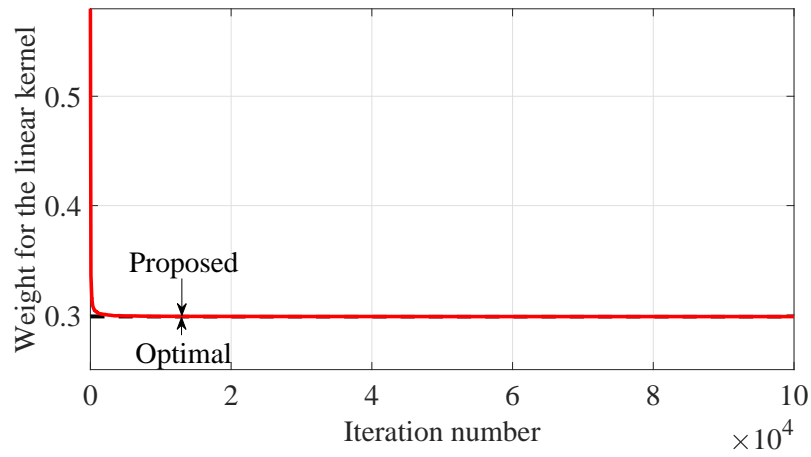


(d)

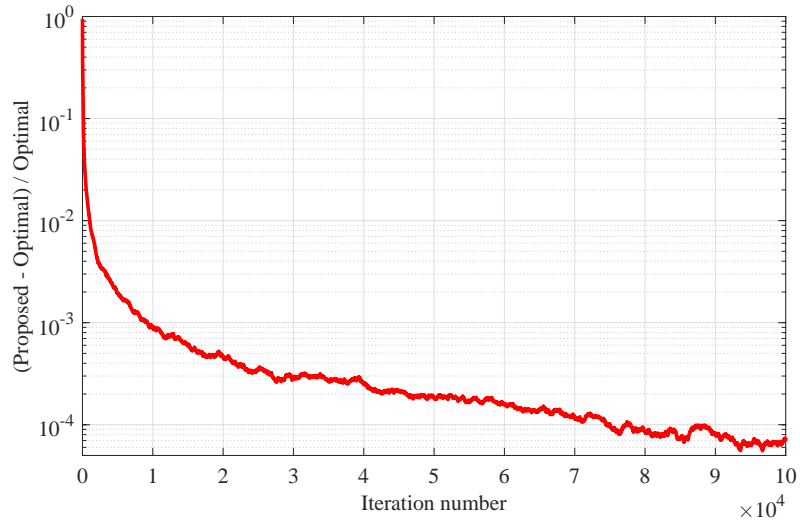
Figure 4.1: Results of Experiment 1a.

Table 4.2: Experiment 1a - Initial settings

Step size	$\mu = 0.1$
Kernel parameter	$\sigma = 0.5$
Coherence threshold	$\delta = 0.8$



(a)



(b)

Figure 4.2: Weight evolution with Gaussian input(Experiment 1b.)

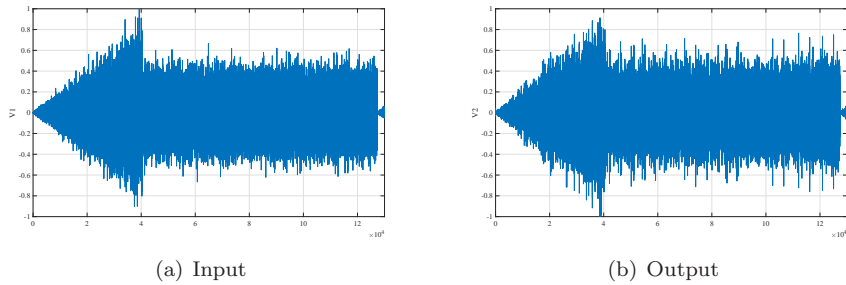


Figure 4.3: Silverbox dataset for Experiment 2a.

Table 4.3: Experiment 2a - Initial settings

Step size	$\mu = 0.7$
Kernel parameter	$\sigma = 0.9$
Discounting parameter (OMKR)	$\beta = 0.5$
Coherence threshold (MKNLMS)	$\delta = 0.9$

where the input  $\mathbf{u}_n \in \mathbb{R}^3$  is randomly generated under the i.i.d. Gaussian distribution  $\mathcal{N}(0, 0.125)$  and  $\boldsymbol{\beta}_1 = [0.5, 0.5, 0]^\top$ ,  $\boldsymbol{\beta}_2 = [0.5, 0, 0.5]^\top$ ,  $\boldsymbol{\beta}_3 = [0, 0.5, 0.5]^\top$ . The target system is estimated using the linear kernel given in Section 4.3.3 and the Gaussian kernel with the kernel parameter  $\sigma = 0.5$  and with the fixed dictionary

$$\mathcal{D}_G := \{\kappa_G(\mathbf{e}_1, \cdot), \kappa_G(\mathbf{e}_2, \cdot), \kappa_G(\mathbf{e}_3, \cdot)\}, \quad (4.53)$$

where  $\kappa_G$  and  $\mathbf{e}_1$ ,  $\mathbf{e}_2$ ,  $\mathbf{e}_3$  are the Gaussian kernel defined in (2.10) and the standard basis defined in Section 4.3.3, respectively. Figure 4.2 shows the behavior of the kernel weight for the linear kernel. If the input distribution and fixed dictionaries are known, the kernel weights at the steady state can be calculated by (4.22). See Section 4.3.4 for detailed calculation for Gaussian input. Under the normalization, the calculated weight is plotted with the dashed line, labeled as *Optimal*, while the empirically obtained weight is plotted with the red solid line, labeled as *Proposed* in Figure 4.2(a).

#### 4.4.2 Experiments with Real Data - System Estimation

**Experiment 2a:** Silverbox dataset [79] is used in this experiment. The dataset consists of the input and output data shown in Figure 4.3. The first 40000 pairs of the input and the output, the triangular part, are used as the test data. The next 80000 pairs, the rectangular part, are used as the training data. Denoting the  $n$ th values of the input data and the output data by  $u_n$  and  $d_n$ , respectively, the observed values  $u_n, u_{n-1}, u_{n-2}, u_{n-3}$  and  $d_{n-1}, d_{n-2}, d_{n-3}$  as the input of the algorithms at the  $n$ th iteration, and hence  $L = 7$ . The initial settings for this experiment are given in Table 4.3. The results of Experiment 2a are shown in Figure 4.4.

**Experiment 2b:** A subset of the Human Sensing Consortium challenge dataset

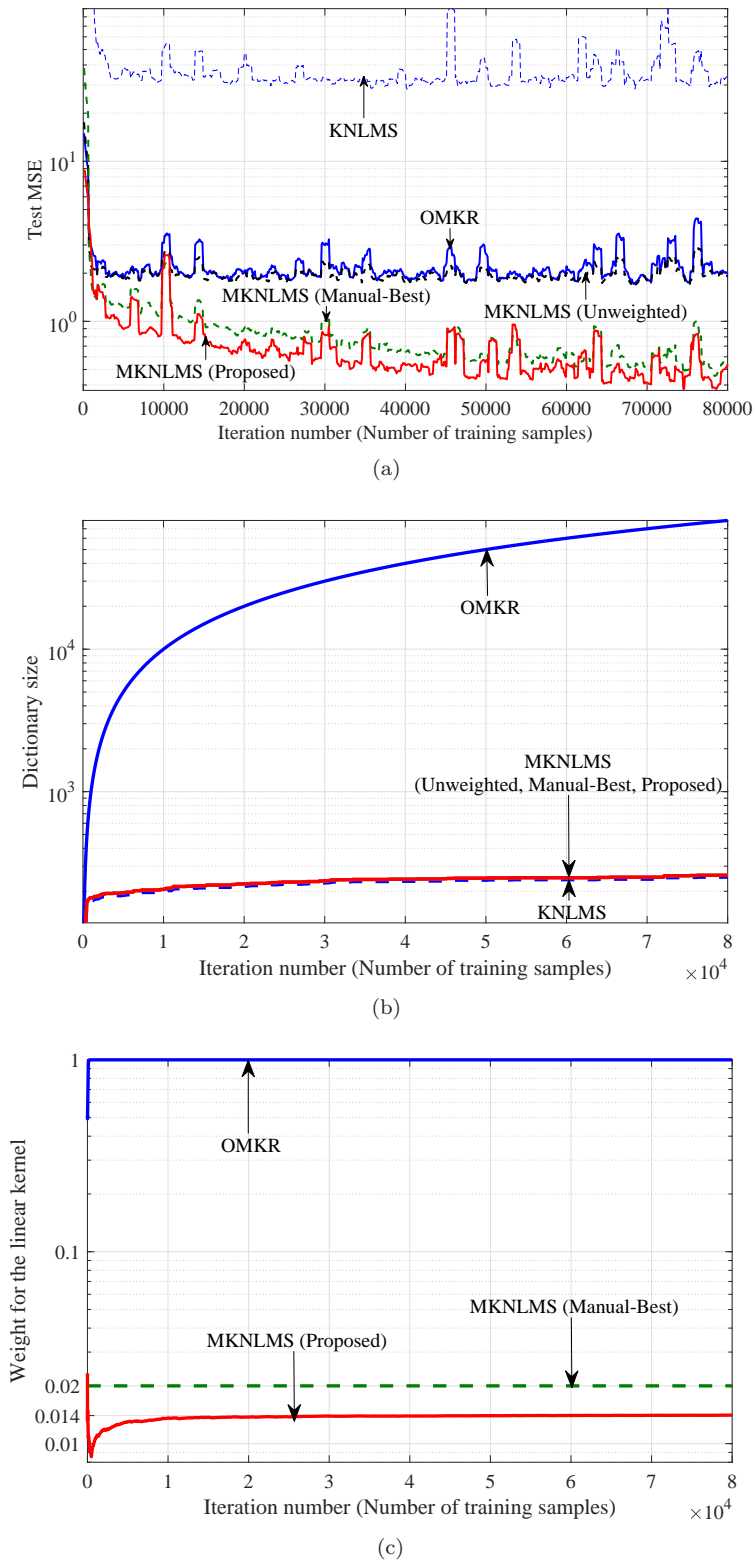
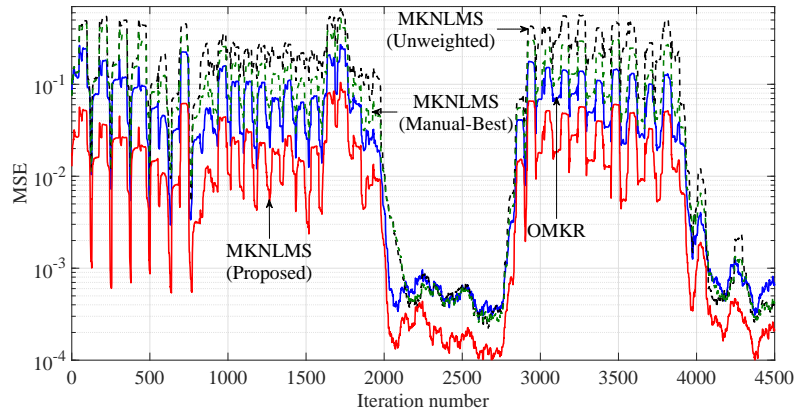


Figure 4.4: Results of Experiment 2a.

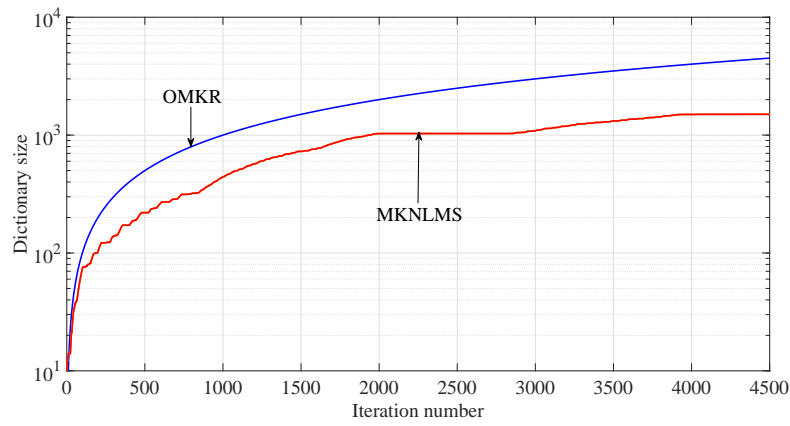


Table 4.4: Experiment 2b - Initial settings

Step size	$\mu = 0.01$
Kernel parameter	$\sigma = 0.2$
Discounting parameter (OMKR)	$\beta = 0.5$
Coherence threshold (MKNLMS)	$\delta = 0.2$



(a)



(b)

Figure 4.5: Results of Experiment 2b.

Table 4.5: Experiment 3a - Initial settings

Step size	$\mu = 0.1$
Kernel parameters	$\sigma^{(1)} = 2, \sigma^{(2)} = 1,$ $\sigma^{(3)} = 2^{-1}, \sigma^{(4)} = 2^{-2},$ $\sigma^{(5)} = 2^{-3}, \sigma^{(6)} = 2^{-4}$
Discounting parameter (OMKR)	$\beta = 0.5$
Coherence threshold (MKNLMS)	$\delta = 0.6$

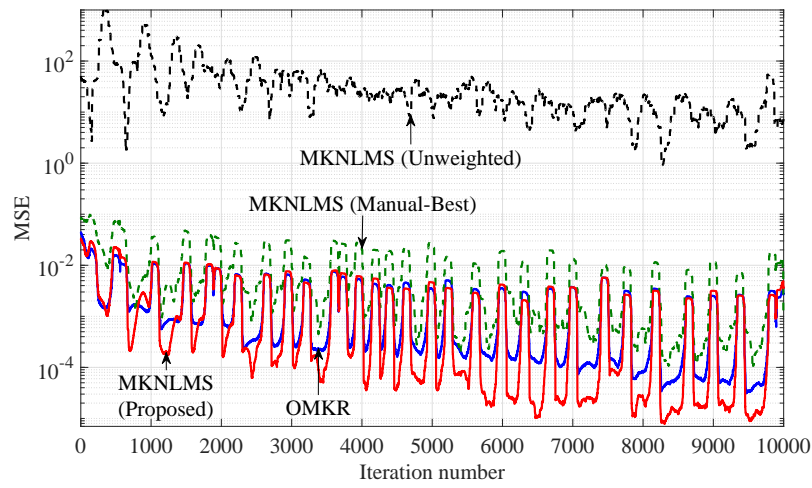
is used for this experiment. (For the detail, visit <http://hasc.jp/hc2011>.) This dataset has been obtained by three-dimensional sensor attached on a person who moves in various patterns, and hence the system to be estimated is non-stationary. The input is a vector of length  $L = 8$  of which the entries are three-dimensional coordinates  $(x, y, z)$  of the recent three samples excluding the present  $z$ , which is used as the output. One linear kernel and one Gaussian kernel are used with the initial parameters given in Table 4.4.

#### 4.4.3 Experiments with Real Data - Time Series Prediction

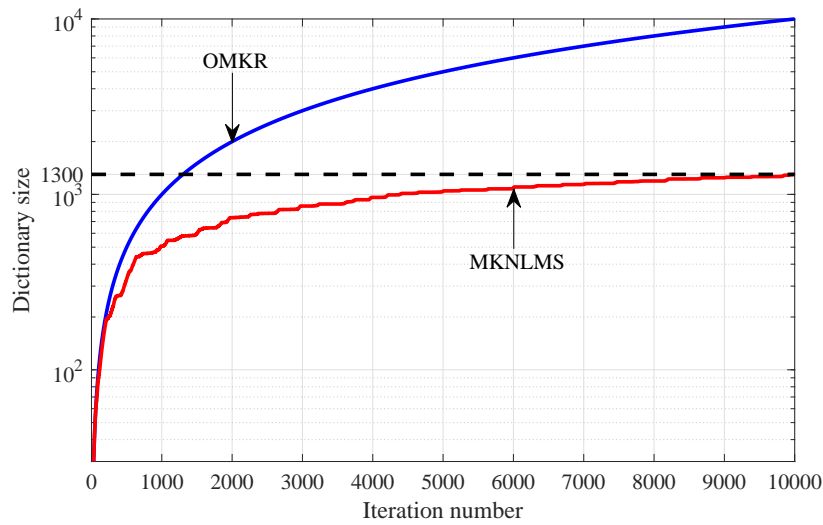
**Experiment 3a:** The Santa Fe laser dataset [80] is given as the time series data to be estimated. At each time, the next value is predicted by the algorithms given the recent nine values as input, *i.e.*,  $L = 9$ . The initial settings for this experiment are given in Table 4.5. As seen in Table 4.5, six Gaussian kernels are used for this task. Figure 4.6 shows the results of the experiment. In addition, effect of the number of kernels on performance of the proposed technique has been examined with the same dataset. See Figure 4.7 and Figure 4.8. The setting for  $Q = 6$  is the same with that in Table 4.5. For  $Q = 3$  and  $Q = 9$ , the maximal value of  $\sigma^{(q)}$  (Figure 4.7) or the minimal value (Figure 4.8) is fixed as that of Table 4.5. It is seen that, when the Gaussian kernels are used, performance of the proposed technique mainly depends on the minimal kernel parameter, rather than the number of Gaussian kernels.

**Experiment 3b:** *Buzz prediction on social media* dataset [81] is used for this experiment and Experiment 3c. The dataset is obtained from two social media: one is Tom's hardware and the other is Twitter. As a characteristic of social media, the dataset has been considered to be non-stationary. For this experiment, the data from Tom's hardware is used. The input length is  $L = 96$  and the other initial settings are given in Table 4.6. Here one linear kernel and four Gaussian kernels are used for the prediction. See Figure 4.9 for the results of the experiment.

**Experiment 3c:** This experiment uses the Twitter dataset, which has been considered to be less stationary than the data from Tom's hardware. The input length is  $L = 77$ . The linear kernel is not used for this dataset, but six Gaussian kernels are employed as shown in Table 4.7. The dictionary size of the OMKR algorithm is here limited to 2600, *i.e.*, the maximal dictionary size of the MKNLMS algorithm as Figure 4.10(b), to reduce computational costs.

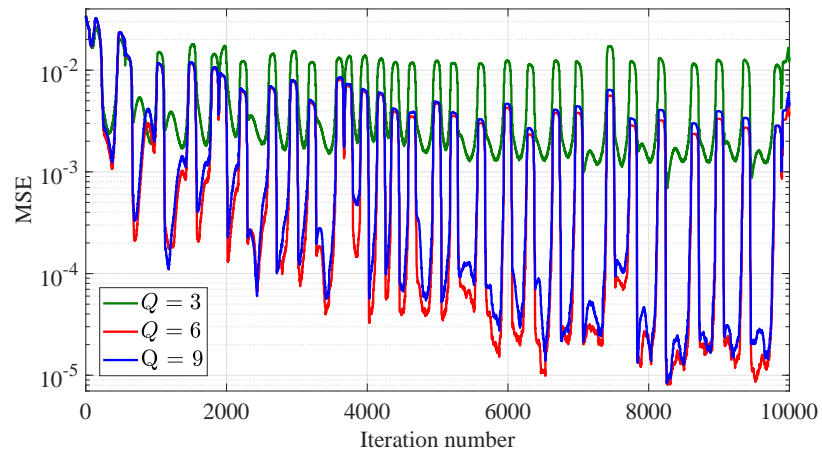


(a)

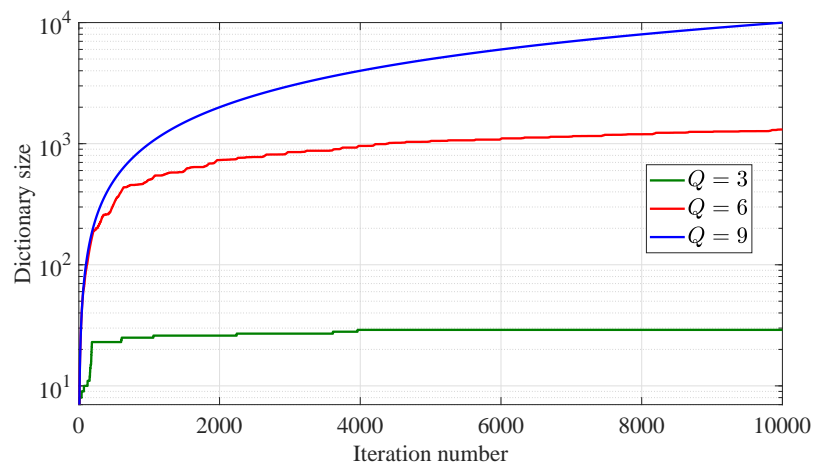


(b)

Figure 4.6: Results of Experiment 3a.

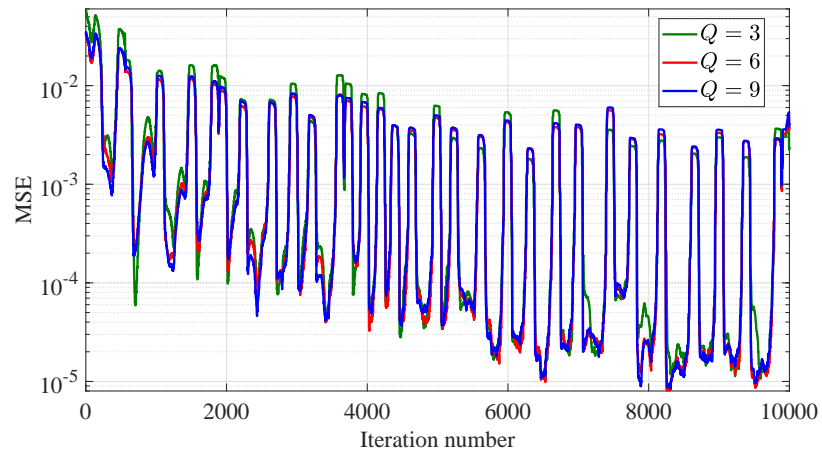


(a)

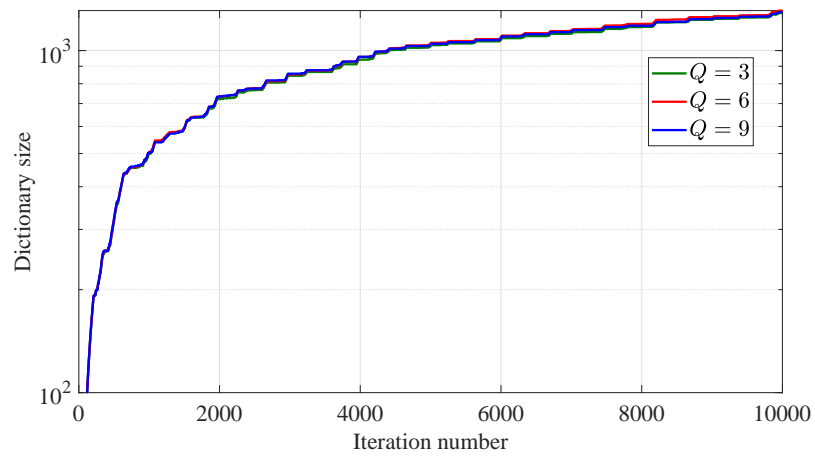


(b)

Figure 4.7: Effect of the number of kernels ( $\max_q \sigma^{(q)} = 2$ ).



(a)

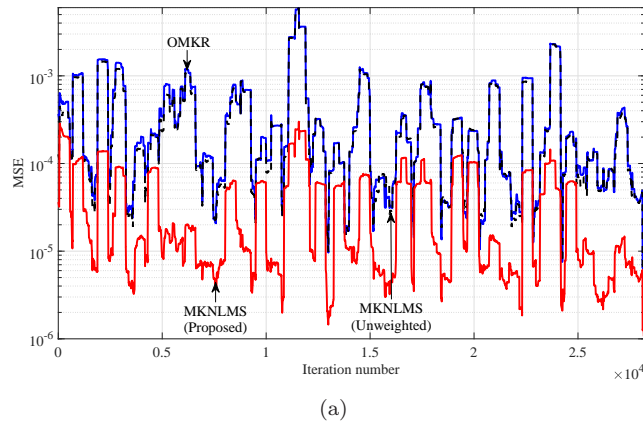


(b)

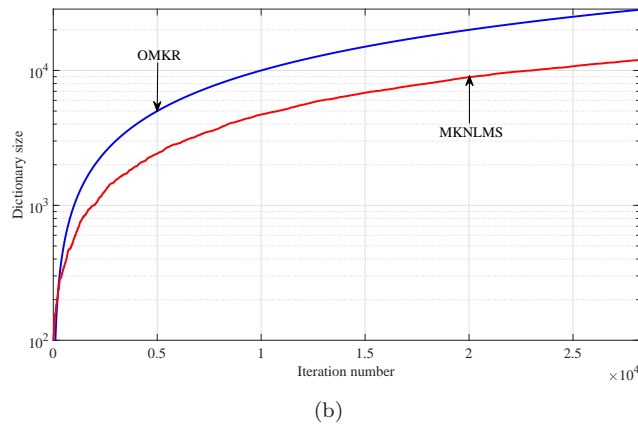
Figure 4.8: Effect of the number of kernels ( $\min_q \sigma^{(q)} = 2^{-4}$ ).

Table 4.6: Experiment 3b - Initial settings

Step size	$\mu = 0.5$
Kernel parameters	$\sigma^{(1)} = 2^{-2}, \sigma^{(2)} = 2^{-3},$ $\sigma^{(3)} = 2^{-4}, \sigma^{(4)} = 2^{-5}$
Discounting parameter (OMKR)	$\beta = 0.5$
Coherence threshold (MKNLMS)	$\delta = 0.5$



(a)

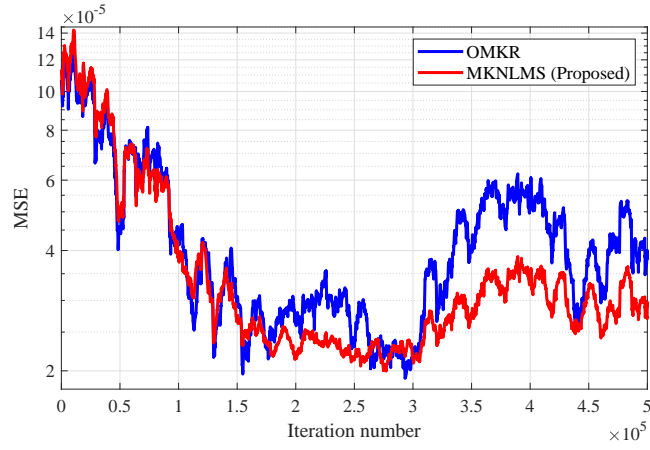


(b)

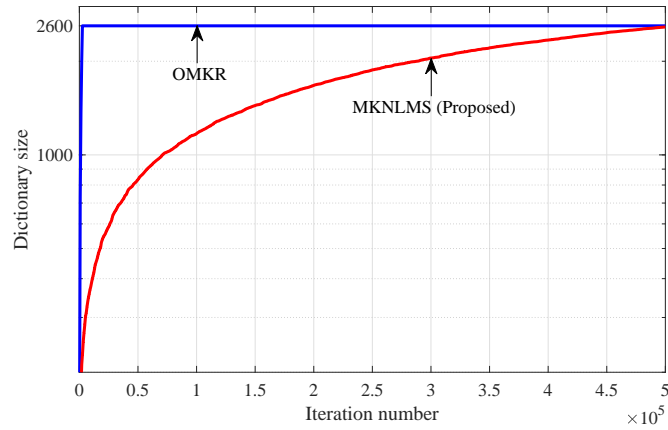
Figure 4.9: Results of Experiment 3b.

Table 4.7: Experiment 3c - Initial settings

Step size	$\mu = 0.01$
Kernel parameters	$\sigma^{(1)} = 2, \sigma^{(2)} = 1,$ $\sigma^{(3)} = 2^{-1}, \sigma^{(4)} = 2^{-2},$ $\sigma^{(5)} = 2^{-3}, \sigma^{(6)} = 2^{-4}$
Discounting parameter (OMKR)	$\beta = 0.5$
Coherence threshold (MKNLMS)	$\delta = 0.0001$



(a)



(b)

Figure 4.10: Results of Experiment 3c.

#### 4.4.4 Discussions on the Numerical Results

##### Proposed technique versus preset kernel weights:

In Figures 4.1(a) and 4.1(b), it is shown that (i) the MKNLMS algorithm with inappropriate kernel weights can be worse than the KNLMS algorithm with a single kernel by the case of (0.7, 0.3), and (ii) the proposed kernel weighting technique achieves a similar performance with the carefully tuned kernel weights (0.37, 0.63). The superior performance of the proposed technique to the manual weighting is also shown in other experiments for system estimation (Experiments 2a and 2b) and time series prediction (Experiments 3a), which use real data.

According to Table 4.1, the proposed automatic weighting requires  $r_n$  times and  $r_n + 3Q$  times of additional addition and multiplication, respectively. Compared with the computational costs of the MKNLMS algorithm itself, the additional calculations for improving the performance are affordable.

##### MKNLMS with the proposed technique versus OMKR:

In Figure 4.4(a), the MKNLMS algorithm with the proposed technique achieves both the fastest convergence and the smallest error, while the error of the OMKR algorithm decreases fast at early iterations at the price of higher steady-state errors. Figure 4.4(c) shows that the kernel weight for the linear kernel used in the OMKR algorithm converges to one approximately. This implies the Gaussian kernel is almost inactive and the OMKR algorithm fails to achieve a good balance between the linear and Gaussian kernels.

The MKNLMS algorithm with the proposed technique achieves a better MSE than the OMKR algorithm also in the other experiments, using smaller dictionaries. In Figure 4.5(a), the MKNLMS algorithm with the proposed technique shows the best tracking ability in estimating the non-stationary system (See the 2000–2500th iterations). In the experiments of time series prediction, the OMKR algorithm achieves similar MSEs with those of the MKNLMS algorithm with preset weights. (See Figures 4.10(a) and 4.9(a).)

From the viewpoint of the computational costs, the MKNLMS algorithms are more efficient than the OMKR algorithm. Focusing on the dictionary size  $r_n$ , the computational costs of the OMKR algorithm are nearly  $Q$  times larger than those of the MKNLMS algorithm even if the same dictionaries are used at every iteration.

##### Behaviors of kernel weights:

As seen in Figure 4.2, in which the dictionaries are fixed, the kernel weights set by the proposed technique converge to the values given in (4.22). In Figures 4.1(d) and 4.4(c), on the other hand, the kernel weight updated by the proposed technique converges to a different value from the manually tuned one. This is because the kernel weights (4.22) depend on the sizes of the dictionaries, which change in time. The proposed technique efficiently tracks the time-varying weights, thereby achieving an even better performance than the manual tuning which is optimized for the entire behaviors. Note that the proposed weights are also different from those of the OMKR algorithm, because the former depends on the kernelized input, whereas the latter depends on the kernel-wise errors.



## 4.5 Concluding Remarks

In this chapter, a recursive kernel weighting technique for multikernel adaptive filtering has been proposed. To deal with the coefficient error, which is unknown in practice, the case when it is a spherically distributed random variable has been considered. The transient behavior of the expected partial coefficient error then depends on the corresponding kernel weight. From the transient behavior, a kernel weight setting which achieves a uniform convergence rate for all the kernel-wise coefficient errors has been derived. A recursive formula to implement the derived kernel weights has also been presented, by the approximation using the simple sample mean. As the time-varying kernel weights make the adaptive algorithm track the corresponding time-varying solution, the derived weights have been implemented through the corresponding metric design, which makes the solution time-invariant. Finally, the numerical results have shown the efficacy of the proposed technique. The proposed technique leads to the better performance than the manually tuned kernel weights in both senses of the steady-state MSE and the convergence rate. Compared with the OMKR algorithm, the proposed approach has achieved a lower MSE, with smaller dictionaries. For the fixed dictionaries, the proposed technique has been shown to yield eventually the kernel weights derived experimentally based on the transient behavior. The kernel weights tuned by the proposed technique depend on the dictionaries, and are thus different from the manually tuned ones in general. The proposed technique will serve as a baseline for solving the kernel weighting problem, and exploiting the full potential of multikernel adaptive filtering.

## Chapter 5

# General Conclusion

This thesis has presented a nonconvex regularization for sparse optimization and a kernel weight design for multikernel adaptive filtering. In Chapter 3, both in overdetermined and underdetermined systems, the existence of a critical-point path which connects the origin and the sparsest OLS solution has been proved. This implies that  $\ell_p$ -regularization ( $0 < p < 1$ ) does yield sparsest solutions of a least squares problem. In addition, in an underdetermined system, a necessary condition for the existence of critical-point path emerging from an OLS solution has been shown. This implies that only few nonsparse OLS solution can be obtained by solving an  $\ell_p$ -regularized least squares problem. The theoretical results given in Chapter 3 will provide useful insights for developing an efficient sparsity-aware learning algorithm using the  $\ell_p$ -regularization. In recent years, nonconvex regularizers other than the  $\ell_p$ -norm have also been studied [82–89]. Some of them have been reported to preserve the convexity of the regularized cost function [82,86,88,89], and the convexity ensures that one can achieve global minima. Although the  $\ell_p$ -norm does not enjoy such a property, considering a critical-point path can be helpful for the other nonconvex regularizers. In Chapter 4, a kernel weight design equalizing the convergence rates of kernel-wise coefficient errors has been proposed. The convergence rates have been considered in the case that the coefficient error vector is a spherically distributed random variable. The kernel weights derived by equalizing the convergence rates have been implemented into the MKNLMS algorithm using the time-varying metric design corresponding to the weights. The numerical results have shown the efficacy of the proposed technique, compared to both of the manually tuned kernel weights and the OMKR algorithm. To derive the above results, we have focused on regression, which is only a major task related to online learning. For the other tasks, the results of Chapter 4 can be a baseline when one tunes kernel weights of multiple kernels. Both results will become sound bases for developing an online learning algorithm which utilizes multiple kernels appropriately, while maintaining its dictionary in an affordable size.

There are many topics beyond the sight of this thesis. Development of a sparsity-aware online learning algorithm exploiting multiple kernels may be only a small part of them. The below can be just the next step which are based on the insights from this thesis.

- Since the existence of a critical-point path connecting the origin and the

sparsest OLS solution has been proven, now a method for achieving the sparsest OLS solution using the  $\ell_p$ -regularization must be developed. Although a greedy approach has been presented in [23], it is not guaranteed to obtain a path connecting the origin and the sparsest OLS solution.

- The soft thresholding has already been used for dictionary sparsification [38]. When a weighted soft thresholding, which virtually realizes the  $\ell_p$ -minimization, is employed instead of the plain soft thresholding, the results in Chapter 4 will lead to an appropriate tuning of the threshold which coincides with the regularization parameter  $\lambda$ .
- The kernel weights can be applied to classification tasks. Since the loss function is not the MSE in this case, the derivation of appropriate weights will be also different from that given in Chapter 4.

# Bibliography

- [1] S. G. Mallat and Z. Zhang, “Matching pursuits with time-frequency dictionaries,” *IEEE Trans. Signal Processing*, vol. 41, no. 12, pp. 3397–3415, Dec. 1993.
- [2] D. L. Donoho, and I. M. Johnstone, “Ideal spatial adaptation by wavelet shrinkage,” *Biometrika*, vol. 81, no. 3, pp. 425–455, Aug. 1994.
- [3] D. L. Donoho, “De-noising by soft thresholding,” *IEEE Trans. Inform. Theory*, vol. 41, no. 3, pp. 713–627, May 1995.
- [4] S. S. Chen, D. L. Donoho and M. A. Saunders, “Atomic decomposition by basis pursuit,” *SIAM J. Sci. Comput.*, vol. 20, no. 1, pp. 33–61, 1998.
- [5] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *J. Royal Statistical Society. Series B*, vol. 58, no. 1, pp. 267–288, 1996.
- [6] Y. C. Pati, R. Rezaifar and P. S. Krishnaprasad, “Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition,” in *Proc. Asilomar*, 1993.
- [7] J. Tropp and A. C. Gilbert, “Signal recovery from partial information via orthogonal matching pursuit,” *IEEE Trans. Inform. Theory*, vol. 53, no. 12, pp. 4655–4666, Dec. 2007.
- [8] D. Needell and J. A. Tropp, “Cosamp: Iterative signal recovery from incomplete and inaccurate samples,” *Appl. Comput. Harmon. Anal.*, vol. 26, pp. 301–321, 2009.
- [9] D. Needell and R. Vershynin, “Signal recovery from incomplete and inaccurate measurements via regularized orthogonal matching pursuit,” *IEEE J. Selected Topics in Signal Process.*, vol. 4, no. 2, pp. 310–316, 2010.
- [10] D. L. Donoho and X. Huo, “Uncertainty principles and ideal atomic decomposition,” *IEEE Trans. Inform. Theory*, vol. 47, no. 7, pp. 2845–2862, Nov. 2001.
- [11] E. J. Candès and T. Tao, “Decoding by linear programming,” *IEEE Trans. Inform. Theory*, vol. 51, no. 12, pp. 4203–4215, Dec. 2005.
- [12] D. L. Donoho, “Compressed sensing,” *IEEE Trans. Inform. Theory*, vol. 52, no. 4, pp. 1289–1306, Apr. 2006.

- [13] E. J. Candès and M. B. Wakin, “An introduction to compressive sampling,” *IEEE Signal Processing Magazine*, vol. 25, no. 2, pp. 21–30, Mar. 2008.
- [14] M. Eldar and G. Kutyniok, *Compressed Sensing: Theory and Applications*, Cambridge, 2012.
- [15] M. R. Osborne, B. Presnell and B. A. Turlach, “A new approach to variable selection in least squares problems,” *IMA Journal of Numer. Anal.*, vol. 20, pp. 389–404, 2000.
- [16] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani, “Least angle regression,” *Annals of Statistics*, vol. 32, no. 2, pp. 407–499, 2004.
- [17] D. L. Donoho and Y. Tsaig, “Fast solution of  $\ell_1$ -norm minimization problems when the solution may be sparse,” *IEEE Trans. Inform. Theory*, vol. 54, no. 11, pp. 4789–4812, Nov. 2008.
- [18] R. Chartrand and V. Staneva, “Restricted isometry properties and non-convex compressive sensing,” *Inverse Problem*, vol. 24, pp. 1–14, 2008.
- [19] S. Foucart and M.-J. Lai, “Sparsest solutions of underdetermined linear systems via  $\ell_q$ -minimization for  $0 < q \leq 1$ ,” *Appl. Comput. Harmon. Anal.*, vol. 26, pp. 395–407, 2009.
- [20] Z. Xu, H. Zhang, Y. Wang, X. Chang and Y. Liang, “L1/2 regularization,” *Science China*, vol. 53, no. 6, pp. 1159–1169, June 2010.
- [21] R. Chartrand and W. Yin, “Iteratively reweighted algorithms for compressive sensing,” in *Proc. 33rd Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2008, pp. 3869–3872.
- [22] Z. Xu, X. Chang, F. Xu and H. Zhang, “ $L_{1/2}$  regularization: A thresholding representation theory and a fast solver,” *IEEE Trans. Neural Networks and Learning Systems*, vol. 23, no. 7, pp. 1013–1027, July 2012.
- [23] M. Yukawa and S. Amari, “ $\ell_p$ -regularized least squares ( $0 < p < 1$ ) and critical path,” *IEEE Trans. Inform. Theory*, vol. 62, no. 1, pp. 488–502, Jan. 2016.
- [24] S. Haykin, *Adaptive Filter Theory*, Prentice Hall, 2002.
- [25] A. Sayed, *Adaptive Filters*, Wiley-IEEE Press, 2008.
- [26] W. Liu, J. C. Príncipe and S. Haykin, *Kernel Adaptive Filtering*, Wiley, Hoboken, NJ, USA, 2010.
- [27] Y. Engel, S. Mannor, and R. Meir, “The kernel recursive least squares algorithm,” *IEEE Trans. Signal Process.*, vol. 52, no. 8, pp. 2275–2285, Aug. 2004.
- [28] W. Liu and J. Principe, “Kernel affine projection algorithms,” *EURASIP Jour. Adv. Signal Process.*, vol. 2008, pp. 1–12, 2008.
- [29] W. Liu, P. P. Pokharel and J. C. Principe, “The kernel least-mean-square algorithm,” *IEEE Trans. Signal Process.*, vol. 56, no. 2, pp. 543–554, Feb. 2008.

- [30] C. Richard, J. C. M. Bermudez, and P. Honeine, “Online prediction of time series data with kernels,” *IEEE Trans. signal Process.*, vol. 57, no. 3, pp. 1058–1067, Mar. 2009.
- [31] Z. Wang, K. Crammer and S. Vucetic, “Breaking the curse of kernelization: Budgeted stochastic gradient descent for large-scale SVM training,” *Journal of Machine Learning Research*, vol. 13, no. 100, pp. 3103–3131, 2012.
- [32] J. Platt, “A resource-allocating network for function interpolation,” *Neural Comput.*, vol. 3, no. 2, pp. 213–225, 1991.
- [33] W. Liu, I. Park and J. Principe, “An information theoretic approach of designing sparse kernel adaptive filters,” *IEEE Trans. Neural Network*, vol. 20, no. 12, pp. 1950–1960, Dec. 2009.
- [34] S. V. Vaerenbergh, J. Fernandez-Bes and V. Elvira, “On the relationship between online gaussian process regression and kernel least mean squares algorithms,” in *IEEE Int. Workshop MLSP*,, 2016.
- [35] L. Csató and M. Opper, “Sparse online gaussian processes,” *Neural Comput.*, vol. 14, no. 3, pp. 641–668, 2002.
- [36] J. A. Tropp, “Greed is good: Algorithmic results for sparse approximation,” *IEEE Trans. Inform. Theory*, vol. 50, no. 10, pp. 2231–2242, 2004.
- [37] B. Chen, S. Zhao, P. Zhu and J. C. Príncipe, “Quantized kernel least mean square algorithm,” *IEEE Trans. Neural Netw. Learn. Sys.*, vol. 23, no. 1, pp. 22–32, Jan. 2012.
- [38] W. Gao, J. Chen, C. Richard, and J. Huang, “Online dictionary learning for kernel LMS,” *IEEE Trans. signal Process.*, vol. 62, no. 11, pp. 2765–2777, Jun. 2014.
- [39] F. R. Bach, G. R. G. Lanckriet, and M. I. Jordan, “Multiple kernel learning, conic duality, and the SMO algorithm,” in *Proc. Int. Conf. Mach. Learn. ACM*, 2004.
- [40] G. Lanckriet, N. Cristianini, P. Bartlett, L. E. Ghaoui, and M. I. Jordan, “Learning the kernel matrix with semi-definite programming,” *J. Mach. Learn. Res.*, vol. 5, pp. 27–72, 2004.
- [41] M. Varma and B. R. Babu, “More generality in efficient multiple kernel learning,” in *Proc. Int. Conf. Mach. Learn. ACM*, 2009.
- [42] A. Rakotomamonjy, F. Bach, S. Canu, and Y. Grandvalet, “Simple MKL,” *J. Mach. Learn. Res.*, vol. 9, pp. 2491–2521, 2008.
- [43] D. Sahoo, S. C. H. Hoi, and B. Li, “Online multiple kernel regression,” in *Proc. the 20th ACM SIGKDD*, 2014, pp. 293–302.
- [44] D. Sahoo, S. C. H. Hoi, and B. Li, “Large scale online multiple kernel regression with application to time-series prediction,” *ACM Trans. Knowledge Discovery from Data*, vol. 13, no. 1, pp. 9:1–9:33, Jan. 2019.

- [45] R. Pokharel, S. Seth and J. C. Principe, “Mixture kernel least mean square,” in *Proc. 2013 Int. Joint Conf. on Neural Networks (IJCNN)*, 2013, pp. 1–7.
- [46] W. Gao, C. Richard, J.-C. M. Bermudez and J. Huang, “Convex combinations of kernel adaptive filters,” in *IEEE Int. Workshop MLSP*, 2014.
- [47] M. Yukawa, “Multikernel adaptive filtering,” *IEEE Trans. Signal Process.*, vol. 60, no. 9, pp. 4672–4682, Sep. 2012.
- [48] M. Yukawa, “Adaptive learning in Cartesian product of reproducing kernel Hilbert spaces,” *IEEE Trans. Signal Process.*, vol. 63, no. 22, pp. 6037–6048, Nov. 2015.
- [49] O. Toda and M. Yukawa, “Online model-selection and learning for nonlinear estimation based on multikernel adaptive filtering,” *IEICE transactions on Fundamentals of Electronics, Communications and Computer Sciences*, vol. E100.A, no. 1, pp. 236–250, 2017.
- [50] M. Kasparick, R. L. G. Cavalcante, S. Valentin, S. Stánczak and M. Yukawa, “Kernel-based adaptive online reconstruction of coverage maps with side information,” *IEEE Trans. Vehicular Tech.*, vol. 65, no. 7, pp. 5461–5473, Jul. 2016.
- [51] D. A. Awan, R. L. G. Cavalcante, M. Yukawa and S. Stanczak, “Detection for 5G-NOMA: An online adaptive machine learning approach,” in *2018 IEEE International Conference on Communications (ICC)*, 2018.
- [52] D. A. Awan, R. L. G. Cavalcante, M. Yukawa and S. Stanczak, “Adaptive learning for symbol detection: A reproducing kernel hilbert space approach,” in *Machine Learning for Future Wireless Communications*, F. Luo, Ed., chapter 11, pp. 197–211. Wiley, 2019.
- [53] B.-S. Shin, M. Yukawa, R. L. G. Cavalcante and A. Dekorsy, “Distributed adaptive learning with multiple kernels in diffusion networks,” *IEEE Trans. Signal Process.*, vol. 66, no. 21, pp. 5505–5519, Nov. 2018.
- [54] X. Chen, F. Xu and Y. Ye, “Lower bound theory of nonzero entries in solutions of  $\ell_2$ - $\ell_p$  minimization,” *SIAM J. Scientific Computing*, vol. 32, no. 5, pp. 2832–2852, 2010.
- [55] M. Wang, W. Xu and A. Tang, “On the performance of sparse recovery via  $\ell_p$ -minimization ( $0 \leq p \leq 1$ ),” *IEEE Trans. Inform. Theory*, vol. 57, no. 11, pp. 7255–7278, Nov. 2011.
- [56] M. Yukawa and S. Amari, “On extensions of LARS by Information Geometry : Convex objectives and  $\ell_p$ -norm,” in *APSIPA Annual Summit and Conference: Special Session on Recent Advances in Adaptive/Sparse Signal Processing*, Oct. 2011.
- [57] M. Yukawa and S. Amari, “ $\ell_p$ -constrained least squares ( $0 < p < 1$ ) and its critical path,” in *Proc. IEEE International Symposium on Information Theory (ISIT)*, 2012, pp. 2231–2235.

- [58] J. Kivinen, A. J. Smola, and R. C. Williamson, “Online learning with kernels,” *IEEE Trans. signal Process.*, vol. 52, no. 8, pp. 2165–2176, Aug. 2004.
- [59] M. Yukawa and R. Ishii, “Online model selection and learning by multi-kernel adaptive filtering,” in *Proc. 21th Eur. Signal Process. Conf. (EUSIPCO)*, 2013, pp. 1–5.
- [60] F. Albu and K. Nishikawa, “The kernel proportionate NLMS algorithm,” in *Proc. 21th Eur. Signal Process. Conf. (EUSIPCO)*, 2013, pp. 1–5.
- [61] M. Takizawa and M. Yukawa, “Adaptive nonlinear estimation based on parallel projection along affine subspaces in reproducing kernel Hilbert space,” *IEEE Trans. signal Process.*, vol. 63, no. 16, pp. 4257–4269, Aug. 2015.
- [62] M. Yukawa and K. R. Müller, “Why does a Hilbertian metric work efficiently in online learning with kernels?,” *IEEE Signal Process. Letter*, vol. 23, no. 10, pp. 1424–1428, Oct. 2016.
- [63] S. Wang, L. Dang, B. Chen, C. Ling, L. Wang and S. Duan, “Kernel online learning algorithm with scale adaptation,” *IEEE Trans. Circuits and Systems*, vol. 65, no. 11, pp. 1788–1792, Nov. 2018.
- [64] F. A. Tobar, S.-Y. Kung, and D. P. Mandic, “Multikernel least mean square algorithm,” *IEEE Trans. Neural Net. Learn. Sys.*, vol. 25, no. 2, pp. 265–277, Feb. 2014.
- [65] K. Nishikawa and F. Albu, “Implementation method of kernel adaptive filter as an add-on for a linear adaptive filter,” in *Proc. 23rd Eur. Signal Process. Conf. (EUSIPCO)*, 2015, pp. 2736–2740.
- [66] S. V. Vaerenbergh, S. Scardapane and I. Santamaria, “Recursive multikernel filters exploiting nonlinear temporal structure,” in *Proc. 25th Eur. Signal Process. Conf. (EUSIPCO)*, 2017, pp. 2674–2678.
- [67] M. Ohnishi and M. Yukawa, “Online nonlinear estimation via iterative  $L^2$ -space projections: Reproducing kernel of subspace,” *IEEE Trans. Signal Process.*, vol. 66, no. 15, pp. 4050–4064, Aug. 2018.
- [68] T. Wada, K. Fukumori and T. Tanaka, “Dictionary learning for Gaussian kernel adaptive filtering with variable kernel center and width,” in *Proc. 43th Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2018, pp. 2766–2770.
- [69] M. Takizawa and M. Yukawa, “Online learning with self-tuned Gaussian kernels: Good kernel-initialization by multiscale screening,” in *Proc. 44th Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2019, pp. 4863–4867.
- [70] S. Ioffe and C. Szegedy, “Batch normalization: accelerating deep network training by reducing internal covariate shift,” in *Proc. the 32th International Conference on Machine Learning (ICML)*, 2015, pp. 448–456.
- [71] S. Santurkar, D. Tsipras, A. Ilyas and A. Madry, “How does batch normalization help optimization?,” in *Proc. Neural Information Processing Systems (NIPS) 2018*, 2018, pp. 2483–2493.



- [72] N. Bjorck, C. P. Gomes, B. Selman and K. Q. Weinberger, “Understanding batch normalization,” in *Proc. Neural Information Processing Systems (NIPS) 2018*, 2018, pp. 7694–7705.
- [73] S. S. Narayan, A. M. Peterson and M. J. Narasimha, “Transform domain LMS algorithm,” *IEEE Trans. Acoust., Speech, Signal Process.*, vol. ASSP-31, pp. 609–615, Jun. 1983.
- [74] D. F. Marshall, W. K. Jenkins and J. J. Murphy, “The use of orthogonal transforms for improving performance of adaptive filters,” *IEEE Trans. Circuits Syst.*, vol. 36, no. 4, pp. 474–484, 1989.
- [75] M. Yukawa, K. Slavakis, and I. Yamada, “Adaptive parallel quadratic-metric projection algorithms,” *IEEE Trans. Audio, Speech and Language Process.*, vol. 15, no. 5, pp. 1665–1680, Jul. 2007.
- [76] M. Yukawa and I. Yamada, “Krylov-proportionate adaptive filtering techniques not limited to sparse systems,” *IEEE Trans. Signal Process.*, vol. 57, no. 3, pp. 927–943, 2009.
- [77] M. Yukawa and I. Yamada, “A unified view of adaptive variable-metric projection algorithms,” *EURASIP Jour. Adv. Signal Process.*, vol. 2009, pp. 1–13, 2009.
- [78] J. Nagumo and A. Noda, “A learning method for system identification,” *IEEE Trans. Automatic Control*, vol. 12, no. 3, pp. 282–287, 1967.
- [79] T. Wigren and J. Schoukens, “Three free data sets for development and benchmarking in nonlinear system identification,” in *European Control Conference (ECC)*, 2013, pp. 2933–2938.
- [80] A.S. Weigend and N.A. Gershenfeld, *Time series prediction: Forecasting the future and understanding the past*, Addison-Wesley, Reading, MA, 1993.
- [81] F. Kawala, A. Douzal-Chouakria, E. Gaussier, and E. Dimert, “Prédictions d’activité dans les réseaux sociaux en ligne,” in *In Actes de la Conférence sur les Modèles et l’Analyse des Réseaux : Approches Mathématiques et Informatique (MARAMI)*, 2013, p. 16.
- [82] H. Mohimani, M. Babaie-Zadeh and C. Jutten, “A fast approach for over-complete sparse decomposition based on smoothed  $\ell^0$  norm,” *IEEE Trans. Signal Process.*, vol. 57, no. 1, pp. 289–301, Jan. 2009.
- [83] D. Malioutov and A. Aravkin, “Iterative log thresholding,” in *Proc. 39th Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2014, pp. 7198–7202.
- [84] Z. Li, S. Ding and Y. Li, “Dictionary learning with log-regularizer for sparse representation,” *IEEE Int’l conf. DSP*, pp. 609–613, 2015.
- [85] I. Bayram, “Penalty functions derived from monotone mappings,” *IEEE Signal Process. Letter*, vol. 22, no. 3, pp. 265–269, Mar. 2015.
- [86] E. Soubies, L. Blanc-Féraud and G. Aubert, “A continuous exact  $\ell_0$  penalty (cel0) for least squares regularized problem,” *SIAM Journal on Imaging Sciences*, vol. 8, no. 3, pp. 1607–1639, 2015.

- [87] M. Sadeghi and M. Babaie-Zadeh, “Iterative sparsification-projection: Fast and robust sparse signal approximation,” *IEEE Trans. Signal Process.*, vol. 64, no. 21, pp. 5536–5548, Nov. 2016.
- [88] I. W. Selesnick and I. Bayram, “Enhanced sparsity by non-separable regularization,” *IEEE Trans. Signal Process.*, vol. 64, no. 9, pp. 2298–2313, May 2016.
- [89] I. Selesnick, “Sparse regularization via convex analysis,” *IEEE Trans. Signal Process.*, vol. 65, no. 17, pp. 4481–4494, Sep. 2017.

# Publications Related to the Dissertation

## Journal Article

1. K. Jeong, M. Yukawa and S. -i. Amari, “Can critical-point paths under  $\ell_p$ -regularization ( $0 < p < 1$ ) reach the sparsest least squares solutions?”, IEEE Transactions on Information Theory, Vol. 60, No. 5, pp. 2960–2968, May 2014.
2. K. Jeong and M. Yukawa, “Kernel weights for equalizing kernel-wise convergence rates of multikernel adaptive filtering”, IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences Vol. E104-A, No. 6, Jun. 2021.

## International Conference Paper

1. K. Jeong and M. Yukawa, “Automatic kernel weighting for multikernel adaptive filtering: Multiscale aspects”, Proceedings of the 44th International Conference on Acoustic, Speech and Signal Processing, pp. 3152–3156, 2019.

# Other Publications

## International Conference Paper

1. K. Jeong, M. Yukawa, M. Yamagishi and I. Yamada, “Automatic shrinkage tuning robust to input correlation for sparsity-aware adaptive filtering”, Proceedings of the 43rd International Conference on Acoustic, Speech and Signal Processing, pp. 4314–4318, 2018.

## Domestic Conference Papers

1. K. Jeong, M. Yukawa and S. -i. Amari, “On relation between  $\ell_p$ -regularization ( $0 < p < 1$ ) and sparsest least squares solution,” in Proceedings of Symposium on Information Theory and Its Applications (SITA), pp. 196–199, Nov. 2013.
2. K. Jeong, M. Yukawa, M. Yamagishi and I. Yamada, “A simple self-tuning technique for adaptive proximal forward-backward splitting method”, Technical Report, IEICE-SIP, pp. 71–76, 2017.
3. K. Jeong, M. Yukawa, M. Yamagishi and I. Yamada, “A time-averaging approach to automatic shrinkage tuning under colored signals for sparsity-aware adaptive filtering”, Proceedings of SIP symposium, Nov. 2017.