

9

Adaptive Kernel Learning for Signal Processing

Adaptive filtering is a central topic in signal processing. An adaptive filter is a filter structure provided with an adaptive algorithm that tunes the transfer function, typically driven by an error signal. Adaptive filters are widely applied in non-stationary environments because they can adapt their transfer function to match the changing parameters of the system that generates the incoming data (Hayes 1996; Widrow *et al.* 1975). They have become ubiquitous in current digital signal processing, mainly due to the increase in computational power and the need to process streamed data. Adaptive filters are now routinely used in all communication applications for channel equalization, array beamforming or echo cancellation, to cite some, and in other areas of signal processing as image processing or medical equipment.

By applying linear adaptive filtering principles in the kernel feature space, powerful nonlinear adaptive filtering algorithms can be obtained. This chapter introduces the wide topic of adaptive signal processing, and it explores the emerging field of kernel adaptive filtering (KAF). Its orientation is different from the preceding ones, as adaptive processing can be used in a variety of scenarios. Attention is paid to kernel LMS/RLS algorithms, to previous taxonomies of adaptive kernel methods, and to emergent kernel methods.

Matlab code snippets are included to illustrate the basic operations of the most common kernel adaptive filters. Tutorial examples are provided on applications including chaotic time-series prediction, respiratory motion prediction, and nonlinear system identification.

9.1 Linear Adaptive Filtering

Let us first define some basic concepts of linear adaptive filtering theory. The goal of adaptive filtering is to model an unknown, possibly time-varying system by observing the inputs and outputs of this system over time. We will denote the input to the system on time instant n as \mathbf{x}_n , and its output as d_n , see Fig. 9.1. The input signal x_n is assumed to be zero-mean. We represent it as a vector, and it will often represent a time-delay vector of L taps of a signal x_n on time instant n as $\mathbf{x}_n = [x_n, x_{n-1}, \dots, x_{n-L+1}]^\top$.

A diagram for a linear adaptive filter is depicted in Fig. 9.2. The input to the adaptive filter on time instant n is \mathbf{x}_n , and its response y_n is obtained through the following linear operation:

$$y_n = \mathbf{w}_n^H \mathbf{x}_n, \quad (9.1)$$

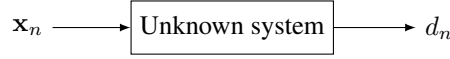


Figure 9.1 An unknown system with input x_n and output d_n at time instant n .

72 where H is the hermitic operator.

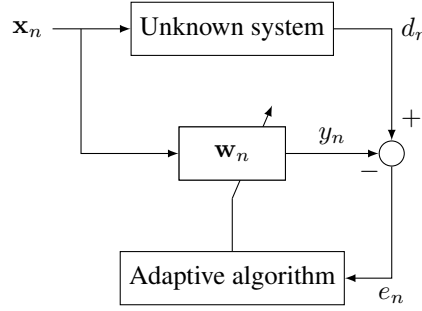


Figure 9.2 A linear adaptive filter for system identification.

73 Linear adaptive filtering follows the *online learning* framework, which consists of two
 74 basic steps that are repeated at each time instant n . First, the online algorithm receives
 75 an observation x_n for which it calculates the estimated image y_n , based on its current
 76 estimate of w_n . Next, the algorithm receives the desired output d_n (also known as *symbol* in
 77 communications), which allows it to calculate the estimation error $e_n = d_n - y_n$ and update
 78 its estimate for w_n . In some situations d_n is known a priori, i.e. the received signal x_n
 79 is one of a set of *training* signals provided with known labels. The procedure is then called
 80 *supervised*. When d_n belongs to a finite set of quantized labels, and it can be assumed that the
 81 error will be likely much smaller than the quantization step or minimum Euclidean distance
 82 between labels, the desired label is estimated by quantizing y_n to the closest label. In these
 83 cases, the algorithm is called *decision directed*.

84 9.1.1 LMS Algorithm

The classical adaptive optimization techniques have their roots in the theoretical approach called the steepest-descent algorithm. Assume that the expectation of the squared error signal, $J_n = E\{|e_n|^2\}$ can be computed. Since this error is a function of the vector w_n , the idea of the algorithm is to modify this vector towards the direction of the steepest descent of J_n . This direction is just opposite to its gradient $\nabla_w J_n$. Indeed, assuming complex stationary signals, the error expectation is

$$\begin{aligned}
 E\{|e_n|^2\} &= E\{|d_n - \mathbf{w}_n^H \mathbf{x}_n|^2\} \\
 &= E\{|d_n|^2 + \mathbf{w}_n^H \mathbf{x}_n \mathbf{x}_n^H \mathbf{w}_n - 2\mathbf{w}_n^H \mathbf{x}_n d_n^*\} \\
 &= \sigma_d^2 + \mathbf{w}_n^H \mathbf{R}_{xx} \mathbf{w}_n - 2\mathbf{w}_n^H \mathbf{p}_{xd},
 \end{aligned} \tag{9.2}$$

where \mathbf{R}_{xx} is the signal autocorrelation matrix, \mathbf{p}_{xd} is the cross-correlation vector between the signal and the filter output, and σ_d^2 is the variance of the system output. Its gradient with

respect to vector \mathbf{w}_n is expressed as

$$\nabla_{\mathbf{w}} J_n = 2\mathbf{R}_{xx}\mathbf{w}_n - 2\mathbf{p}_{xd}. \quad (9.3)$$

The adaptation rule based on steepest descent thus becomes

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta \nabla_{\mathbf{w}} J_n, \quad (9.4)$$

85 where η is the *step size* or *learning rate* of the algorithm.

The Least Mean Squares (LMS) algorithm, introduced in 1960 by Widrow (Widrow *et al.* 1975), is a very simple and elegant method of training a linear adaptive system to minimize the mean square error that approximates the gradient $\nabla_{\mathbf{w}} J_n$ using an instantaneous estimate of the gradient. From Eq. (9.3), an approximation can be written as

$$\nabla_{\mathbf{w}} J_n \approx 2\mathbf{x}_n \mathbf{x}_n^H \mathbf{w}_n - 2\mathbf{x}_n d_n^* \quad (9.5)$$

Using this approximation in Eq. (9.4) leads to the well-known stochastic gradient descent update rule which is the core of the LMS algorithm:

$$\begin{aligned} \mathbf{w}_{n+1} &= \mathbf{w}_n - \eta \mathbf{x}_n (\mathbf{x}_n^H \mathbf{w}_n - d_n^*) \\ &= \mathbf{w}_n + \eta \mathbf{x}_n e_n^*. \end{aligned} \quad (9.6)$$

86 This optimization procedure is also the basis for tuning nonlinear filter structures such as
87 neural networks Haykin (2001) and some of the kernel-based adaptive filters discussed later
88 in this chapter. A detailed analysis including that of convergence and misadjustment is given
89 in (Haykin 2001). The Matlab code for the LMS training step on a new data pair (\mathbf{x}, d) is
90 displayed in Listing 9.1.

```
91 y = x'*w; % evaluate filter output
92 err = d - y; % instantaneous error
93
94 w = w + mu*x*err'; % update filter coefficients
```

Listing 9.1 Training step of the LMS algorithm on a new datum (\mathbf{x}, d) .

95 Under the stationarity assumption, the LMS algorithm converges to the Wiener solution
96 in mean, but the weight vector \mathbf{w}_n shows a variance that converges to a value that is a
97 function of η . Therefore, low variances are only achieved at low adaptation speed. A more
98 sophisticated approach with faster convergence is found in the the Recursive Least-Squares
99 (RLS) algorithm.

100 9.1.2 RLS Algorithm

101 The RLS algorithm was first introduced by Plackett in 1950 (Plackett 1950). In a stationary
102 scenario, it converges to the Wiener solution in mean and variance, improving also the slow
103 rate of adaptation of the LMS algorithm. Nevertheless, this gain in convergence speed comes
104 at the price of a higher complexity, as we will see below.

105 Recursive update

The basis of the RLS algorithm consists of recursively updating the vector \mathbf{w} that minimizes a regularized version of the cost function J_n

$$J_n = \sum_{i=1}^n |d_i - \mathbf{x}_i^H \mathbf{w}|^2 + \delta \mathbf{w}^H \mathbf{w}, \quad (9.7)$$

where δ is a positive constant *regularization factor*. The regularization factor penalizes the squared norm of the solution vector so that this solution does not apply too much weight to any specific dimension¹. The solution that minimizes the least-squares cost function (9.7) is well known and given by

$$\mathbf{w} = (\mathbf{R}_{xx} + \delta \mathbf{I})^{-1} \mathbf{p}_{xd}. \quad (9.8)$$

106 The regularization δ guarantees that the inverse in Eq. (9.8) exists. In the absence of
107 regularization, i.e. for $\delta = 0$, the solution requires to invert the matrix \mathbf{R}_{xx} , which may be
108 rank-deficient.

109 For a detailed derivation of the RLS algorithm we refer the reader to Haykin (2001);
110 Sayed (2003). In the sequel, we will provide its update equations and a short discussion
111 of its properties compared to LMS.

We denote the autocorrelation matrix for the data \mathbf{x}_1 till \mathbf{x}_n as \mathbf{R}_n ,

$$\mathbf{R}_n = \sum_{i=1}^n \mathbf{x}_i^H \mathbf{x}_i. \quad (9.9)$$

RLS requires the introduction of an *inverse autocorrelation matrix* \mathbf{P}_n , defined as

$$\mathbf{P}_n = (\mathbf{R}_n + \delta \mathbf{I})^{-1}. \quad (9.10)$$

At step $n - 1$ of the recursion, the algorithm has processed $n - 1$ data, and its estimate \mathbf{w}_{n-1} is the optimal solution for minimizing the squared cost function (9.7) at time step $n - 1$. When a new datum \mathbf{x}_n is obtained, the inverse autocorrelation matrix is updated as

$$\mathbf{P}_n = \mathbf{P}_{n-1} - \mathbf{g}_n \mathbf{g}_n^H \mathbf{P}_{n-1}, \quad (9.11)$$

where \mathbf{g}_n is the *gain vector* of the RLS algorithm,

$$\mathbf{g}_n = \frac{\mathbf{P}_{n-1} \mathbf{x}_n}{1 + \mathbf{x}_n^H \mathbf{P}_{n-1} \mathbf{x}_n}. \quad (9.12)$$

The update of the solution itself reads

$$\mathbf{w}_n = \mathbf{w}_{n-1} + \mathbf{g}_n e_n, \quad (9.13)$$

in which e_n represents the usual prediction error $e_n = d_n - \mathbf{x}_n^H \mathbf{w}_{n-1}$. The RLS algorithm starts by initializing its solution to

$$\mathbf{w}_0 = 0, \quad (9.14)$$

and the estimate of the inverse autocorrelation matrix \mathbf{P}_n to

$$\mathbf{P}_0 = \delta^{-1} \mathbf{I}. \quad (9.15)$$

112 Then, it recursively updates its solution by including one datum \mathbf{x}_i at a time and performing
113 the calculations from Eqs. 9.11 to (9.13).

114 Due to the matrix multiplications involved in the RLS updates, the computational
115 complexity per time step for RLS is quadratic in terms of the data dimension, $\mathcal{O}(L^2)$, while
116 the LMS algorithm has only linear complexity, $\mathcal{O}(L)$.

¹A slightly more general formulation involves a *regularization matrix* which can penalize the individual elements of the solution differently Sayed (2003).

117 Exponentially-weighted RLS

118 The RLS algorithm takes into account all previous data when it updates its solution with
 119 a new datum. This kind of update yields a faster convergence than LMS, which guides its
 120 update based only on the performance on the newest datum. Nevertheless, by guaranteeing
 121 that its solution is valid for all previous data, the RLS algorithm is in essence looking for
 122 a stationary solution, and thus it cannot adapt to nonstationary scenarios, where a tracking
 123 algorithm is required. LMS, on the other hand, deals correctly with nonstationary scenarios,
 124 thanks to the instantaneous nature of its update which forgets older data and only adapts to
 125 the newest datum.

A tracking version of RLS can be obtained by including a *forgetting factor* $\lambda \in (0, 1]$ in its cost function, as follows

$$J_n = \sum_{i=1}^n \lambda^{n-i} \|d_i - \mathbf{x}_i^H \mathbf{w}\|^2 + \lambda^n \delta \mathbf{w}^H \mathbf{w}. \quad (9.16)$$

The resulting algorithm is called exponentially-weighted RLS [Haykin \(2001\)](#); [Sayed \(2003\)](#). The inclusion of the forgetting factor assigns lower weights to older data, which allows the algorithm to adapt gradually to changes. The update for the inverse autocorrelation matrix becomes

$$\mathbf{P}_n = \lambda^{-1} \mathbf{P}_{n-1} - \lambda^{-1} \mathbf{g}_n \mathbf{x}_n^H \mathbf{P}_{n-1}, \quad (9.17)$$

and the new gain vector becomes

$$\mathbf{g}_n = \frac{\lambda^{-1} \mathbf{P}_{n-1}}{1 + \lambda^{-1} \mathbf{x}_n^H \mathbf{P}_{n-1} \mathbf{x}_n}. \quad (9.18)$$

126 The Matlab code for the training step of the exponentially-weighted RLS algorithm is
 127 displayed in Listing 9.2.

```

128 y = x'*w; % evaluate filter output
129 err = d - y; % instantaneous error
130
131 g = P*x/(lambda+x'*P*x); % gain vector
132 w = w + g*err; % update filter coefficients
133 P = lambda\(P - g*x'*P); % update inverse autocorrelation matrix

```

Listing 9.2 Training step of the exponentially-weighted RLS algorithm on a new datum (\mathbf{x}, d) .

134 Recursive estimation algorithms play a crucial role for many problems in adaptive control,
 135 adaptive signal processing, system identification, and general model building and monitoring
 136 problems [Ljung \(1999\)](#). In the signal processing literature, great attention has been paid to
 137 their efficient implementation. Linear autoregressive models require relatively few parameters
 138 and allow closed-form analysis, while ladder or lattice implementation of linear filters has
 139 long been studied in signal theory. However, when the system generating the data is driven
 140 by nonlinear dynamics, the model specification and parameter estimation problems increase
 141 their complexity, and hence nonlinear adaptive filtering becomes strictly necessary.

142 Note that, in the field of control theory, a range of sequential algorithms for nonlinear
 143 filtering have been proposed since the sixties, notably the extended Kalman filter [Lewis et al.](#)
 144 [\(2007\)](#) and the unscented Kalman filter [Julier and Uhlmann \(1996\)](#), which are both nonlinear
 145 extensions of the celebrated Kalman filter [Kalman \(1960\)](#), and particle filters [Del Moral](#)
 146 [\(1996\)](#). These methods generally require knowledge of a state-space model, and while some
 147 of them are related to adaptive filtering, we will not deal with them explicitly in this chapter.

148 9.2 Kernel Adaptive Filtering

149 The nonlinear filtering problem and the online adaptation of model weights were first
 150 addressed by neural networks in the nineties [Dorffner \(1996\)](#); [Narendra and Parthasarathy](#)
 151 [\(1990\)](#). Throughout the last decade, a great interest has been devoted to developing nonlinear
 152 versions of linear adaptive filters by means of kernels ([Liu et al. 2010](#)). The goal is to develop
 153 machines that learn over time in changing environments, and at the same time adopt the nice
 154 characteristics of convexity, convergence and reasonable computational complexity, which
 155 was not successfully implemented in neural networks.

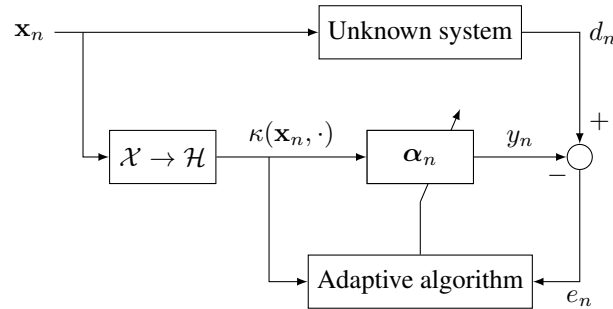


Figure 9.3 A kernel adaptive filter for nonlinear system identification.

156 Kernel adaptive filtering aims to formulate the classic linear adaptive filters in rkHs, such
 157 that a series of convex least-squares problems is solved. Several basic kernel adaptive filters
 158 can be obtained by applying a linear adaptive filter directly on the transformed data, as
 159 illustrated in Fig. 9.3. This requires the reformulation of scalar-product based operations in
 160 terms of *kernel evaluations*. The resulting algorithms typically consist of algebraically simple
 161 expressions, though they feature powerful nonlinear filtering capabilities. Nevertheless, the
 162 design of these online kernel methods requires to deal with some of the challenges that
 163 typically arise when dealing with kernels, such as overfitting and computational complexity
 164 issues.

165 In the sequel we will discuss two families of kernel adaptive filters in detail, namely kernel
 166 least mean squares and kernel recursive least-squares algorithms. Several related kernel
 167 adaptive filters will be reviewed briefly as well.

168 9.3 Kernel Least Mean Squares

169 The early approach to kernel adaptive filtering introduced a kernel version of the celebrated
 170 ADALINE in ([Frieß and Harrison 1999](#)), though this method was not online. Kivinen et al.
 171 proposed an algorithm to perform stochastic gradient descent in rkHs [Kivinen et al. \(2004\)](#);
 172 The so-called Naive Online regularized Risk Minimization Algorithm (NORMA) introduces
 173 a regularized risk that can be solved online and can be shown to be equivalent to a kernel
 174 version of leaky LMS, which itself is a regularized version of LMS.

175 9.3.1 Derivation of KLMS

As an illustrative guiding example of a kernel adaptive filter, we will take the kernelization
 of the standard LMS algorithm, known as Kernel Least Mean Squares (KLMS) ([Liu et al.](#)

2008). The approach employs the traditional kernel trick. Essentially, a nonlinear function $\phi(\cdot)$ maps the data \mathbf{x}_n from the input space to $\phi(\mathbf{x}_n)$ in the feature space. Let $\mathbf{w}_{\mathcal{H}}$ be the weight vector in this space such that the filter output is $y_n = \mathbf{w}_{\mathcal{H},n}^\top \mathbf{x}_n$, where $\mathbf{w}_{\mathcal{H},n}$ is the estimate of $\mathbf{w}_{\mathcal{H}}$ at time instant n . Note that we will be taking scalar products of real-valued vectors from now on. Given a desired response d_n we wish to minimize squared loss, $J_{\mathbf{w}_{\mathcal{H},n}}$, with respect to $\mathbf{w}_{\mathcal{H}}$. Similar to Eq. (9.6), the obtained stochastic gradient descent update rule reads

$$\mathbf{w}_{\mathcal{H},n} = \mathbf{w}_{\mathcal{H},n-1} + \eta e_n \phi(\mathbf{x}_n). \quad (9.19)$$

By initializing the solution as $\mathbf{w}_{\mathcal{H},0} = 0$ (and hence $e_0 = d_0 = 0$), the solution after n iterations can be expressed in closed form as

$$\mathbf{w}_{\mathcal{H},n} = \eta \sum_{i=1}^n e_i \phi(\mathbf{x}_i). \quad (9.20)$$

By exploiting the kernel trick, one obtains the prediction function

$$y_* = \eta \sum_{i=1}^n e_i \phi(\mathbf{x}_i) \phi(\mathbf{x}_*) = \eta \sum_{i=1}^n e_i k(\mathbf{x}_i, \mathbf{x}_*), \quad (9.21)$$

176 where \mathbf{x}_* represents an arbitrary input datum and $k(\cdot, \cdot)$ is the kernel function, for instance
 177 the commonly used Gaussian kernel $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_k)$ with *kernel width*
 178 σ_k . Note that the weights $\mathbf{w}_{\mathcal{H},n}$ of the nonlinear filter are not used explicitly in the KLMS
 179 algorithm. Also, since the present output is determined solely by previous inputs and all
 180 the previous errors, it can be readily computed in the input space. These error samples are
 181 similar to innovation terms in sequential state estimation (Haykin 2001), since they add new
 182 information to improve the output estimate. Each new input sample results in an output, and
 183 hence a corresponding error, which is never modified further and incorporated in the estimate
 184 of the next output. This recursive computation makes KLMS especially useful for online
 185 (adaptive) nonlinear signal processing.

186 In Liu *et al.* (2008) it was shown that the KLMS algorithm is well-posed in rkHs
 187 without the need of an extra regularization term in the finite training data case, because
 188 the solution is always forced to lie in the subspace spanned by the input data. The lack
 189 of an explicit regularization term leads to two important advantages. First of all, it has a
 190 simpler implementation than NORMA, as the update equations are straightforward kernel
 191 versions of the original linear ones. Second, it can potentially provide better results because
 192 regularization biases the optimal solution. In particular, it was shown that a small enough
 193 step-size can provide a sufficient “self-regularization” mechanism. Moreover, since the space
 194 spanned by the mapped samples is possibly infinite-dimensional, the projection error of the
 195 desired signal d_n could be very small, as is well known from Cover’s theorem Haykin (1999).
 196 On the downside, the speed of convergence and the misadjustment also depend upon the
 197 step-size. As a consequence, they conflict with the generalization ability.

198 9.3.2 Implementation challenges and dual formulation

Another important drawback of the KLMS algorithm becomes apparent when analyzing its update Eq. (9.21). In order to make a prediction, the algorithm requires to store all previous errors e_i and all processed input data \mathbf{x}_i , for $i = 1, 2, \dots, n$. In online scenarios where data is continuously being received, the size of the KLMS network will continuously grow, posing

implementation challenges. This becomes even more evident if we recast the weight update from Eq. (9.19) into a more standard filtering formulation, by relying on the Representer theorem Schölkopf *et al.* (2001). This theorem states that the solution $\mathbf{w}_{\mathcal{H},n}$ can be expressed as a linear combination of the transformed input data,

$$\mathbf{w}_{\mathcal{H},n} = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i). \quad (9.22)$$

This allows the prediction function to be written as the familiar kernel expansion

$$y_* = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x}_*). \quad (9.23)$$

The expansion coefficients α_i are called the *dual variables* and the reformulation of the filtering problem in terms of α_i is called the *dual formulation*. The update from Eq. (9.19) now becomes

$$\sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i) = \sum_{i=1}^{n-1} \alpha_i \phi(\mathbf{x}_i) + \eta e_n \phi(\mathbf{x}_n), \quad (9.24)$$

and after multiplying both sides with the new datum $\phi(\mathbf{x}_n)$ and adopting a vector notation, we obtain

$$\boldsymbol{\alpha}_n^\top \mathbf{k}_n = \boldsymbol{\alpha}_{n-1}^\top \mathbf{k}_{n-1} + \eta e_n k_{n,n}, \quad (9.25)$$

where $\boldsymbol{\alpha}_n = [\alpha_1, \alpha_2, \dots, \alpha_n]^\top$, the vector \mathbf{k}_n contains the kernels of the n data and the newest point, $\mathbf{k}_n = [k(\mathbf{x}_1, \mathbf{x}_n), k(\mathbf{x}_2, \mathbf{x}_n), \dots, k(\mathbf{x}_n, \mathbf{x}_n)]$, and $k_{n,n} = k(\mathbf{x}_n, \mathbf{x}_n)$. KLMS resolves this relationship by updating $\boldsymbol{\alpha}_n$ as

$$\boldsymbol{\alpha}_n = \begin{bmatrix} \boldsymbol{\alpha}_{n-1} \\ \eta e_n \end{bmatrix}. \quad (9.26)$$

199 The Matlab code for the complete KLMS training step on a new data pair (\mathbf{x}, d) is displayed
200 in Listing 9.3.

```
201 k = kernel(dict, x, kerneltype, kernelpar); % kernels between dictionary and x
202 y = k'*alpha; % evaluate function output
203 err = d - y; % instantaneous error
204
205 kaf.dict = [kaf.dict; x]; % add base to dictionary
206 kaf.alpha = [kaf.alpha; kaf.eta*err]; % add new coefficient
```

Listing 9.3 Training step of the KLMS algorithm on a new datum (\mathbf{x}, d) .

207 The update (9.26) emphasizes the growing nature of the KLMS network, which precludes
208 its direct implementation in practice. In order to design a practical KLMS algorithm, the
209 number of terms in the kernel expansion (9.23) should stop growing over time. This can be
210 achieved by implementing an *online sparsification technique*, whose aim is to identify terms
211 in the kernel expansion that can be omitted without degrading the solution. We will discuss
212 several different sparsification approaches in Section 9.5.

213 Finally, observe that the computational complexity and memory complexity of the KLMS
214 algorithm are both linear in terms of the number of data it stores, $\mathcal{O}(n)$. Recall that the
215 complexity of the LMS algorithm is also linear, though not in terms of the number of data but
216 in terms of the data *dimension*.

217 9.3.3 Example: Prediction of the Mackey-Glass time series

We demonstrate the online learning capabilities of the KLMS kernel adaptive filter by predicting the Mackey-Glass time series, which is a classic benchmark problem [Liu et al. \(2010\)](#). The Mackey-Glass time series is well-known for its strong non-linearity. It corresponds to a high-dimensional chaotic system, and its output is generated by the following time-delay differential equation:

$$\frac{dx_n}{dn} = -bx_n + \frac{ax_{n-\Delta}}{1 + x_{n-\Delta}^{10}}. \quad (9.27)$$

218 We focus on the sequence with parameters $b = 0.1$, $a = 0.2$ and time delay $\Delta = 30$, better
 219 known as the MG30 time series. The prediction problem consists in predicting the n -th
 220 sample, given all samples of the time series up till the $n - 1$ -th sample.

221 Time-series prediction with kernel adaptive filters is typically performed by considering
 222 a time-delay vector $\mathbf{x}_n = [x_n, x_{n-1}, \dots, x_{n-L+1}]^\top$ as the input and the next sample of the
 223 time series as the desired output, $d_n = x_{n+1}$. This approach casts the prediction problem
 224 into the well-know filtering framework². Prediction of several steps ahead can be obtained
 225 by choosing a prediction horizon $h > 1$, and $d_n = x_{n+h}$. For time series generated by a
 226 deterministic process, a principled tool to find the optimal embedding is Takens' theorem
 227 [Takens \(1981\)](#). In the case of the MG30 time series, Takens' theorem indicates that the
 228 optimal embedding is around $L = 7$ [Van Vaerenbergh et al. \(2012a\)](#).

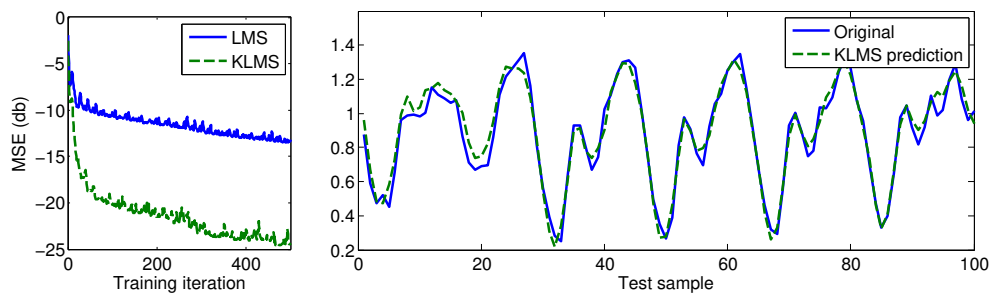


Figure 9.4 KLMS predictions on the Mackey-Glass time series. Left: Learning curve over 500 training iterations. Right: Test samples of the Mackey-Glass time-series and the predictions provided by KLMS.

229 We consider 500 samples for online training of the KLMS algorithm and use the next
 230 100 data for testing. The step size of KLMS is fixed to 0.2, and we use the Gaussian kernel
 231 with $\sigma = 1$. Fig. 9.4 displays the prediction results after 500 training steps. The left plot
 232 shows the learning curve of the algorithm, obtained as the mean squared error (MSE) of the
 233 prediction on the test set, at each iteration of the online training process. As a reference,
 234 we include the learning curve of the linear LMS algorithm with a suitable step size. The
 235 right plot shows the 100 test samples of the original time series, as the full line, and KLMS'
 236 prediction on these test samples after 500 training steps. These predictions are calculated by
 237 evaluating the prediction equation (9.21) on the test samples. The code for this experiment
 238 and all subsequent ones is included in the accompanying material.

²Note that, recently, some authors have proposed to model time-series through a different approach based on explicit recursivity in the rKFs [Li and Príncipe \(2016\)](#); [Tuia et al. \(2014\)](#), as we will see later on in this chapter.

239 9.3.4 Practical KLMS algorithms

240 In the described Mackey-Glass experiment, the KLMS algorithm requires to store 500
 241 coefficients α_i and the 500 corresponding data \mathbf{x}_i . The stored data \mathbf{x}_i are referred to as its
 242 *dictionary*. If the online learning process were to continue indefinitely, the algorithm would
 243 require *evergrowing* memory and computation per time step. This issue has been identified
 244 as a major roadblock early on in the research on kernel adaptive filters, and it has led to the
 245 development of several procedures to slow down the dictionary growth by *sparsifying* the
 246 dictionary.

247 A sparsification procedure based on Gaussian elimination steps on the gram matrix was
 248 proposed in (Pokharel *et al.* 2009). This method is successful in limiting the dictionary size
 249 in the n -th training step to some $m < n$, but in order to do so it requires $\mathcal{O}(m^2)$ complexity,
 250 which defeats the purpose of using a KLMS algorithm.

251 Kernel normalized least-mean squares and the coherence criterion

Around the same time the KLMS algorithm was published, a kernelized version of the Affine
 Projection (AP) algorithm was proposed Richard *et al.* (2009). AP algorithms hold the middle
 ground between LMS and RLS algorithms by calculating an estimate of the correlation matrix
 based on the p last data. For $p = 1$ the algorithm reduces to a kernel version of the normalized
 LMS algorithm Haykin (2001), called Kernel Normalized Least-Mean Squares (KNLMS),
 and its update reads

$$\alpha_n = \begin{bmatrix} \alpha_{n-1} \\ 0 \end{bmatrix} + \frac{\eta}{\epsilon + \|\mathbf{k}_n\|} e_n \mathbf{k}_n. \quad (9.28)$$

252 Note that this algorithm updates all coefficients in each iteration, in contrast to KLMS which
 253 updates just one coefficient.

The kernel affine projection (KAP) and KNLMS algorithms introduced in Richard *et al.*
 (2009) also included an efficient dictionary sparsification procedure, called the *coherence*
criterion. Coherence is a measure to characterize a dictionary in sparse approximation
 problems, defined in a kernel context as

$$\mu = \max_{i \neq j} |k(\mathbf{x}_i, \mathbf{x}_j)|. \quad (9.29)$$

The coherence of a dictionary will be large if it contains two bases \mathbf{x}_i and \mathbf{x}_j that are very
 similar, in terms of the kernel function. Due to their similarity, such bases contribute almost
 identical information to the algorithm, and one of them may be considered redundant. The
 online dictionary sparsification procedure based on coherence operates by only including a
 new datum \mathbf{x}_n into the current dictionary \mathcal{D}_{n-1} if it maintains the dictionary coherence below
 a certain threshold,

$$\max_{j \in \mathcal{D}_{n-1}} |k(\mathbf{x}_n, \mathbf{x}_j)| < \mu_0. \quad (9.30)$$

If the new datum fulfills this criterion, it is included in the dictionary, and the KNLMS
 coefficients are updated through Eq. (9.28). If the coherence criterion (9.30) is not fulfilled,
 the new datum is not included in the dictionary, and a *reduced* update of the KNLMS
 coefficients is performed,

$$\alpha_n = \alpha_{n-1} + \frac{\eta}{\epsilon + \|\mathbf{k}_n\|} e_n \mathbf{k}_n. \quad (9.31)$$

254 This update does not increase the number of coefficients and therefore it maintains the
 255 algorithm's computational complexity fixed during that iteration. The Matlab code for the
 256 complete KNLMS training step on a new data pair (\mathbf{x}, d) is displayed in Listing 9.4.

```

257 k = kernel(dict,x,kerneltype,kernelpar); % kernels between dictionary and x
258 if (max(k) <= mu0), % coherence criterion
259     dict = [dict; x]; % add base to dictionary
260     alpha = [alpha; 0]; % reserve spot for new coefficient
261 end
262
263 k = kernel(dict,x,kerneltype,kernelpar); % kernels with new dictionary
264 y = k'*alpha; % evaluate function output
265 err = d - y; % instantaneous error
266
267 alpha = alpha + eta/(eps + k'*k)*err*k; % update coefficients

```

Listing 9.4 Training step of the KNLMS algorithm on a new datum (x, d) .

268 The coherence criterion is computationally efficient in that it has a complexity that does
 269 not exceed that of the kernel adaptive filter itself, and it has demonstrated to be successful in
 270 practical situations [Van Vaerenbergh and Santamaría \(2013\)](#).

271 Quantized kernel least-mean squares

Recently, a kernel LMS algorithm was proposed that combines elements from the original KLMS algorithm and the coherence criterion, called Quantized Kernel Least Mean Squares (QKLMS) [Chen et al. \(2012\)](#). In particular, when the sparsification criterion decides to include a datum into the dictionary, the algorithm updates its coefficients as follows:

$$\alpha_n = \begin{bmatrix} \alpha_{n-1} \\ \eta e_n \end{bmatrix}. \quad (9.32)$$

When the datum does not fulfil the coherence criterion, it is not included in the dictionary. Instead, the closest dictionary element is retrieved, and the corresponding coefficient is updated as follows

$$\alpha_{n,j} = \alpha_{n-1,j} + \eta e_n, \quad (9.33)$$

272 where j is the dictionary index of the element that is closest. Though conceptually very
 273 simple, this algorithm obtains state-of-the-art results in several applications when only a low
 274 computational budget is available. The Matlab code for the complete QKLMS training step
 275 on a new data pair (x, d) is displayed in Listing 9.5.

```

276 k = kernel(dict,x,kerneltype,kernelpar); % kernels between dictionary and x
277 y = k'*alpha; % evaluate function output
278 err = d - y; % instantaneous error
279
280 [d2,j] = min(sum((dict - repmat(x,m,1)).^2,2)); % distance to dictionary
281 if d2 <= epsu^2,
282     alpha(j) = alpha(j) + eta*err; % reduced coefficient update
283 else
284     dict = [dict; x]; % add base to dictionary
285     alpha = [alpha; eta*err]; % add new coefficient
286 end

```

Listing 9.5 Training step of the QKLMS algorithm on a new datum (x, d) .

287 9.4 Kernel recursive least squares

288 In linear adaptive filtering, the RLS algorithm represents an alternative to LMS, with faster
 289 convergence and typically lower bias, at the expense of a higher computational complexity.

290 RLS is obtained by designing a recursive solution to the least-squares problem. Analogously,
 291 a recursive solution can be designed for the kernel ridge regression problem, yielding kernel
 292 recursive least-squares (KRLS) algorithms.

293 9.4.1 Kernel ridge regression

In order to obtain the kernel-based version of the regularized least-squares cost function (9.7), we first transform the data into the kernel feature space,

$$\begin{aligned} J_n &= \sum_{i=1}^n |d_i - \phi(\mathbf{x}_i)^\top \mathbf{w}_{\mathcal{H}}|^2 + \delta \mathbf{w}_{\mathcal{H}}^\top \mathbf{w}_{\mathcal{H}} \\ &= \|\mathbf{d} - \mathbf{K}\boldsymbol{\alpha}\|^2 + \delta \boldsymbol{\alpha}^\top \mathbf{K}\boldsymbol{\alpha}, \end{aligned} \quad (9.34)$$

where we have applied the kernel trick to obtain the second equality. Here, vector \mathbf{d} contains the n desired values, $\mathbf{d} = [d_1, d_2, \dots, d_n]^\top$, and \mathbf{K} is the kernel matrix with elements $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$. Eq. (9.34) represents the *kernel ridge regression* problem [Saunders et al. \(1998\)](#), and its solution is given by

$$\boldsymbol{\alpha} = (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{d}. \quad (9.35)$$

The prediction for a new datum \mathbf{x}_* is obtained as

$$y_* = \mathbf{k}_*^\top \boldsymbol{\alpha} = \mathbf{k}_*^\top (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{d}. \quad (9.36)$$

294 9.4.2 Derivation of KRLS

The KRLS algorithm [Engel et al. \(2004\)](#) formulates a recursive procedure to obtain the solution of the regression problem (9.34) in the absence of regularization, $\delta = 0$. Without regularization, the solution (9.35) reads

$$\boldsymbol{\alpha} = \mathbf{K}^{-1} \mathbf{d}. \quad (9.37)$$

295 KRLS guarantees the invertibility of the kernel matrix \mathbf{K} by excluding those data \mathbf{x}_i from
 296 from the dictionary that are linearly dependent on the already included data, in the feature
 297 space. As we will see, this is achieved by applying a specific online sparsification procedure,
 298 which guarantees both that \mathbf{K} is invertible and that the algorithm's dictionary stays compact.

Assume the solution after processing $n - 1$ data is available, given by

$$\boldsymbol{\alpha}_{n-1} = \mathbf{K}_{n-1}^{-1} \mathbf{d}_{n-1}, \quad (9.38)$$

In the next iteration, n , a new data pair (\mathbf{x}_n, d_n) is received and we wish to obtain the new solution $\boldsymbol{\alpha}_n$ by applying a low-complexity update on the previous solution (9.38). We first calculate the predicted output

$$y_n = \mathbf{k}_n^\top \boldsymbol{\alpha}_{n-1}, \quad (9.39)$$

and we obtain the a-priori error for this datum, $e_n = d_n - y_n$. The updated kernel matrix can be written as

$$\mathbf{K}_n = \begin{bmatrix} \mathbf{K}_{n-1} & \mathbf{k}_n \\ \mathbf{k}_n^\top & k_{nn} \end{bmatrix}. \quad (9.40)$$

By introducing the variables

$$\mathbf{a}_n = \mathbf{K}_{n-1}^{-1} \mathbf{k}_n, \quad (9.41)$$

and

$$\gamma_n = k_{nn} - \mathbf{k}_n^\top \mathbf{a}_n, \quad (9.42)$$

the update for the inverse kernel matrix can be written as

$$\mathbf{K}_n^{-1} = \frac{1}{\gamma_n} \begin{bmatrix} \gamma_n \mathbf{K}_{n-1}^{-1} + \mathbf{a}_n \mathbf{a}_n^\top & -\mathbf{a}_n \\ -\mathbf{a}_n & 1 \end{bmatrix}. \quad (9.43)$$

Eq. (9.43) is obtained by applying the Sherman-Morrison-Woodbury formula for matrix inversion, see for instance [Golub and Van Loan \(2012\)](#). Finally, the updated solution $\boldsymbol{\alpha}_n$ is obtained as

$$\boldsymbol{\alpha}_n = \begin{bmatrix} \boldsymbol{\alpha}_{n-1} \\ 0 \end{bmatrix} - e_n / \gamma_n \begin{bmatrix} \mathbf{a}_n \\ -1 \end{bmatrix}. \quad (9.44)$$

299 Equations (9.43) and (9.44) are efficient updates that allow to obtain the new solution in
 300 $\mathcal{O}(n^2)$ time and memory, based on the previous solution. Directly applying Eq. (9.37)
 301 at iteration n would require $\mathcal{O}(n^3)$ cost, so the recursive procedure is preferred in
 302 online scenarios. A detailed derivation of this result can be found in [Engel et al. \(2004\)](#);
 303 [Van Vaerenbergh et al. \(2012b\)](#).

304 Online sparsification by approximate linear dependency

305 The KRLS algorithm from [Engel et al. \(2004\)](#) follows the described recursive solution. In
 306 order to slow down the dictionary growth, shown in Eqs. (9.43) and (9.44), it introduces a
 307 sparsification criterion based on *approximate linear dependency* (ALD). According to this
 308 criterion, a new datum \mathbf{x}_n should only be included in the dictionary if $\phi(\mathbf{x}_n)$ cannot be
 309 approximated sufficiently well in feature space by a *linear combination* of the already present
 310 data.

Given a dictionary \mathcal{D} of data \mathbf{x}_j and a new training point \mathbf{x}_n , we need to find a set of coefficients $\mathbf{a} = [a_1, a_2, \dots, a_m]^\top$ that satisfy the approximate linear dependency condition

$$\min_{\mathbf{a}} \left\| \sum_{j=1}^m a_j \phi(\mathbf{x}_j) - \phi(\mathbf{x}_n) \right\|^2 \leq \nu \quad (9.45)$$

where m is the cardinality of the dictionary. Interestingly, it can be shown that these coefficients are already calculated by the KRLS update itself, and they are available at each iteration n as $\mathbf{a}_n = \mathbf{K}_{n-1}^{-1} \mathbf{k}_n$, see Eq. (9.41). The ALD condition can therefore be verified by simply comparing γ_n to the ALD threshold,

$$\gamma_n = k_{nn} - \mathbf{k}_n^\top \mathbf{a}_n \leq \nu. \quad (9.46)$$

311 If $\gamma_n > \nu$, then we must add the newest datum \mathbf{x}_n to the dictionary, $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{\mathbf{x}_n\}$,
 312 before updating the solution through Eq.(9.44). If $\gamma_n \leq \nu$ then the datum \mathbf{x}_n is already
 313 represented sufficiently well by the dictionary. In this case the dictionary is not expanded,
 314 $\mathcal{D}_n = \mathcal{D}_{n-1}$, and a *reduced* update of the solution is performed, see [Engel et al. \(2004\)](#) for
 315 details. The Matlab code for the complete KRLS training step on a new data pair (\mathbf{x}, d) is
 316 displayed in Listing 9.6.

```

317 k = kernel(dict,x,kerneltype,kernelpar); % kernels between dictionary and x
318 kxx = kernel(x,x,kaf.kerneltype,kaf.kernelpar); % kernel on x
319
320 a = Kinv*k; % coefficients of closest linear combination in feature space
321 gamma = kxx - k'*a; % residual of linear approximation in feature space
322
323 y = k'*alpha; % evaluate function output
324 err = d - y; % instantaneous error
325
326 if gamma>nu % new datum is not approximately linear dependent
327     dict = [dict; x]; % add base to dictionary
328     Kinv = 1/gamma*[gamma*Kinv+a*a',-a;-a',1]; % update inv. kernel matrix
329     Z = zeros(size(P,1),1);
330     P = [P Z; Z' 1]; % add linear combination coeff. to projection matrix
331     ode = 1/gamma*err;
332     alpha = [alpha - a*ode; ode]; % full update of coefficients
333 else % perform reduced update of alpha
334     q = P*a/(1+a'*P*a);
335     P = P - q*(a'*P); % update projection matrix
336     alpha = alpha + Kinv*q*err; % reduced update of coefficients
337 end

```

Listing 9.6 Training step of the KRLS algorithm on a new datum (x, d) .

9.4.3 Prediction of the Mackey-Glass time series with KRLS

The update equations for KRLS require substantially more computation than KLMS. In particular, KRLS has quadratic complexity, $\mathcal{O}(m^2)$, in terms of its dictionary size, and KLMS has linear complexity, $\mathcal{O}(m)$. On the other hand, KRLS has faster convergence and lower bias. We illustrate these properties by applying KRLS on the Mackey-Glass prediction experiment from Section 9.3.3.

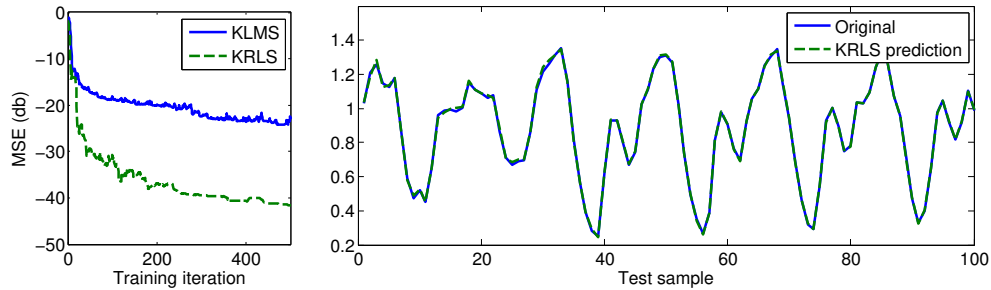


Figure 9.5 KRLS predictions on the Mackey-Glass time series. Left: Learning curve over 500 training iterations, including comparison to KLMS. Right: test samples and KRLS predictions.

Fig. 9.5 shows the results of training the KRLS algorithm on the Mackey-Glass time series. KRLS is applied with a Gaussian kernel with $\sigma_k = 1$, and its precision parameter was fixed to $\nu = 10^{-4}$. The left plot compares the learning curves of KLMS and KRLS, demonstrating a slightly faster initial convergence rate for KRLS, after which the algorithm converges to a much lower MSE than KLMS. The low bias is also visible in the right plot, which shows the prediction results on the test data.

350 9.4.4 Beyond the stationary model

351 One important limitation of the KRLS algorithm is that it always assumes a stationary model,
 352 and therefore it cannot track changes in the true underlying data model. This is a somewhat
 353 odd property for an *adaptive* filter, though note that this is also the case for the original RLS
 354 algorithm, see Section 9.1.2.

355 In order to enable tracking and make a truly adaptive KRLS algorithm, several
 356 modifications have been presented in the literature. An exponentially-weighted KRLS
 357 algorithm was proposed by including a forgetting factor, and an extended KRLS algorithm
 358 was designed by assuming a simple state-space model, though both algorithms show
 359 numerical instabilities in practice Liu *et al.* (2010). In the sequel we briefly discuss two
 360 different approaches that successfully allow KRLS to adapt to changing environments.

361 Sliding-Window KRLS

362 The KRLS algorithm summarizes past information into a compact formulation that does
 363 not allow easy manipulation. For instance, there does not exist a straightforward manner
 364 to include a forgetting factor to exclude the influence of older data.

365 In Van Vaerenbergh *et al.* (2006), a sliding-window based version of KRLS was proposed,
 366 called Sliding-Window Kernel Recursive Least-Squares (SW-KRLS). This algorithm stores
 367 a window of the last m data as its dictionary, and once a datum is older than m time
 368 steps it is simply discarded. In each step the algorithm adds the new datum and discards
 369 the oldest datum, leading to a sliding-window approach. The algorithm stores the inverse
 370 regularized kernel matrix, $(\mathbf{K}_n + \delta\mathbf{I})^{-1}$, calculated on its current dictionary, and a vector of
 371 the corresponding desired outputs, \mathbf{d}_n . By storing these variables it can calculate the solution
 372 vector by simply evaluating $\boldsymbol{\alpha}_n = (\mathbf{K}_n + \delta\mathbf{I})^{-1}\mathbf{d}_n$, see Eqs. (9.35) and (9.36).

The inverse kernel matrix is updated in two steps: First, the new datum is added, which
 requires expanding the matrix with one row and one column. This is carried out by performing
 the operation from Eq. (9.43), similar as in the KRLS algorithm. Second, the oldest datum is
 discarded, which requires removing one row and column from the inverse kernel matrix. This
 can be achieved by writing the kernel matrix and its inverse as follows,

$$\mathbf{K}_{n-1} = \begin{bmatrix} a & \mathbf{b}^T \\ \mathbf{b} & \mathbf{D} \end{bmatrix}, \quad \mathbf{K}_{n-1}^{-1} = \begin{bmatrix} e & \mathbf{f}^T \\ \mathbf{f} & \mathbf{G} \end{bmatrix}, \quad (9.47)$$

after which the inverse (regularized) kernel matrix is found as

$$\mathbf{D}^{-1} = \mathbf{G} - \mathbf{f}\mathbf{f}^T/e. \quad (9.48)$$

373 Details can be found in Van Vaerenbergh *et al.* (2006). Figure 9.6 illustrates the kernel matrix
 374 updates when using a sliding window, compared to the classical growing-window approach of
 375 KRLS. The Matlab code for the SW-KRLS training step on a new data pair (\mathbf{x}, d) is displayed
 376 in Listing 9.7.

```

377 dict = [dict; x]; % add base to dictionary
378 dict_d = [dict_d; d]; % add d to output dictionary
379 k = kernel(dict, x, kerneltype, kernelpar); % kernels between dictionary and x
380 Kinv = grow_kernel_matrix(Kinv, k, c); % calculate new inverse kernel matrix
381
382 if (size(dict,1) > M) % prune
383     dict(1, :) = []; % remove oldest base from dictionary
384     dict_d(1) = []; % remove oldest d from output dictionary
  
```

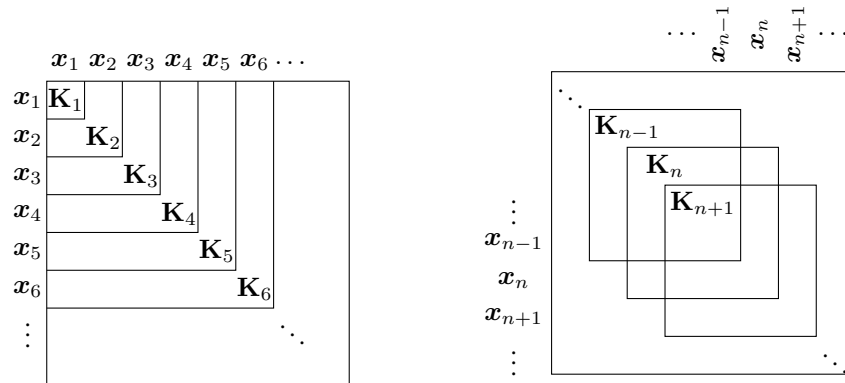


Figure 9.6 Different forms of updating the kernel matrix during online learning. In KRLS-type algorithms the update involves calculating the inverse of each kernel matrix, given the inverse of the previous matrix. Left: Growing kernel matrix, as constructed in KRLS (omitting sparsification for simplicity). Right: Sliding-window kernel matrix of a fixed size, as constructed in SW-KRLS.

```

385     Kinv = prune_kernel_matrix(Kinv); % prune inverse kernel matrix
386 end
387
388 alpha = Kinv*dict_d; % obtain new filter coefficients

```

Listing 9.7 Training step of the SW-KRLS algorithm on a new datum (\mathbf{x}, d) . The functions “grow_kernel_matrix” and “prune_kernel_matrix” implement the operations (9.43) and (9.48).

SW-KRLS is a conceptually very simple algorithm that obtains reasonable performance in a wide range of scenarios, most notably in non-stationary environments. Nevertheless, its performance is limited by the quality of the bases in its dictionary, over which it has little control. In particular, it has no means to avoid redundancy in its dictionary or to maintain older bases that are relevant to its kernel expansion. In order to improve this performance, a Fixed-Budget KRLS (FB-KRLS) algorithm was proposed in Van Vaerenbergh *et al.* (2010). Instead of discarding the oldest data point in each iteration, FB-KRLS discards the data point that causes the least error upon being discarded, using a least a-posteriori error based pruning criterion we will discuss in Section 9.5. In stationary scenarios, FB-KRLS obtains significantly better results.

401 KRLS Tracker

402 The tracking limitations of previous KRLS algorithms were overcome by development of
 403 the Kernel Recursive Least-Squares Tracker (KRLS-T) algorithm Van Vaerenbergh *et al.*
 404 (2012b), which has its roots in the probabilistic theory of Gaussian Process (GP) regression.
 405 Similar to the FB-KRLS algorithm, this algorithm uses a fixed memory size and has a
 406 criterion to discard which data to discard in each iteration. But unlike FB-KRLS, the KRLS-T
 407 algorithm incorporates a forgetting mechanism to gradually downweigh older data.

408 We will provide more details on this algorithm in the discussion on probabilistic kernel
 409 adaptive filtering of Section 9.6. As a reference, we list the Matlab code for the KRLS-T
 410 training step on a new data pair (\mathbf{x}, d) in Listing 9.8. While this algorithm has similar
 411 complexity as other KRLS-type algorithms, its implementation is more complex due its
 412 fully probabilistic treatment of the regression problem. Note that some additional checks to

413 avoid numerical problems have been left out; The complete code can be found in the Kernel
414 Adaptive Filtering toolbox, discussed in Section 9.8.

```

415 % perform one forgetting step
416 Sigma = lambda*Sigma + (1-lambda)*K; % forgetting on covariance matrix
417 mu = sqrt(lambda)*mu; % forgetting on mean vector
418
419 k = kernel(dict,x,kerneltype,kernelpar); % kernels between dictionary and x
420 kxx = kernel(x,x,kaf.kerneltype,kaf.kernelpar); % kernel on x
421
422 q = Q*k;
423 y_mean = q'*mu; % predictive mean of new datum
424 gamma2 = kxx - k'*q; % projection uncertainty
425 h = Sigma*q;
426 sf2 = gamma2 + q'*h; % noiseless prediction variance
427 sy2 = sn2 + sf2; % unscaled predictive variance of new datum
428 y_var = s02*sy2; % predictive variance of new datum
429
430 % include new sample and add a basis
431 Qold = Q; % old inverse kernel matrix
432 p = [q; -1];
433 Q = [Q zeros(m,1); zeros(1,m) 0] + 1/gamma2*(p*p'); % updated inverse matrix
434
435 err = d - y_mean; % instantaneous error
436 p = [h; sf2];
437 mu = [mu; y_mean] + err/sy2*p; % posterior mean
438 Sigma = [Sigma h; h' sf2] - 1/sy2*(p*p'); % posterior covariance
439 dict = [dict; x]; % add base to dictionary
440
441 % estimate scaling power s02 via ML
442 nums02ML = nums02ML + lambda*(y - y_mean)^2/sy2;
443 dens02ML = dens02ML + lambda;
444 s02 = nums02ML/dens02ML;
445
446 % delete a basis if necessary
447 m = size(dict,1);
448 if m>M
449     % MSE pruning criterion
450     errors = (Q*mu)./diag(Q);
451     criterion = abs(errors);
452
453     [~, r] = min(criterion); % remove element which incurs in the min. err.
454     smaller = 1:m; smaller(r) = [];
455
456     if r == m, % remove the element we just added (perform reduced update)
457         Q = Qold;
458     else
459         Qs = Q(smaller, r);
460         qs = Q(r,r); Q = Q(smaller, smaller);
461         Q = Q - (Qs*Qs')/qs; % prune inverse kernel matrix
462     end
463     mu = mu(smaller); % prune posterior mean
464     Sigma = Sigma(smaller, smaller); % prune posterior covariance
465     dict = dict(smaller,:); % prune dictionary
466 end

```

Listing 9.8 Training step of the KRLS-T algorithm on a new datum (x, d) .

9.4.5 Example: Nonlinear channel identification and reconvergence

In order to demonstrate the tracking capabilities of some of the reviewed kernel adaptive filters we perform an experiment similar to the setup described in [Lázaro-Gredilla et al. \(2011\)](#); [Van Vaerenbergh et al. \(2006\)](#). Specifically, we consider the problem of online identification of a communication channel in which an abrupt change (switch) is triggered at some point.

A signal $x_n \in \mathcal{N}(0, 1)$ is fed into a nonlinear channel that consists of a linear finite impulse response (FIR) channel followed by the nonlinearity $y = \tanh(z)$, where z is the output of the linear channel. During the first 500 iterations the impulse response of the linear channel is chosen as $\mathcal{H}_1 = [1, -0.3817, -0.1411, 0.5789, 0.191]$, and at iteration 501 it is switched to $\mathcal{H}_2 = [1, -0.0870, 0.9852, -0.2826, -0.1711]$. Finally, 20dB of Gaussian white noise is added to the channel output.

We perform an online identification experiment with the algorithms LMS, QKLMS, SW-KRLS, and KRLS-T. Each algorithm performs online learning of the nonlinear channel, processing one input datum (with a time-embedding of 5 taps) and one output sample per iteration. At each step, the MSE performance is measured on a set of 100 data points that are generated with the current channel model. The results are averaged out over 10 simulations.

The kernel adaptive filters use a Gaussian kernel with $\sigma_k = 1$. LMS and QKLMS use a learning rate $\eta = 0.5$. The sparsification threshold of QKLMS is set to $\epsilon_{\text{U}} = 0.3$, which leads to a final dictionary of size around $m = 300$ at the end of the experiment. The regularization of SW-KRLS and KRLS-T is set to match the true value of the noise-to-signal ratio, 0.01. Regarding memory, SW-KRLS and KRLS-T are given a maximum dictionary size of $m = 50$. Finally, KRLS-T uses a forgetting factor of $\lambda = 0.998$.

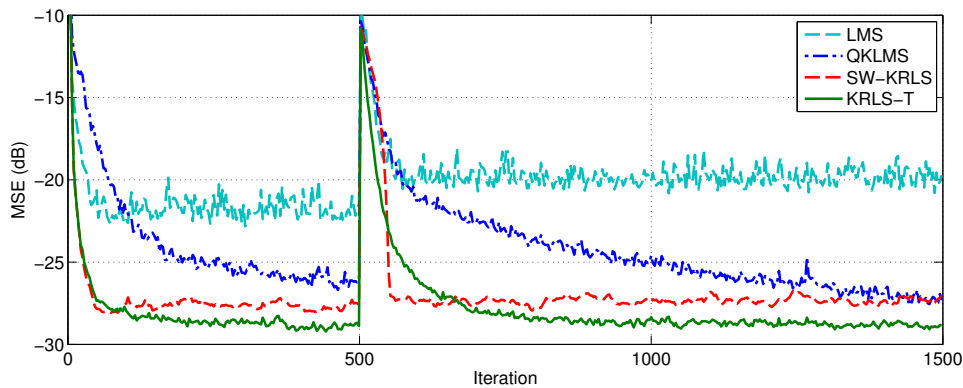


Figure 9.7 MSE learning curves of different kernel adaptive filters on a communications channel that shows an abrupt change at iteration 500.

The results are shown in Fig. 9.7. LMS performs worst, as it is not capable of modeling the nonlinearities in the system. QKLMS shows good results, given its low complexity, but a slow convergence. SW-KRLS and KRLS-T converge to a value which is mostly limited by its dictionary size, $m = 50$, and both show fast convergence rates. All algorithms are capable of reconverging after the switch, though their convergence rate is typically slower at that point.

9.5 Online Sparsification with Kernels

The idea behind sparsification methods is to construct a sparse dictionary of bases that represent the remaining data sufficiently well. As a general rule in learning theory, it is desirable to design a network with as few processing elements as possible. Sparsity reduces the complexity in terms of computation and memory, and it usually gives better generalization ability to unseen data Platt (1991); Vapnik (1995). In the context of kernel methods, sparsification aims to identify the bases in the kernel expansion $y_* = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x}_*)$, see Eq (9.23), that can be discarded without incurring a significant performance loss.

Online sparsification is typically performed by starting with an empty dictionary, $\mathcal{D}_0 = \emptyset$, and, in each iteration, adding the input datum \mathbf{x}_i if it fulfills a chosen sparsification criterion. We denote the dictionary at time instant $n - 1$ as $\mathcal{D}_{n-1} = \{\mathbf{c}_i\}_{i=1}^{m_{n-1}}$, where \mathbf{c}_i is the i -th stored center, taken from the input data \mathbf{x} received up till this instant, and m_{n-1} is the dictionary cardinality at this instant. When a new input-output pair (\mathbf{x}_n, d_n) is received, a decision is made whether or not \mathbf{x}_n should be added to the dictionary as a center. If the sparsification criterion is fulfilled, \mathbf{x}_n is added to the dictionary, $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{\mathbf{x}_n\}$. If the criterion is not fulfilled, the dictionary is maintained, $\mathcal{D}_n = \mathcal{D}_{n-1}$, to preserve its sparsity.

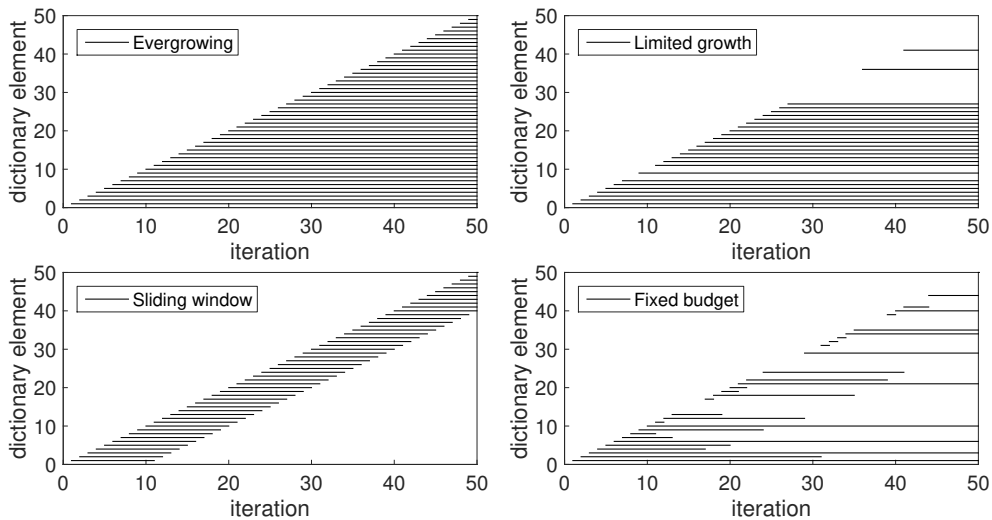


Figure 9.8 Dictionary construction processes for different sparsification approaches. Each horizontal line marks the presence of a center in the dictionary. Top left: The evergrowing dictionary construction, in which the dictionary contains n elements in iteration n ; Top right: Online sparsification by slowing down the dictionary growth, as obtained by the coherence and ALD criteria; Bottom left: Sliding-window approach, displayed with 10 elements in the dictionary; Bottom right: Fixed-budget approach, in which the pruning criterion discards one element per iteration, displayed with dictionary size 10.

Fig. 9.8 illustrates the dictionary construction process for different sparsification approaches. Each horizontal line represents the presence of a center in the dictionary. At any given iteration, the elements in the dictionary are indicated by the horizontal lines that are present at that iteration. Below we discuss each approach in detail.

515

9.5.1 Sparsity by construction

516 We will first give an general overview of the different online sparsification methods in the
 517 literature, some of which we have already introduced in the context of the algorithms for
 518 which they were proposed. We distinguish three criteria that achieve sparsity by construction:
 519 novelty criterion, approximate linear dependency criterion, and coherence criterion. If the
 520 dictionary is not allowed to grow beyond a specified maximum size, it may be necessary to
 521 discard bases at some point. This process is referred to as *pruning*, and we will review the
 522 most important pruning criteria later.

523 *Novelty Criterion.* The novelty criterion is a data selection method introduced by Platt
 524 Platt (1991). It was used to construct resource allocating networks (RAN), which
 525 are essentially growing radial basis function networks. When a new data point \mathbf{x}_n is
 526 obtained by the network, the novelty criterion calculates the distance of this point to the
 527 current dictionary, $\min_{j \in \mathcal{D}_{n-1}} \|\mathbf{x}_n - \mathbf{c}_j\|$. If this distance is smaller than some preset
 528 threshold, \mathbf{x}_n is added to the dictionary. Otherwise, it computes the prediction error,
 529 and only if this error e_n is larger than another preset threshold, the datum \mathbf{x}_n will be
 530 accepted as a new center.

Approximate Linear Dependency Criterion. A more sophisticated dictionary growth
 criterion was introduced for the KRLS algorithm in Engel *et al.* (2004): Each time
 a new datum \mathbf{x}_n is observed, the approximate linear dependency (ALD) criterion
 measures how well the datum can be approximated in the feature space as a linear
 combination of the dictionary bases in that space. It does so by checking if the ALD
 condition holds, see Eq. (9.45),

$$\min_{\mathbf{a}} \left\| \sum_{j=1}^m a_j \phi(\mathbf{c}_j) - \phi(\mathbf{x}_n) \right\|^2 \leq \nu.$$

531 Evaluating the ALD criterion requires quadratic complexity, $\mathcal{O}(m^2)$ and therefore it is
 532 not suitable for algorithms with linear complexity such as KLMS.

Coherence Criterion. The coherence criterion is a straightforward criterion to check whether
 the newly arriving datum is sufficiently informative. It was introduced in the context
 of the KNLMS algorithm Richard *et al.* (2009). Given the dictionary D_{n-1} at iteration
 $n - 1$ and the newly arriving datum \mathbf{x}_n , the coherence criterion to include the datum
 reads

$$\max_{j \in \mathcal{D}_{n-1}} |k(\mathbf{x}_n, \mathbf{c}_j)| < \mu_0. \quad (9.49)$$

533 In essence, the coherence criterion checks the similarity, as measured by the kernel
 534 function, between the new datum and the most similar dictionary center. Only if this
 535 similarity is below a certain threshold μ_0 , the datum is inserted into the dictionary. The
 536 higher the threshold μ_0 is chosen, the more data will be accepted in the dictionary. It is
 537 an effective criterion that has linear computational complexity in each iteration: it only
 538 requires to calculate m kernel functions, making it suitable for KLMS-type algorithms.

539 In Chen *et al.* (2012) a similar criterion was introduced, $\min_{j \in \mathcal{D}_{n-1}} \|\mathbf{x}_n - \mathbf{c}_j\| > \epsilon_u$,
 540 which is essentially equivalent to the coherence criterion with a Euclidean distance
 541 based kernel.

542 9.5.2 Sparsity by pruning

543 In practice, it is often necessary to specify a maximum dictionary size m , or *budget*, that may
 544 not be exceeded, for instance due to limitations on hardware or execution time. In order to
 545 avoid exceeding this budget, one could simply stop including any data in the dictionary once
 546 the budget is reached, hence *locking* the dictionary. Nevertheless, it is very probable that at
 547 some point after locking the dictionary a new datum is received that is very informative. In
 548 this case, the quality of the algorithm's solution may improve by pruning the least relevant
 549 center of the dictionary and replacing it with the new, more informative datum.

550 The goal of a pruning criterion is to select a datum out of a given set, such that the
 551 algorithm's performance is least affected. This makes pruning criteria conceptually different
 552 from the previously discussed online sparsification criteria, whose goal is to decide whether
 553 or not to include a datum. Pruning techniques have been studied in the context of neural
 554 network design [Hassibi et al. \(1993\)](#); [LeCun et al. \(1989\)](#) and kernel methods [De Kruif and
 555 De Vries \(2003\)](#); [Hoegaerts et al. \(2004\)](#). We briefly discuss the two most important pruning
 556 criteria that appear in kernel adaptive filtering: sliding-window criterion and error criterion.

557 *Sliding Window.* In time-varying environments, it may be useful to discard the oldest bases,
 558 as these were observed when the underlying model was most different from the current
 559 model. This strategy is at the core of sliding-window algorithms such as NORMA
 560 [Kivinen et al. \(2004\)](#) and SW-KRLS [Van Vaerenbergh et al. \(2006\)](#). In every iteration,
 561 these algorithms include the new datum in the dictionary and discard the oldest datum,
 562 thereby maintaining a dictionary of fixed size.

Error Criterion. Instead of simply discarding the oldest datum, the error based criterion
 determines the datum that will cause the least increase of the squared-error
 performance after it is pruned. This is a more sophisticated pruning strategy that was
 introduced in [Csató and Opper \(2002\)](#); [De Kruif and De Vries \(2003\)](#) and requires
 quadratic complexity to evaluate, $\mathcal{O}(m^2)$. Interestingly, if the inverse kernel matrix is
 available, it is straightforward to evaluate this criterion. Given the i -th element on the
 diagonal of the inverse kernel matrix, $[\mathbf{K}^{-1}]_{ii}$, and the i -th expansion coefficient α_i , the
 squared error after pruning the i -th center from a dictionary is $\alpha_i/[\mathbf{K}^{-1}]_{ii}$. The error
 based pruning criterion therefore selects the index for which this quantity is minimized,

$$\arg \min_i \frac{\alpha_i}{[\mathbf{K}^{-1}]_{ii}}. \quad (9.50)$$

563 This criterion is used in the fixed-budget algorithms FB-KRLS [Van Vaerenbergh et al.
 564 \(2010\)](#) and KRLS-T [Van Vaerenbergh et al. \(2012b\)](#). An analysis performed in [Lázaro-
 565 Gredilla et al. \(2011\)](#) shows that the results obtained by this criterion are very close
 566 to the optimal approach, which is based on minimization of the Kullback–Leibler
 567 divergence between the original and the approximate posterior distributions.

Currently, the most successful pruning criteria used in the kernel adaptive filtering literature
 have quadratic complexity, $\mathcal{O}(m^2)$ and therefore they can only be used in KRLS-type
 algorithms. Optimal pruning in KLMS is a particularly challenging problem, as it is hard
 to define a pruning criterion that can be evaluated with linear computational complexity. A
 simple criterion is found in [Rzepka \(2012\)](#), where the center with the least weight is pruned,
 and weight is determined by the associated expansion coefficient,

$$\arg \min_i |\alpha_i|. \quad (9.51)$$

568 The design of more sophisticated pruning strategies is currently an open topic in KLMS
 569 literature. Some recently proposed criteria can be found in [Zhao *et al.* \(2013, 2016\)](#).

570 9.6 Probabilistic Approaches to Kernel Adaptive Filtering

571 In many signal processing applications, the problem of signal estimation is addressed.
 572 Probabilistic models have proven to be very useful in this context [Arulampalam *et al.* \(2002\)](#);
 573 [Rabiner \(1989\)](#). One of the advantages of probabilistic approaches is that they force the
 574 designer to specify all the prior assumptions of the model, and that they make a clear
 575 separation between the model and the applied algorithm. Another benefit is that they typically
 576 provide a measure of uncertainty about the estimation. Such an uncertainty estimate is not
 577 provided by classical kernel adaptive filtering algorithms, which produce a point estimate
 578 without any further guarantees.

579 In this section, we will review how the probabilistic framework of Gaussian Processes
 580 (GP) allows to extend kernel adaptive filters to probabilistic methods. The resulting GP-based
 581 algorithms not only produce an estimate of an unknown function, but an entire probability
 582 distribution over functions, see [Fig. 9.9](#).

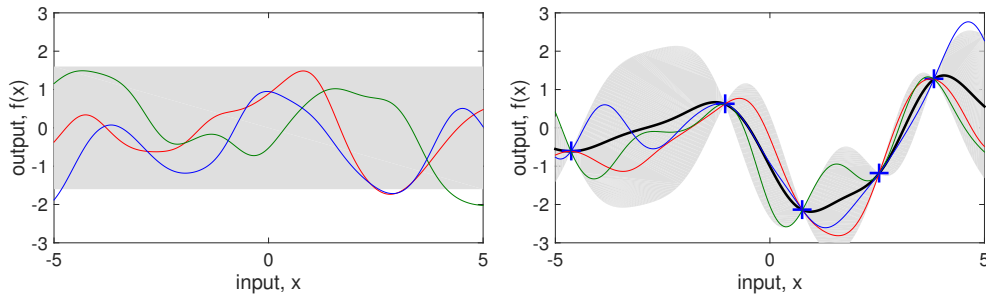


Figure 9.9 Functions drawn from a Gaussian process with a squared exponential covariance $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2/2\sigma_k^2)$. The 95% confidence interval is plotted as the shaded area. Left: Draws from the prior function distribution. Right: Draws from the posterior function distribution, which is obtained after 5 data points (blue crosses) are observed. The predictive mean is displayed in black.

583 Before we describe any probabilistic kernel adaptive filtering algorithms, it is instructive to
 584 take a step back to the non-adaptive setting, and consider the kernel ridge regression problem
 585 [\(9.34\)](#). We will adopt the GP framework to analyze this problem from a probabilistic point of
 586 view.

587 9.6.1 Gaussian Processes and Kernel Ridge Regression

Let us assume that the observed data in a regression problem can be described by the following model,

$$d_n = f(\mathbf{x}_n) + \epsilon_n, \quad (9.52)$$

in which f represents an unobservable *latent function* and $\epsilon_n \sim \mathcal{N}(0, \sigma^2)$ is zero-mean Gaussian noise. We will furthermore assume a zero-mean GP prior on $f(\mathbf{x})$

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')), \quad (9.53)$$

and a Gaussian prior on the noise ϵ ,

$$\epsilon \sim \mathcal{N}(0, \sigma^2). \quad (9.54)$$

588 In the GP literature, the kernel function $k(\mathbf{x}, \mathbf{x}')$ is referred to as the *covariance*, since it
 589 specifies the a-priori relationship between values $f(\mathbf{x})$ and $f(\mathbf{x}')$ in terms of their respective
 590 locations, and its parameters are called *hyperparameters*.

By definition, the marginal distribution of a GP at a finite set of points is a joint Gaussian distribution, with its mean and covariance being specified by the functions $m(\mathbf{x})$ and $k(\mathbf{x}, \mathbf{x}')$ evaluated at those points [Rasmussen and Williams \(2006\)](#). Thus, the joint distribution of outputs $\mathbf{d} = [d_1, \dots, d_n]^\top$ and the corresponding latent vector $\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)]^\top$ is

$$\begin{bmatrix} \mathbf{d} \\ \mathbf{f} \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma^2 \mathbf{I} & \mathbf{K} \\ \mathbf{K} & \mathbf{K} \end{bmatrix}\right). \quad (9.55)$$

By conditioning on the observed outputs \mathbf{y} , the posterior distribution over the latent vector can be inferred

$$\begin{aligned} p(\mathbf{f}|\mathbf{d}) &= \mathcal{N}(\mathbf{f}|\mathbf{K}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{d}, \mathbf{K} - \mathbf{K}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{K}) \\ &= \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}, \boldsymbol{\Sigma}). \end{aligned} \quad (9.56)$$

Assuming this posterior is obtained for the data up till time instant $n - 1$, the predictive distribution of a new output d_n at location \mathbf{x}_n is computed as

$$p(d_n|\mathbf{x}_n, \mathbf{d}_{n-1}) = \mathcal{N}(d_n|\mu_{\text{GP},n}, \sigma_{\text{GP},n}^2) \quad (9.57a)$$

$$\mu_{\text{GP},n} = \mathbf{k}_n^\top (\mathbf{K}_{n-1} + \sigma^2 \mathbf{I})^{-1} \mathbf{d}_{n-1} \quad (9.57b)$$

$$\sigma_{\text{GP},n}^2 = \sigma^2 + k_{nn} - \mathbf{k}_n^\top (\mathbf{K}_{n-1} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_n. \quad (9.57c)$$

591 The mode of the predictive distribution, given by $\mu_{\text{GP},n}$ in Eq. (9.57b), coincides with the
 592 prediction of kernel ridge regression, given by Eq. (9.36), showing that the regularization in
 593 KRR can be interpreted as a noise power σ^2 . Furthermore, the variance of the predictive
 594 distribution, given by $\sigma_{\text{GP},n}^2$ in Eq. (9.57c), coincides with Eq. (9.42), which is used by the
 595 ALD dictionary criterion for KRLS.
 596

597 9.6.2 Online recursive solution for Gaussian process regression

A recursive update of the complete GP (9.57) was proposed in [Csató and Opper \(2002\)](#), as the Sparse Online Gaussian Process (SOGP) algorithm. We will follow the notation of [Van Vaerenbergh et al. \(2012b\)](#), whose solution is equivalent to SOGP but whose choice of variables allows for an easier interpretation. Specifically, the predictive mean and covariance of the GP solution (9.57) can be updated as

$$p(\mathbf{f}_n|\mathbf{X}_n, \mathbf{d}_n) = \mathcal{N}(\mathbf{f}_n|\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n) \quad (9.58a)$$

$$\boldsymbol{\mu}_n = \begin{bmatrix} \boldsymbol{\mu}_{n-1} \\ \hat{d}_n \end{bmatrix} + \frac{e_n}{\hat{\sigma}_{dn}^2} \begin{bmatrix} \mathbf{h}_n \\ \hat{\sigma}_{fn}^2 \end{bmatrix} \quad (9.58b)$$

$$\boldsymbol{\Sigma}_n = \begin{bmatrix} \boldsymbol{\Sigma}_{n-1} & \mathbf{h}_n \\ \mathbf{h}_n^\top & \hat{\sigma}_{fn}^2 \end{bmatrix} - \frac{1}{\hat{\sigma}_{dn}^2} \begin{bmatrix} \mathbf{h}_n \\ \hat{\sigma}_{fn}^2 \end{bmatrix} \begin{bmatrix} \mathbf{h}_n \\ \hat{\sigma}_{fn}^2 \end{bmatrix}^\top, \quad (9.58c)$$

598 where \mathbf{X}_n contains the n input data, $\mathbf{h}_n = \Sigma_{n-1} \mathbf{K}_{n-1}^{-1} \mathbf{k}_n$, and $\hat{\sigma}_{fn}^2$ and $\hat{\sigma}_{dn}^2$ are the predictive
 599 variances of the latent function and the new output, respectively, calculated at the new input.
 600 Details can be found in [Lázaro-Gredilla *et al.* \(2011\)](#); [Van Vaerenbergh *et al.* \(2012b\)](#). In
 601 particular, the update of the predictive mean can be shown to be equivalent to the KRLS
 602 update. The advantage of using a full GP model is that not only does it allow to update the
 603 predictive mean, as does KRLS, but it keeps track of the entire predictive distribution of
 604 the solution. This allows, for instance, to establish confidence intervals when predicting new
 605 outputs.

606 Similar to KRLS, this online GP update assumes a stationary model. Interestingly however,
 607 the Bayesian approach (and in particular its handling of the uncertainty) does allow for a
 608 principled extension that performs tracking, as we briefly discuss in the sequel.

609 9.6.3 Kernel Recursive Least Squares Tracker

610 In [Van Vaerenbergh *et al.* \(2012b\)](#), a KRLS Tracker (KRLS-T) algorithm was presented that
 611 explicitly handles uncertainty about the data, based on the probabilistic GP framework. In
 612 stationary environments it operates identically to the earlier proposed Sparse Online GP
 613 algorithm (SOGP) from [Csató and Opper \(2002\)](#), though it includes a *forgetting mechanism*
 614 that enables it to handle non-stationary scenarios as well.

During each iteration, KRLS-T performs a forgetting operation in which the mean and
 covariance are replaced through

$$\boldsymbol{\mu} \leftarrow \sqrt{\lambda} \boldsymbol{\mu} \quad (9.59a)$$

$$\boldsymbol{\Sigma} \leftarrow \lambda \boldsymbol{\Sigma} + (1 - \lambda) \mathbf{K}. \quad (9.59b)$$

615 The effect of this operation on the predictive distribution is shown in Fig. 9.10. For illustration
 616 purposes, the forgetting factor is chosen unusually low, $\lambda = 0.9$.

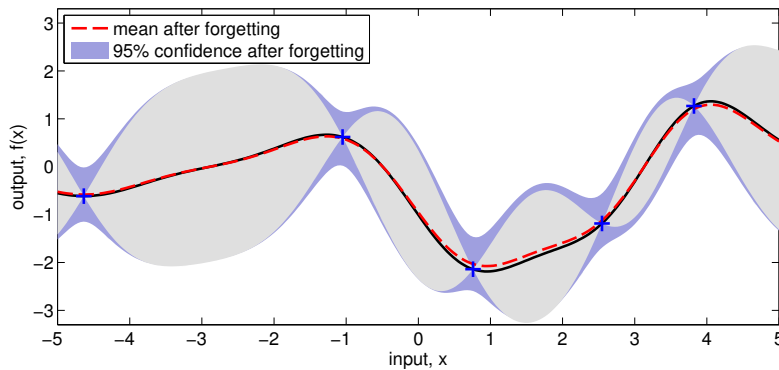


Figure 9.10 Forgetting operation of KRLS-T. The original predictive mean and variance are indicated as the black line and shaded grey area, as in Fig. 9.9. After one forgetting step, the mean becomes the dashed red curve, and the new 95% confidence interval is indicated in blue.

617 This particular form of forgetting corresponds to blending the informative posterior with
 618 a “noise” distribution that uses the same color as the prior. In other words, forgetting occurs
 619 by taking a step back towards the prior knowledge. Since the prior has zero mean, the

mean is simply scaled by the square root of the forgetting factor λ . The covariance, which represents the posterior uncertainty on the data, is pulled towards the covariance of the prior. Interestingly, a regularized version of RLS (known as *extended RLS*) can be obtained by using a linear kernel with the B2P forgetting procedure. Standard RLS can be obtained by using a different forgetting rule, see [Van Vaerenbergh et al. \(2012b\)](#).

The KRLS-T algorithm can be seen as a probabilistic extension of KRLS that obtains confidence intervals and is capable of adapting to time-varying environments. It obtains state-of-the-art performance in several nonlinear adaptive filtering problems, see [Van Vaerenbergh and Santamaría \(2013\)](#) and the results of Fig. 9.7, though it has a more complex formulation than most other kernel adaptive filters and it requires a higher computational complexity. We will explore these aspects through additional examples in Section 9.8.

9.6.4 Probabilistic KLMS

The success of the probabilistic approach for KRLS-like algorithms has lead several researchers to investigate the design of probabilistic KLMS algorithms. The low complexity of KLMS-type algorithms makes them very popular in practical solutions. Nevertheless, this low computational complexity is also a limitation that makes the design of a probabilistic KLMS algorithm a particularly hard research problem.

Some advances have already been made in this direction. Specifically, in [Park et al. \(2014\)](#) a probabilistic KLMS algorithm was proposed, though it only considered the maximum-a-posteriori (MAP) estimate. In [Van Vaerenbergh et al. \(2016a\)](#), it was shown that several KLMS algorithms can be obtained by imposing a simplifying restriction on the full SOGP model, thereby linking KLMS algorithms and online GP approaches directly.

9.7 Further Reading

A myriad of different kernel adaptive filtering algorithms have appeared in the literature. We described the most prominent algorithms, which represent the state of the art. While we only focused on their online learning operation, several other aspects are worth studying. In this section, we briefly introduce the most interesting topics that are the subject of current research.

9.7.1 Selection of Kernel Parameters

A typical problem in kernel methods in general and kernel adaptive filtering in particular is the determination of the optimal kernel and other parameters, such as regularization, forgetting factor, embedding size, etc. These parameters are often referred to as *hyperparameters* in order to distinguish them from the kernel expansion coefficients α_i . A standard approach to determine the optimal hyperparameters is to perform a grid search with cross-validation. Nevertheless, this approach quickly becomes infeasible when more than a few hyperparameters or parameter values are to be considered, due to the combinatorial explosion of possible grid points to evaluate.

A more efficient and principled method is offered by the Gaussian process framework. Specifically, the optimal hyperparameters of GP regression can be found by maximizing the log marginal likelihood, which has an analytic expression given by

$$\log p(\mathbf{d}|\mathbf{X}) = -\frac{1}{2}\mathbf{d}^\top (\mathbf{K} + \sigma^2\mathbf{I})^{-1} \mathbf{d} - \frac{1}{2} \log \|\mathbf{K}\| - \frac{n}{2} \log 2\pi. \quad (9.60)$$

657 It is straightforward to compute this log marginal likelihood and its gradients, and one can
658 choose any of the existing nonlinear optimization methods to perform the maximization. This
659 procedure is commonly referred to as type-II maximum likelihood (ML). Details can be found
660 in [Rasmussen and Williams \(2006\)](#).

661 The optimal hyperparameters obtained by type-II ML correspond to the optimal choices
662 for kernel ridge regression, due to the correspondence between GP regression and KRR, and
663 for several kernel adaptive filters. A case study for KRLS and KRLS-T can be found in
664 [Van Vaerenbergh et al. \(2012a\)](#).

665 Finally, in online scenarios it would be interesting to perform an online estimation of the
666 optimal hyperparameters. This, however, is a difficult open research problem for which only a
667 handful of methods have been proposed, see for instance [Soh and Demiris \(2015\)](#). In practice
668 it is still more appropriate to perform type-II ML offline on a batch of training data, before
669 running the online learning procedure using the found hyperparameters.

670 9.7.2 Multi-Kernel Adaptive Filtering

671 In the last decade, several methods have been proposed to consider multiple kernels instead
672 of a single one [Bach et al. \(2004\)](#); [Sonnenburg et al. \(2006\)](#). The different kernels may
673 correspond to different notions of similarity, or they may address information coming from
674 multiple, heterogeneous data sources.

675 On the other hand, in the field of linear adaptive filtering, it was recently shown that
676 a convex combination of adaptive filters can improve the convergence rate and tracking
677 performance [Arenas-García et al. \(2006\)](#) compared to running a single adaptive filter.

678 Multi-kernel adaptive filtering combines ideas from the above two approaches [Yukawa
679 \(2012\)](#). Its learning procedure activates those kernels whose hyperparameters correspond
680 best to the currently observed data, which could be interpreted as a form of hyperparameter
681 learning. Furthermore, the adaptive nature of these algorithms allow them to track the
682 importance of each kernel in time-varying scenarios, possibly giving them an advantage over
683 single-kernel adaptive filtering.

684 Several multi-kernel adaptive filtering algorithms have been proposed in the recent
685 literature, for instance [Gao et al. \(2014\)](#); [Ishida and Tanaka \(2013\)](#); [Pokharel et al. \(2013\)](#);
686 [Yukawa \(2012\)](#). While they show promising performance gains over single-kernel adaptive
687 filtering algorithms, their computational complexity is much higher. This is an important
688 aspect inherent to the combination of multiple kernel methods, and it is a topic of current
689 research.

690 9.7.3 Recursive Filtering in Kernel Hilbert Spaces

691 The modeling and prediction of time series with kernel adaptive filters is usually addressed
692 by time-embedding the data, thus considering each time lag as a different input dimension.
693 This approach presents some drawbacks: First, the optimal filter order may change over time,
694 which would require an additional tracking mechanism; Second, if the optimal filter order
695 is high, as for instance in audio applications [Van Vaerenbergh et al. \(2016b\)](#), the method be
696 affected by the curse of dimensionality. For some problems, the concept of an optimal filter
697 order may not even make sense.

698 An alternative approach to modeling and predicting time series is to construct *recursive*
699 kernel machines, which implement recursive models explicitly in the rkHs. A preliminary
700 work in this direction considered the design of a *recursive kernel* in the context of infinite
701 recurrent neural networks [Hermans and Schrauwen \(2012\)](#). More recently, recursive versions

of the autoregressive, moving-average and gamma filters in rkHs were proposed [Tuia et al. \(2014\)](#). By exploiting properties of functional analysis and recursive computation, this approach avoids the reduced-rank approximations that are required in standard kernel adaptive filters. Finally, a kernel version of the autoregressive-moving-average filter was presented in [Li and Príncipe \(2016\)](#).

9.8 Tutorial Examples

This section presents experiments in which kernel adaptive filters are applied to time series prediction and nonlinear system identification. These experiments are implemented using code based on the Kernel Adaptive Filtering Toolbox [Van Vaerenbergh and Santamaría \(2013\)](#), which is available at <https://github.com/steven2358/kafbox/>.

9.8.1 Kernel Adaptive Filtering Toolbox

The Kernel Adaptive Filtering Toolbox (KAFBOX) is a Matlab benchmarking toolbox to evaluate and compare kernel adaptive filtering algorithms. It includes a large list of algorithms that have appeared in the literature, and additional tools for hyperparameter estimation and algorithm profiling, among others.

The kernel adaptive filtering algorithms in KAFBOX are implemented as objects using the `classdef` syntax. Since all KAF algorithms are online methods, each of them includes two basic operations: 1) Obtaining the filter output, given a new input \mathbf{x}_* ; and 2) Training on a new data pair (\mathbf{x}_n, d_n) . These operations are implemented as the methods `evaluate` and `train`, respectively.

As an example, we list the code for the KLMS algorithm in Listing 9.9. The object definition contains two sets of properties, one for the hyperparameters and one for the variables it will learn. The first method is the object's constructor method, which copies the specified hyperparameter settings. The second method is the `evaluate` function, which performs the operation $y_* = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x}_*)$. It includes an `if` clause to check if the algorithm has at least performed one training step yet. If not, zeroes are returned as predictions. Finally, the `train` method implements a single training step of the online learning algorithm. This method typically handles algorithm initialization as well, such that functions that operate on a KAF object do not have to worry about initializing. The training step itself is summarized in very few lines of Matlab code, for many algorithms.

```

732 % Kernel Least-Mean-Square algorithm
733 %
734 % W. Liu, P.P. Pokharel, and J.C. Principe, "The Kernel Least-Mean-Square
735 % Algorithm," IEEE Transactions on Signal Processing, vol. 56, no. 2, pp.
736 % 543-554, Feb. 2008, http://dx.doi.org/10.1109/TSP.2007.907881
737 %
738 % Remark: implementation includes a maximum dictionary size M
739 %
740 % This file is part of the Kernel Adaptive Filtering Toolbox for Matlab.
741 % https://github.com/steven2358/kafbox/
742
743 classdef klms < handle
744
745     properties (GetAccess = 'public', SetAccess = 'private')
746         eta = .5; % learning rate
747         M = 10000; % maximum dictionary size
748         kerneltype = 'gauss'; % kernel type

```

```

749     kernelpar = 1; % kernel parameter
750 end
751
752 properties (GetAccess = 'public', SetAccess = 'private')
753     dict = []; % dictionary
754     alpha = []; % expansion coefficients
755 end
756
757 methods
758
759     function kaf = klms(parameters) % constructor
760         if (nargin > 0) % copy valid parameters
761             for fn = fieldnames(parameters)',
762                 if ismember(fn,fieldnames(kaf)),
763                     kaf.(fn{1}) = parameters.(fn{1});
764                 end
765             end
766         end
767     end
768
769     function y_est = evaluate(kaf,x) % evaluate the algorithm
770         if size(kaf.dict,1)>0
771             k = kernel(kaf.dict,x,kaf.kerntype,kaf.kernelpar);
772             y_est = k'*kaf.alpha;
773         else
774             y_est = zeros(size(x,1),1);
775         end
776     end
777
778     function train(kaf,x,y) % train the algorithm
779         if (size(kaf.dict,1)<kaf.M), % avoid infinite growth
780             y_est = kaf.evaluate(x);
781             err = y - y_est;
782             kaf.alpha = [kaf.alpha; kaf.eta*err]; % grow
783             kaf.dict = [kaf.dict; x]; % grow
784         end
785     end
786
787 end
788 end

```

Listing 9.9 Matlab code for the KLMS algorithm object class, from KAFBOX.

789 For the experiments below, we will use v2.0 of KAFBOX, which can be downloaded from
790 the “releases” page <https://github.com/steven2358/kafbox/releases/>.

791 9.8.2 Prediction of a Respiratory Motion Time Series

792 In the first experiment we apply KAF algorithms to predict a bio-medical time series, more
793 specifically a respiratory motion trace. These data come from robotic radiosurgery, in which
794 a photon beam source is used to ablate tumors. The beam is operated by a robot arm that aims
795 to move the beam source to compensate for the motion of internal organs. Traditionally,
796 this is achieved by recording the motion of markers applied to the body surface and by
797 using this motion to draw conclusions about the tumor position. Although this method
798 significantly increases the targeting accuracy, the system delay arising from data processing
799 and positioning of the beam results in a systematic error. This error can be decreased by
800 predicting the motion of the body surface Ernst (2012).

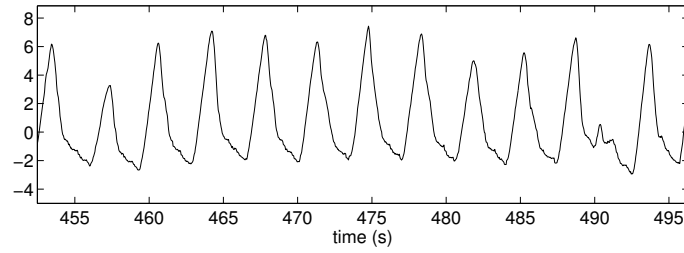


Figure 9.11 A snapshot of the respiratory motion trace.

801 The data was recorded at Georgetown University Hospital using CyberKnife® equipment,
 802 and it represents the recorded position of one of the markers attached to the body surface³. A
 803 snapshot of this motion trace is shown in Fig. 9.11. The delay to compensate totals 115 ms,
 804 which, at a sampling frequency of 26 Hz, corresponds to 3 samples. The task thus consists
 805 in three-step ahead prediction. We use a time-embedding of 8 samples. Since the breathing
 806 pattern may change over time, we employed only tracking algorithms. Their parameters are
 807 listed in Table 9.1. The MSE results of the four algorithms are displayed in the last column
 808 of this table. A comparison of the original series and the predictions of one of the algorithms
 809 is shown in Fig. 9.12. The code to reproduce these results can be found in Listing 9.10.

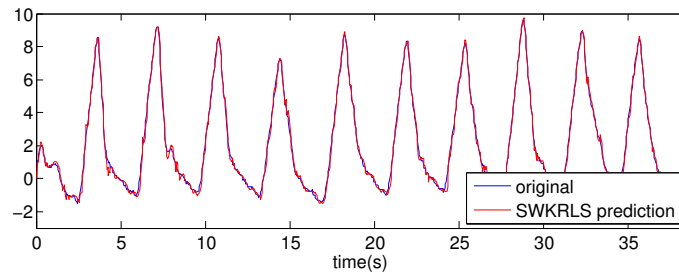


Figure 9.12 The respiratory motion trace and the 3-step ahead prediction of a KAF algorithm.

Table 9.1 Parameters used for predicting the respiratory motion trace, MSE result for 3-step ahead prediction, and measured training time.

Algorithm	Parameters	MSE performance	Training time
NORMA	$\lambda = 10^{-4}, \tau = 30$	-5.78 dB	0.16 s
QKLMS	$\eta = 0.99, \epsilon_{\text{U}} = 1$	-8.14 dB	0.18 s
SWKRLS	$c = 10^{-4}, m = 50$	-13.35 dB	0.29 s
KRLS-T	$\sigma_n^2 = 10^{-4}, m = 50, \lambda = 0.999$	-18.16 dB	0.54 s

³Data available at <http://signals.rob.uni-luebeck.de/>

```

810 % 3-step ahead prediction on the respiratory motion time series.
811 % Requires KAFBOX toolbox.
812 close all; clear
813
814 %% PARAMETERS
815 h = 3; % prediction horizon
816 L = 8; % embedding
817 n = 1000; % number of data
818
819 sigma = 7; % kernel parameter
820
821 % Uncomment one of the following algorithms. All use a Gaussian kernel.
822 % kaf = norma(struct('lambda',1E-4,'tau',30,'kernelpar',sigma,'eta',0.99));
823 % kaf = qklms(struct('epsu',1,'kernelpar',sigma,'eta',0.99));
824 kaf = swkrls(struct('M',50,'kernelpar',sigma,'c',1E-4));
825 % kaf = krlst(struct('M',50,'lambda',0.999,'sn2',1E-4,'kernelpar',sigma));
826
827 %% PREPARE DATA
828 data = load('respiratorymotion3.dat');
829 X = zeros(n,L);
830 for i = 1:L,
831     X(i:n,i) = data(1:n-i+1,1); % time embedding
832 end
833 y = data((1:n)+h);
834
835 %% RUN ALGORITHM
836 MSE = zeros(n,1);
837 y_est_all = zeros(n,1);
838
839 title_ = upper(class(kaf)); % store algorithm name
840 fprintf('Training %s',title_)
841 for i=1:n,
842     if ~mod(i,floor(n/10)), fprintf('.'); end % progress indicator
843
844     xi = X(i,:);
845     yi = y(i);
846
847     y_est = kaf.evaluate(xi); % evaluate on test data
848     MSE(i) = (yi-y_est)^2; % test error
849     y_est_all(i) = y_est;
850
851     kaf.train(xi,yi); % train with one input-output pair
852 end
853 fprintf('\n');
854
855 %% OUTPUT
856 fprintf('Mean MSE: %.2fdB\n',10*log10(mean(MSE)));
857
858 figure; hold all;
859 t = (1:n)/26; % sample rate is 26 Hz
860 plot(t,y)
861 plot(t,y_est_all,'r')
862 legend({'original',sprintf('%s prediction',title_)},'Location','SE');

```

Listing 9.10 Matlab code for running KAF algorithms on the respiratory motion prediction problem.

863 9.8.3 Online Regression on the KIN40K Data Set

864 In the second experiment we train the online algorithms to perform regression of the KIN40K
 865 data set Ghahramani (1996)⁴, which is a standard regression problem in the machine learning
 866 literature. The KIN-40K data set is obtained from the forward kinematics of an 8-link all-
 867 revolute robot arm, similar to the one depicted in Fig. 9.13. It contains 40000 examples, each
 868 consisting of an 8-dimensional input vector and a scalar output. KIN40K was generated with
 869 maximum nonlinearity and little noise, representing a very difficult regression test.

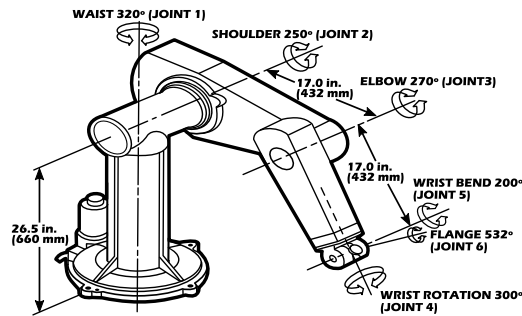


Figure 9.13 Sketch of a 5-link all-revolute robot arm. The data used in the KIN40K experiment were generated by simulating an 8-link extension of this arm.

870 In this experiment, we first determine the optimal hyperparameters for the kernel adaptive
 871 filters by running the tool `kafbox_parameter_estimation`, which is based on the
 872 GPML toolbox from Rasmussen and Williams (2006). We use 1000 randomly selected data
 873 points for the hyperparameter optimization. In the literature, an anisotropic kernel function,
 874 which has a different kernel width per dimension, is commonly used on these data. For
 875 simplicity, though, we employ an isotropic Gaussian kernel. The hyperparameters found
 876 by the optimization procedure are listed in Table 9.2. The forgetting factor is only used by
 877 KRLS-T, though its optimal value is determined to be 1, which indicates that no forgetting
 878 takes place in practice.

Table 9.2 Optimal hyperparameters found for the KIN40K regression problem.

Parameter	Optimal value
kernel width σ	1.66
regularization	$2.47 \cdot 10^{-6}$
forgetting factor λ	1

879 From the remaining data we randomly select 5000 data points for training and 5000 data for
 880 testing the regression. Apart from the hyperparameters that are determined automatically in
 881 this experiment, the kernel adaptive filters have some parameters relating to memory size and
 882 learning rate. The values chosen for these parameters are listed in Table 9.3. The values for

⁴Data available at <http://www.cs.toronto.edu/~delve/data/datasets.html>

883 QKLMS are chosen such that it obtains optimal performance after training with a dictionary
 884 size that is one order of magnitude larger than that of the KRLS algorithms. The precision
 885 parameter ν of KRLS is tuned to yield a dictionary size of around $m = 500$ at the end of the
 886 experiment, which is the budget of SWKRLS and KRLS-T.

887 The learning curves for this experiment are shown in Fig. 9.14. The code for reproducing
 888 this experiment is displayed in Listing 9.11.

Table 9.3 Additional parameters used in the KIN40K regression experiment, final dictionary size, and measured training time.

Algorithm	Parameters	Final dictionary size	Training time
QKLMS	$\eta = 0.99, \epsilon_{\mathbb{U}} = 1.2$	2750	18.51 s
KRLS	$\nu = 0.32$	510	12.09 s
FBKRLS	$m = 500$	500	33.49 s
KRLS-T	$m = 500$	500	86.41 s

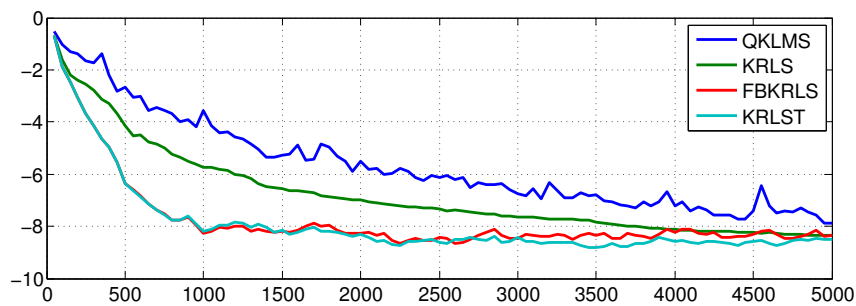


Figure 9.14 Learning curves of different algorithms on the KIN40K data.

```

889 % This demo estimates the optimal hyperparameters on the KIN40K data and
890 % performs online regression learning. The estimated parameters are:
891 % forgetting factor lambda, regularization c and Gaussian kernel width.
892
893 close all; clear
894 rng(1); % fix random seed, for reproducibility
895
896 %% PARAMETERS
897
898 n_hyper = 1000; % number of data to use for hyperparameter estimation
899 n_train = 5000; % number of data for online training
900 n_test = 5000; % number of data for testing
901 test_every = 50; % run test every time this number of iterations passes
902
903 % selected algorithms are defined below
904
905 %% PROGRAM
906 tic
907

```



```

908 fprintf('Loading KIN40K data...\n')
909 load('kin40k'); % data and hyperparameters
910
911 %% PREPARE DATA
912 indp = randperm(length(y)); % random permutation
913 ind_hyper = indp(1:n_hyper);
914 ind_train = indp(n_hyper+1:n_hyper+n_train);
915 ind_test = indp(n_hyper+n_train+1:n_hyper+n_train+n_test);
916 X_hyper = X(ind_hyper,:); % data for hyperparameter estimation
917 y_hyper = y(ind_hyper);
918 X_train = X(ind_train,:); % training data
919 y_train = y(ind_train);
920 X_test = X(ind_test,:); % test data
921 y_test = y(ind_test);
922
923 fprintf('Estimating KRLS-T parameters...\n\n')
924 [sigma,reg,ff] = kafbox_parameter_estimation(X_hyper,y_hyper);
925
926 % select algorithms
927 i=1;
928 algos{i} = qklms(struct('eta',0.5,'epsu',1.2,'kernelpar',sigma)); i=i+1;
929 algos{i} = krls(struct('nu',.32,'kernelpar',sigma)); i=i+1;
930 algos{i} = fbkrls(struct('M',500,'lambda',reg,'kernelpar',sigma)); i=i+1;
931 algos{i} = ...
932     krlst(struct('lambda',ff,'M',500,'sn2',reg,'kernelpar',sigma)); i=i+1;
933
934 n_algos = length(algos);
935 MSE = nan*zeros(n_train,1);
936 titles = cell(n_algos,1);
937 final_dict_size = zeros(n_algos,1);
938
939 for j=1:n_algos
940     kaf = algos{j};
941
942     titles{j} = upper(class(kaf));
943     fprintf(sprintf('Running %s with estimated parameters...\n',titles{j}))
944
945     for i=1:n_train,
946         if ~mod(i,floor(n_train/10)), fprintf('.'); end % progress ...
947             indicator, 10 dots
948
949         kaf.train(X_train(i,:),y_train(i)); % train with one input-output ...
950             pair
951
952         if mod(i,test_every) == 0 % run test only every
953             y_est = kaf.evaluate(X_test); % predict on test set
954             MSE(i,j) = mean((y_test-y_est).^2);
955         end
956     end
957     final_dict_size(j) = size(kaf.dict,1);
958     fprintf('\n');
959 end
960
961 toc
962 %% OUTPUT
963
964 figure;
965 xs=find(~isnan(MSE(:,1)));
966 plot(xs,10*log10(MSE(xs,:)))
967 legend(titles)
968
969 fprintf('\n');

```

```

970 fprintf('      Estimated\n');
971 fprintf('sigma:  %.4f\n', sigma)
972 fprintf('c:      %e\n', reg)
973 fprintf('lambda: %.4f\n\n', ff)
974
975 final_dict_size

```

976 **Listing 9.11** Matlab code for determining the optimal hyperparameters and running online regression
 977 on the KIN40K data.

978 9.8.4 The Mackey-Glass Time Series

979 The Mackey-Glass time-series prediction is a benchmarking problem in nonlinear time-series
 980 modelling. We discussed this time series and the prediction results for KLMS and KRLS in
 981 Sections 9.3.3 and 9.4.3, respectively.

982 The learning curves, shown in Fig. 9.5, indicate that KLMS converges very slowly, and that
 983 KRLS can obtain a much lower MSE in less iterations. On the other hand, KRLS requires
 984 an order of magnitude more computation and memory. These results are in line with the
 985 intuitions from linear adaptive filtering, in which LMS and RLS represent two different
 986 choices in the compromise between complexity and convergence rate.

987 Nevertheless, there is a fundamental difference between the complexity analyses of linear
 988 and kernel adaptive filtering algorithms. While in linear adaptive filters the complexity
 989 depends on the data dimension, in KAF algorithms it depends on the dictionary size. And,
 990 importantly, the latter is a parameter that can be controlled.

991 A KRLS-type algorithm with a large dictionary can converge faster than a KLMS-
 992 type algorithm with a similarly sized dictionary, at the expense of a higher computational
 993 complexity. But it would be instructive to ask how a KRLS-type algorithm with a small
 994 dictionary compares to a KLMS-type algorithm with a large dictionary. Can the KRLS
 995 algorithm obtain similar complexity as KLMS, while maintaining its better convergence
 996 rate? This question is answered in the diagrams of Fig. 9.15. We have included two KAF
 997 algorithms, QKLMS and KRLS-T, that allow easy control over their dictionary size.

Table 9.4 Parameters used in the Mackey-Glass time series prediction. A Gaussian kernel with $\sigma_k = 1$ was used.

Algorithm	Fixed parameters	Varying parameter
QKLMS	$\eta = 0.5$	$\epsilon_U \in \{10^{-4}, 10^{-3}, 10^{-2}, 0.1, 0.2, 0.3, 0.5, 0.7, 1\}$
KRLS-T	$\lambda = 1, \sigma_n^2 = 10^{-6}$	$m \in \{3, 5, 7, 10, 20, 30, 50, 150\}$

998 Fig. 9.15 represents the results obtained by the KAF profiler tool included in KAFBOX.
 999 The Matlab code to reproduce this figure is displayed in Listing 9.12. The profiler tool runs
 1000 each algorithm several times with different configurations, whose parameters are shown in
 1001 Table 9.4, producing one point in the plot per algorithm configuration. It calculates several
 1002 variables, such as the number of floating point operations (FLOPS), the used memory (in
 1003 bytes), and execution time.

1004 By plotting the MSE vs. the FLOPS or memory, we get an idea of the resources required to
 1005 obtain a desired MSE result. If, for instance, we are working in a scenario with a restriction
 1006 on computational complexity, we should select the algorithm that performs best under this
 1007 restriction by determining which performance curve is most to the left for the amount of

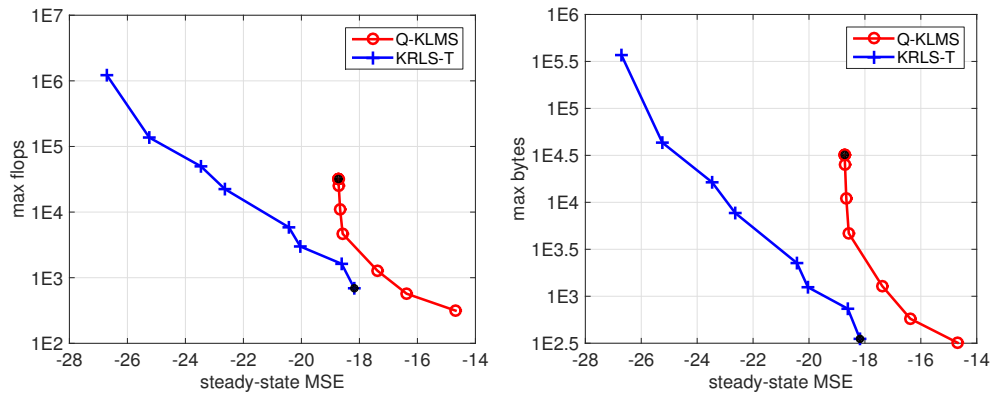


Figure 9.15 MSE vs. complexity trade-off comparisons for prediction of the Mackey-Glass time-series. Left: Maximum number of FLOPS per iteration as a function of the steady-state MSE. Right: Maximum number of bytes per iteration as a function of the steady-state MSE. Each marker represents a single run of one of the algorithms with a single set of parameters. The start of each parameter sweep is indicated by a black dot.

1008 FLOPS available. In the same manner, by fixing a maximum on MSE we obtain the FLOPS
 1009 and memory required by each algorithm. In the left plot of Fig. 9.15 we observe that if
 1010 the available computational complexity is very limited, it may be more interesting to use
 1011 QKLMS. In other cases, KRLS-T is preferred as it obtains better MSE for the same amount
 1012 of FLOPS. In terms of memory used, it appears that it is always advantageous to use KRLS-T,
 1013 as can be seen in the right plot.

```

1014 % Experiment: kernel adaptive filter algorithm profiler.
1015 % Compares the cost vs prediction error tradeoffs and convergence speeds
1016 % for several algorithms on the MG30 data set.
1017 % Requires KAFBOX toolbox.
1018
1019 clear
1020 close all
1021
1022 %% PARAMETERS
1023
1024 % data and algorithm setup
1025 data.name = 'mg30';
1026 data.n_train = 500; % number of data points
1027 data.n_test = 100; % number of data points
1028 data.embedding = 7; % time embedding
1029 data.offset = 50; % apply offset per simulation
1030
1031 sim_opts.numsim = 5; % 10 seconds per simulation on a 2016 PC
1032 sim_opts.error_measure = 'MSE';
1033
1034 i=0; % initialize setups
1035
1036 %% QKLMS
1037 i=i+1;
1038 algorithms{i}.name = 'QKLMS';
1039 algorithms{i}.class = 'qklms';
1040 algorithms{i}.figstyle = struct('color',[1 0 0], 'marker','o');
```

```

1041 algorithms{i}.options = ...
1042     struct('eta',0.5,'sweep_par','epsu','sweep_val',[1E-4 1E-3 1E-2 .1 ...
1043         .2 .3 .5 .7 1],...
1044         'kerneltype','gauss','kernelpar',1);
1045
1046 %% KRLS-T
1047 i=i+1;
1048 algorithms{i}.name = 'KRLS-T';
1049 algorithms{i}.class = 'krlst';
1050 algorithms{i}.figstyle = struct('color',[0 0 1],'marker','+');
1051 algorithms{i}.options = struct('sn2',1E-6,'lambda',1,'sweep_par','M',...
1052     'sweep_val',[3 5 7 10 20 30 50 150],...
1053     'kerneltype','gauss','kernelpar',1);
1054
1055 %% PROGRAM
1056
1057 fprintf('Running profiler for %d algorithms on %s data.\n',i,data.name);
1058 output_dir = fullfile(mfilename('fullpath'),'..','results');
1059
1060 t1 = tic;
1061 [data,algorithms,results] = ...
1062     kafbox_profiler(data,sim_opts,algorithms,output_dir);
1063 t2 = toc(t1);
1064
1065 fprintf('Elapsed time: %d seconds\n',ceil(t2));
1066
1067 %% OUTPUT
1068
1069 mse_curves = kafbox_profiler_msecurves(results);
1070
1071 kafbox_profiler_plotresults(algorithms,mse_curves,results,{'ssmse','flops'});
1072
1073 kafbox_profiler_plotresults(algorithms,mse_curves,results,{'ssmse','bytes'});
1074
1075 resinds = [1,2;2,8]; % result indices
1076 kafbox_profiler_plotconvergence(algorithms,mse_curves,resinds);

```

Listing 9.12 Matlab code for profiling QKLS and KRLS-T on the Mackey-Glass time series.

1077 9.9 Questions and Problems

1078 **Exercise 9.9.1** *When is it be more useful to employ a KLMS-like algorithm and when a*
1079 *KRLS-like algorithm?*

1080 **Exercise 9.9.2** *Demonstrate that the computational complexity of KLMS is $\mathcal{O}(m)$, where m*
1081 *is the number of data in its dictionary.*

1082 **Exercise 9.9.3** *Demonstrate that the computational complexity of KRLS is $\mathcal{O}(m^2)$, where m*
1083 *is the number of data in its dictionary.*

1084 **Exercise 9.9.4** *List the advantages and disadvantages of using a sliding-window approach*
1085 *for determining a filter's dictionary, as used for instance by SW-KRLS and NORMA.*

1086 **Exercise 9.9.5** *In the tracking experiment of Section 9.4.5, the slope of the learning curve*
1087 *for some algorithms is less steep just after the switch than in the beginning of the experiment.*
1088 *Identify for which algorithms this happens, in Fig. 9.7, and explain for each of these*
1089 *algorithms why this is the case.*

References

- 1090
- 1091 Arenas-García, J., Figueiras-Vidal, A. R., and Sayed, A. H. (2006). Mean-square performance of a convex
1092 combination of two adaptive filters. *IEEE transactions on signal processing*, **54**(3), 1078–1090.
- 1093 Arulampalam, M. S., Maskell, S., Gordon, N., and Clapp, T. (2002). A tutorial on particle filters for online
1094 nonlinear/non-Gaussian Bayesian tracking. *IEEE Trans. Signal Process.*, **50**(2), 174–188.
- 1095 Bach, F. R., Lanckriet, G. R., and Jordan, M. I. (2004). Multiple kernel learning, conic duality, and the SMO
1096 algorithm. In *Proceedings of the twenty-first international conference on Machine learning*, page 6. *ACM*.
- 1097 Chen, B., Zhao, S., Zhu, P., and Príncipe, J. C. (2012). Quantized kernel least mean square algorithm. *IEEE*
1098 *Transactions on Neural Networks and Learning Systems*, **23**(1), 22–32.
- 1099 Csató, L. and Opper, M. (2002). Sparse online Gaussian processes. *Neural Computation*, **14**(3), 641–668.
- 1100 De Kruif, B. J. and De Vries, T. J. A. (2003). Pruning error minimization in least squares support vector machines.
1101 *IEEE Transactions on Neural Networks*, **14**(3), 696–702.
- 1102 Del Moral, P. (1996). Non-linear filtering: interacting particle resolution. *Markov processes and related fields*, **2**(4),
1103 555–581.
- 1104 Dorffner, G. (1996). *Neural networks for time series processing*. *Neural Network World*, **6**, 447–468.
- 1105 Engel, Y., Mannor, S., and Meir, R. (2004). The kernel recursive least squares algorithm. *IEEE Transactions on*
1106 *Signal Processing*, **52**(8), 2275–2285.
- 1107 Ernst, F. (2012). Compensating for quasi-periodic motion in robotic radiosurgery. *Springer*.
- 1108 Frieß, T.-T. and Harrison, R. F. (1999). A kernel-based adaline. In *Proceedings of the 7th European Symposium on*
1109 *Artificial Neural Networks (ESANN 1999)*, pages 245–250, *Bruges, Belgium*.
- 1110 Gao, W., Richard, C., Bermudez, J.-C. M., and Huang, J. (2014). Convex combinations of kernel adaptive filters. In
1111 2014 IEEE International Workshop on Machine Learning for Signal Processing (MLSP), pages 1–5. *IEEE*.
- 1112 Ghahramani, Z. (1996). The kin datasets. URL: [http://www.cs.utoronto.ca/~delve/data/kin/
1113 kin.ps.gz](http://www.cs.utoronto.ca/~delve/data/kin/kin.ps.gz).
- 1114 Golub, G. H. and Van Loan, C. F. (2012). *Matrix computations, volume 3*. *JHU Press*.
- 1115 Hassibi, B., Stork, D. G., and Wolff, G. J. (1993). Optimal brain surgeon and general network pruning. In *Neural*
1116 *Networks, 1993.*, *IEEE International Conference on*, pages 293–299. *IEEE*.
- 1117 Hayes, M. H. M. H. (1996). *Statistical digital signal processing and modeling / Monson H. Hayes*. *New York : John*
1118 *Wiley & Sons*.
- 1119 Haykin, S. (1999). *Neural Networks – A Comprehensive Foundation*. *Prentice Hall, 2nd edition*.
- 1120 Haykin, S. (2001). *Adaptive Filter Theory (4th Edition)*. *Prentice Hall*.
- 1121 Hermans, M. and Schrauwen, B. (2012). Recurrent kernel machines: Computing with infinite echo state networks.
1122 *Neural Computation*, **24**(1), 104–133.
- 1123 Hoegaerts, L., Suykens, J., Vandewalle, J., and De Moor, B. (2004). A comparison of pruning algorithms for sparse
1124 least squares support vector machines. *Lecture Notes in Computer Science*, pages 1247–1253.
- 1125 Ishida, T. and Tanaka, T. (2013). Multikernel adaptive filters with multiple dictionaries and regularization. In *Signal*
1126 *and Information Processing Association Annual Summit and Conference (APSIPA), 2013 Asia-Pacific*, pages
1127 1–6. *IEEE*.
- 1128 Julier, S. J. and Uhlmann, J. K. (1996). A general method for approximating nonlinear transformations of probability
1129 distributions. *Technical report, Technical report, Robotics Research Group, Department of Engineering Science,*
1130 *University of Oxford*.
- 1131 Kalman, R. E. (1960). A new approach to linear filtering and prediction problems. *Journal of basic Engineering*,
1132 **82**(1), 35–45.
- 1133 Kivinen, J., Smola, A. J., and Williamson, R. C. (2004). Online learning with kernels. *IEEE Transactions on Signal*
1134 *Processing*, **52**(8), 2165–2176.
- 1135 Lázaro-Gredilla, M., Van Vaerenbergh, S., and Santamaría, I. (2011). A Bayesian approach to tracking with kernel
1136 recursive least-squares. In 2011 IEEE International Workshop on Machine Learning for Signal Processing
1137 (MLSP), pages 1–6.
- 1138 LeCun, Y., Denker, J. S., Solla, S. A., Howard, R. E., and Jackel, L. D. (1989). Optimal brain damage. In *NIPs,*
1139 *volume 2*, pages 598–605.
- 1140 Lewis, F. L., Xie, L., and Popa, D. (2007). *Optimal and robust estimation: with an introduction to stochastic control*
1141 *theory, volume 29*. *CRC press*.
- 1142 Li, K. and Principe, J. C. (2016). The kernel adaptive autoregressive-moving-average algorithm. *IEEE Transactions*
1143 *on Neural Networks and Learning Systems*, **27**(2), 334–346.
- 1144 Liu, W., Pokharel, P. P., and Principe, J. C. (2008). The kernel least-mean-square algorithm. *IEEE Transactions on*
1145 *Signal Processing*, **56**(2), 543–554.
- 1146 Liu, W., Principe, J. C., and Haykin, S. (2010). *Kernel Adaptive Filtering: A Comprehensive Introduction*. *Wiley*.
- 1147 Ljung, L. (1999). *System identification: theory for the user*. *Prentice-Hall, Upper Saddle River, NJ, USA*.
- 1148 Narendra, K. S. and Parthasarathy, K. (1990). Identification and control of dynamical systems using neural networks.
1149 *IEEE Trans. Neur. Net.*, **1**(1), 4–27.
- 1150 Park, I. M., Seth, S., and Van Vaerenbergh, S. (2014). Probabilistic kernel least mean squares algorithms. In 2014
1151 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 8272–8276. *IEEE*.
- 1152 Plackett, R. L. (1950). Some theorems in least squares. *Biometrika*, **37**, 149–157.
- 1153 Platt, J. (1991). A resource-allocating network for function interpolation. *Neural computation*, **3**(2), 213–225.

- 1154 Pokharel, P. P., Liu, W., and Príncipe, J. C. (2009). Kernel least mean square algorithm with constrained growth.
1155 Signal Processing, **89**(3), 257–265.
- 1156 Pokharel, R., Seth, S., and Príncipe, J. C. (2013). Mixture kernel least mean square. In The 2013 International Joint
1157 Conference on Neural Networks (IJCNN), pages 1–7. IEEE.
- 1158 Rabiner, L. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. Proc.
1159 IEEE, **77**(2), 257–286.
- 1160 Rasmussen, C. E. and Williams, C. K. I. (2006). Gaussian Processes for Machine Learning. MIT Press.
- 1161 Richard, C., Bermúdez, J. C. M., and Honeine, P. (2009). Online prediction of time series data with kernels. IEEE
1162 Transactions on Signal Processing, **57**(3), 1058–1067.
- 1163 Rzepka, D. (2012). Fixed-budget kernel least mean squares. In Proceedings of 2012 IEEE 17th International
1164 Conference on Emerging Technologies & Factory Automation (ETFA 2012), pages 1–4. IEEE.
- 1165 Saunders, C., Gammernan, A., and Vovk, V. (1998). Ridge regression learning algorithm in dual variables. In
1166 Proceedings of the 15th International Conference on Machine Learning (ICML), pages 515–521, Madison, WI,
1167 USA.
- 1168 Sayed, A. H. (2003). Fundamentals of Adaptive Filtering. Wiley-IEEE Press.
- 1169 Schölkopf, B., Herbrich, R., and Smola, A. J. (2001). A generalized representer theorem. In Computational learning
1170 theory, pages 416–426. Springer.
- 1171 Soh, H. and Demiris, Y. (2015). Spatio-temporal learning with the online finite and infinite echo-state gaussian
1172 processes. IEEE transactions on neural networks and learning systems, **26**(3), 522–536.
- 1173 Sonnenburg, S., Rätsch, G., Schäfer, C., and Schölkopf, B. (2006). Large scale multiple kernel learning. Journal of
1174 Machine Learning Research, **7**(Jul), 1531–1565.
- 1175 Takens, F. (1981). Detecting strange attractors in turbulence. Dynamical Systems and Turbulence, **898**, 366–381.
- 1176 Tuia, D., Muñoz-Marí, J., Rojo-Álvarez, J. L., Martínez-Ramón, M., and Camps-Valls, G. (2014). Explicit recursive
1177 and adaptive filtering in reproducing kernel hilbert spaces. IEEE Transactions on Neural Networks and Learning
1178 Systems, **25**(7), 1413–1419.
- 1179 Van Vaerenbergh, S. and Santamaría, I. (2013). A comparative study of kernel adaptive filtering algorithms. In 2013
1180 IEEE Digital Signal Processing (DSP) Workshop and IEEE Signal Processing Education (SPE), Napa, CA, USA.
- 1181 Van Vaerenbergh, S., Vía, J., and Santamaría, I. (2006). A sliding-window kernel RLS algorithm and its application
1182 to nonlinear channel identification. In 2006 IEEE International Conference on Acoustics, Speech, and Signal
1183 Processing (ICASSP), volume 5, pages 789–792, Toulouse, France.
- 1184 Van Vaerenbergh, S., Santamaría, I., Liu, W., and Príncipe, J. C. (2010). Fixed-budget kernel recursive least-squares.
1185 In 2010 IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP), Dallas, TX, USA.
- 1186 Van Vaerenbergh, S., Santamaría, I., and Lázaro-Gredilla, M. (2012a). Estimation of the forgetting factor in kernel
1187 recursive least squares. In 2012 IEEE International Workshop on Machine Learning for Signal Processing
1188 (MLSP).
- 1189 Van Vaerenbergh, S., Lázaro-Gredilla, M., and Santamaría, I. (2012b). Kernel recursive least-squares tracker for
1190 time-varying regression. IEEE Transactions on Neural Networks and Learning Systems, **23**(8), 1313–1326.
- 1191 Van Vaerenbergh, S., Fernandez-Bes, J., and Elvira, V. (2016a). On the relationship between online Gaussian process
1192 regression and kernel least mean squares algorithms. In 2016 IEEE International Workshop on Machine Learning
1193 for Signal Processing (MLSP), Salerno, Italy. IEEE.
- 1194 Van Vaerenbergh, S., Communiello, D., and Azpicueta-Ruiz, L. A. (2016b). A split kernel adaptive filtering
1195 architecture for nonlinear acoustic echo cancellation. In 24th European Signal Processing Conference (EUSIPCO
1196 2016), Budapest, Hungary.
- 1197 Vapnik, V. N. (1995). The Nature of Statistical Learning Theory. Springer-Verlag New York, Inc., New York, NY,
1198 USA.
- 1199 Widrow, B., McCool, J., and Ball, M. (1975). The complex LMS algorithm. Proceedings of the IEEE, **63**(4), 719–720.
- 1200 Yukawa, M. (2012). Multikernel adaptive filtering. IEEE Transactions on Signal Processing, **60**(9), 4672–4682.
- 1201 Zhao, S., Chen, B., Zhu, P., and Príncipe, J. C. (2013). Fixed budget quantized kernel least-mean-square algorithm.
1202 Signal Processing, **93**(9), 2759–2770.
- 1203 Zhao, S., Chen, B., Cao, Z., Zhu, P., and Principe, J. C. (2016). Self-organizing kernel adaptive filtering. EURASIP
1204 Journal on Advances in Signal Processing, **2016**(1), 106.

INDEX

- 1205 *A-posteriori error based pruning*
1206 *criterion, 23*
1207 *Adaptive filtering, 3, 4*
1208 *Approximate Linear Dependency*
1209 *criterion, 15, 22*
- 1210 *Bayesian recursive update, 25*
1211 *Budget, 18*
1212 *Bytes, 36*
- 1213 *Coherence criterion, 12, 22*
- 1214 *Dictionary growing, 21*
1215 *Dictionary pruning, 22*
- 1216 *FB-KRLS algorithm, 18*
1217 *Fixed-budget online learning, 18*
1218 *FLOPS, 36*
1219 *Forgetting, 26*
1220 *forgetting, 26*
- 1221 *Gaussian process regression, 24*
- 1222 *hyperparameter optimization, 27*
1223 *Hyperparameters, 24*
- 1224 *Kernel Adaptive Filtering toolbox, 29*
1225 *Kernel ridge regression, 14*
1226 *KIN40K data set, 32*
1227 *KLMS algorithm, 8*
1228 *KNLMS algorithm, 12*
1229 *KRLS algorithm, 14*
1230 *KRLS-T algorithm, 18, 26*
- 1231 *Linear adaptive filtering, 3*
1232 *Linear complexity, 16*
1233 *LMS algorithm, 5*
- 1234 *Mackey-Glass time series, 11, 16, 36*
- 1235 *Naive Online regularized Risk*
1236 *Minimization Algorithm, 8*
1237 *Nonlinear system identification, 8*
1238 *Novelty criterion, 21*
- 1239 *Online dictionary pruning, 22*
1240 *Online Gaussian processes, 25*
1241 *Online learning, 4*
1242 *Online sparsification, 15, 20*
1243 *Online system identification, 8*
- 1244 *Pruning, 22*
- 1245 *QKLMS algorithm, 13*
1246 *Quadratic complexity, 22*
- 1247 *Respiratory motion traces dataset, 30*
1248 *RLS algorithm, 5*
- 1249 *Sliding-window online learning, 23*
1250 *Sparse online Gaussian process*
1251 *algorithm, 25*
1252 *Stochastic gradient descent, 5*
1253 *SW-KRLS algorithm, 17*
1254 *system identification, 3*
- 1255 *Time-series prediction, 11*
1256 *Tracking, 7, 16*
1257 *Type-II ML, 27*

Cite as: Steven Van Vaerenbergh, "Adaptive Kernel Learning for Signal Processing". In J. L. Rojo-Álvarez, M. Martínez-Ramón, J. Muñoz-Mari, G. Camps-Valls (Eds.), *Digital Signal Processing with Kernel Methods*, pp. 387–431, Wiley-IEEE Press: Hoboken, NJ, USA, 2018.