



جامعة خليفة
Khalifa University

SENIOR RESEARCH PROJECT

On a Generalization of Kernel RLS to Nonlinear State-space Systems

Author

Ahmed A. RADHI¹
(100049714)

Supervisor

Prof. Ibrahim M. ELFADEL²

Co-Supervisor

Prof. Jorge P. ZUBELLI³

May 8, 2023

¹Department of Mathematics, Khalifa University, Abu Dhabi P.O. Box 127788, United Arab Emirates (Email: 100049714@ku.ac.ae)

²Department of Electrical Engineering and Computer Science, Khalifa University, Abu Dhabi P.O. Box 127788, United Arab Emirates (Email: ibrahim.elfadel@ku.ac.ae)

³Department of Mathematics, Khalifa University, Abu Dhabi P.O. Box 127788, United Arab Emirates (Email: jorge.zubelli@ku.ac.ae)

Abstract

In this paper, I provide a clear introduction to show how the recursive least squares (RLS) algorithm can be enhanced by utilizing a kernel function to produce a kernel nonlinear version known as the kernel recursive least squares (KRLS) algorithm. In addition to explaining the algorithm mathematically, the paper tests the KRLS algorithm by applying it in the domain of time-series prediction (TSP). Also, the paper compares the results obtained from the TSP task and compares it with other results reported previously. Furthermore, to extend the idea of using the kernel functions to produce kernel methods, the paper introduces a kernel adaptive filtering algorithm, namely the kernel adaptive autoregressive moving average (KAARMA) algorithm, which is an algorithm that can be used to predict the state variables of a system based on previous inputs and outputs. The paper explains analytically the KAARMA algorithm and compares it with the KRLS algorithm.

Keywords— kernel methods, kernel adaptive filtering (KAF), kernel adaptive autoregressive moving average (KAARMA), least squares, recursive least squares (RLS), kernel recursive least squares (KRLS)

Contents

1	Introduction	3
1.1	Overview	3
1.2	Contributions	3
2	ARMA Model	4
2.1	Learning Algorithm	4
2.2	Parameters Estimation	10
2.2.1	Least Squares Method	11
2.2.2	Recursive Least Squares (RLS) Method	13
3	KRLS Algorithm	16
3.1	Learning Algorithm	16
3.2	Example	18
4	KAARMA Algorithm	21
4.1	Learning Algorithm	21
4.2	Comparison with KRLS	23
5	Conclusion	26
6	References	27
7	Appendix	28
7.1	MATLAB Code for the Learning Algorithm	28
7.2	MATLAB Codes for the Least Squares Method	31
7.2.1	Least Squares Method Function	31
7.2.2	Least Squares Method Estimations Analysis	34
7.3	MATLAB Codes for the RLS Method	36
7.3.1	RLSE_Online.m Function	36
7.3.2	RLS Algorithm Analysis	37
7.4	MATLAB Codes for the KRLS Algorithm	40
7.4.1	KRLS MATLAB Functions	40
7.4.2	KRLS Santa Fe TSP	45

1 Introduction

1.1 Overview

Kernel methods represent a class of learning algorithms that uses the Mercer kernel function to produce nonlinear versions of familiar linear supervised and unsupervised learning algorithms. For instance, support vector machines (SVMs), which are well-known supervised learning algorithms used in classification and regression tasks, use kernel methods. Based on [1], Vladimir Vapnik, the developer of the SVM algorithm, observed that mapping data into a high-dimensional space may allow problems that are difficult to solve in low dimensions to be easily solved. Therefore, Vapnik introduced the “kernel trick”, which implies that input vectors that are applied to a Mercer kernel function can be mapped as an inner product into a high-dimensional space, usually called the Hilbert space \mathcal{H} or the feature space.

One of the main algorithms that the paper will be addressing is a kernel adaptive filtering (KAF) algorithm known as the kernel adaptive autoregressive moving average (KAARMA) algorithm. Similar to the kernel methods, the kernel adaptive filtering algorithms are learning algorithms that use the Mercer kernel function to develop a nonlinear version of the adaptive filtering algorithms. However, compared with the kernel methods, kernel adaptive filtering algorithms have different purposes as they are mainly used in signal processing tasks. In this case, kernel functions are used to produce a nonlinear version of the famous time series model called the autoregressive moving average (ARMA) model. Furthermore, the authors of [2] trained the adaptive filter using the stochastic gradient descent in the feature space; as a result, the algorithm was able to approximate any dynamical or nonlinear time-dependent relationships in the input space. Moreover, the authors compared the KAARMA algorithm to the recurrent neural network (RNN) by using them in solving grammatical inference problems. Due to the results they obtained, the authors observed that the KAARMA algorithm outperforms the first-order and second-order RNNs.

1.2 Contributions

This paper aims firstly, to derive the “linear version” of the KAARMA learning algorithm, the ARMA model, analyze the linear learning algorithm, and show that the learning algorithm is solvable. The paper will provide two ways to show

that the learning algorithm is solvable, particularly the least squares algorithm and the recursive least squares (RLS) algorithm. Secondly, before discussing the “nonlinear version” of the ARMA model, the paper aims to discuss and analyze the “kernelized” or “nonlinear version” of the recursive least squares algorithm, namely the kernel recursive least squares (KRLS) algorithm. Moreover, the KRLS algorithm will be experimentally tested on a famous time series data set that will be introduced later in the paper. Finally, the paper will provide a detailed comparison between both algorithms, the KAARMA algorithm and the KRLS algorithm, to show how they are related or unrelated.

2 ARMA Model

Before addressing the KAARMA algorithm, the “non-kernelized” version of the KAARMA algorithm will be addressed first. Therefore, this section of the paper will discuss the ARMA time-series model thoroughly by providing a detailed derivation of the learning algorithm and analyzing the obtained results.

2.1 Learning Algorithm

Define $x_i \in \mathbb{R}$ to be the internal state, $u_i \in \mathbb{R}$ to be the external input, $y_i \in \mathbb{R}$ to be the output, and $s_i \in \mathbb{R}^2$ to be the state vector, where $i = 0, 1, 2, 3, \dots$

Let $f, g, h \in \mathbb{R}$ be the parameters we are estimating.

Then, x_i, y_i , and s_i can be defined

$$x_i = fx_{i-1} + gu_i \tag{1}$$

$$y_i = hx_i \tag{2}$$

$$s_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix} = \begin{bmatrix} fx_{i-1} + gu_i \\ hx_i \end{bmatrix} \tag{3}$$

In the second row of equation (3), we replace x_i by equation (1) producing equation (4).

$$s_i = \begin{bmatrix} fx_{i-1} + gu_i \\ hfx_{i-1} + hgu_i \end{bmatrix} \quad (4)$$

Subsequently, we replace hx_{i-1} by y_{i-1} as indicated in equation (2) resulting in equation (5).

$$s_i = \begin{bmatrix} fx_{i-1} + gu_i \\ fy_{i-1} + hgu_i \end{bmatrix} \quad (5)$$

Obviously, we can take the common parameters f and g outside the vector leading to equation (6).

$$s_i = f \begin{bmatrix} x_{i-1} \\ y_{i-1} \end{bmatrix} + g \begin{bmatrix} 1 \\ h \end{bmatrix} u_i \quad (6)$$

Due to equation (3), we can replace the vector $\begin{bmatrix} x_{i-1} \\ y_{i-1} \end{bmatrix}$ by s_{i-1} resulting in equation (7).

$$\boxed{s_i = fs_{i-1} + g \begin{bmatrix} 1 \\ h \end{bmatrix} u_i} \quad (7)$$

Going back to equation (5), we can also express s_i in the following way:

$$s_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix} = \underbrace{\begin{bmatrix} f & 0 & g \\ 0 & f & gh \end{bmatrix}}_{\Omega^T} \underbrace{\begin{bmatrix} x_{i-1} \\ y_{i-1} \\ u_i \end{bmatrix}}_{\psi(s_{i-1}, u_i)} \quad (8)$$

Note that x_i in equation (1) represents the first-order ARMA algorithm with external input u_i , and s_i in equation (7) represents the ARMA algorithm operating on two variables.

Next, we define the cost function and error vector as follows, respectively,

$$\epsilon_i = \frac{1}{2}e_i^2 \quad (9)$$

$$e_i = d_i - y_i \quad (10)$$

where d_i is the desired output.

Our goal is to minimize the cost function with respect to the parameters f , g , and h ; hence, we apply the gradient descent as shown in the following steps.

First, we calculate $\nabla\epsilon_i$

$$\nabla\epsilon_i = \begin{bmatrix} \frac{\partial\epsilon_i}{\partial f} \\ \frac{\partial\epsilon_i}{\partial g} \\ \frac{\partial\epsilon_i}{\partial h} \end{bmatrix} \quad (11)$$

Note that $y_i = \underbrace{\begin{bmatrix} 0 & 1 \end{bmatrix}}_{\mathbb{I}} \underbrace{\begin{bmatrix} x_i \\ y_i \end{bmatrix}}_{s_i} = \mathbb{I}s_i$, then we calculate the partial derivatives indicated in equation (11).

$$\frac{\partial\epsilon_i}{\partial f} = \frac{\partial}{\partial f} \left(\frac{1}{2}e_i^2 \right) \quad (12)$$

$$= e_i \frac{\partial e_i}{\partial f} \quad (13)$$

$$= e_i \frac{\partial(d_i - y_i)}{\partial f} \quad (14)$$

$$= -e_i \frac{\partial y_i}{\partial f} \quad (15)$$

$$= -e_i \left(\frac{\partial y_i}{\partial s_i} \cdot \frac{\partial s_i}{\partial f} \right) \quad (16)$$

$$= -e_i \left(\mathbb{I} \underbrace{\begin{pmatrix} x_{i-1} \\ y_{i-1} \end{pmatrix}}_{s_{i-1}} \right) \quad (17)$$

$$= -e_i (\mathbb{I}s_{i-1}) \quad (18)$$

$$= -e_i y_{i-1} \quad (19)$$

$$\frac{\partial \epsilon_i}{\partial g} = \frac{\partial}{\partial g} \left(\frac{1}{2} e_i^2 \right) \quad (20)$$

$$= e_i \frac{\partial e_i}{\partial g} \quad (21)$$

$$= e_i \left(\frac{\partial e_i}{\partial y_i} \cdot \frac{\partial y_i}{\partial g} \right) \quad (22)$$

$$= e_i \left(\frac{\partial (d_i - y_i)}{\partial y_i} \cdot \frac{\partial (f y_{i-1} + h g u_i)}{\partial g} \right) \quad (23)$$

$$= e_i (-h u_i) \quad (24)$$

$$= -e_i (h u_i) \quad (25)$$

$$\frac{\partial \epsilon_i}{\partial h} = \frac{\partial}{\partial h} \left(\frac{1}{2} e_i^2 \right) \quad (26)$$

$$= e_i \frac{\partial e_i}{\partial h} \quad (27)$$

$$= e_i \left(\frac{\partial e_i}{\partial y_i} \cdot \frac{\partial y_i}{\partial h} \right) \quad (28)$$

$$= e_i \left(\frac{\partial (d_i - y_i)}{\partial y_i} \cdot \frac{\partial (h f x_{i-1} + h g u_i)}{\partial h} \right) \quad (29)$$

$$= e_i (-(f x_{i-1} + g u_i)) \quad (30)$$

$$= -e_i (f x_{i-1} + g u_i) \quad (31)$$

$$= -e_i (x_i) \quad (32)$$

$$= -e_i \left(\frac{y_i}{h} \right) \quad (33)$$

Therefore, $\nabla \epsilon_i$ can now be written as shown in equation (34):

$$\nabla \epsilon_i = \begin{bmatrix} \frac{\partial \epsilon_i}{\partial f} \\ \frac{\partial \epsilon_i}{\partial g} \\ \frac{\partial \epsilon_i}{\partial h} \end{bmatrix} = -e_i \begin{bmatrix} y_{i-1} \\ h u_i \\ \frac{y_i}{h} \end{bmatrix} \quad (34)$$

Now, the gradient descent will have the following form:

$$\begin{bmatrix} \Delta f_i \\ \Delta g_i \\ \Delta h_i \end{bmatrix} = -\eta \begin{bmatrix} \frac{\partial \epsilon_i}{\partial f} \\ \frac{\partial \epsilon_i}{\partial g} \\ \frac{\partial \epsilon_i}{\partial h} \end{bmatrix} \quad (35)$$

Finally, the learning algorithm can be formulated as the following explicit numerical scheme:

$$\begin{aligned} f_i &= f_{i-1} + \eta e_{i-1} y_{i-1} \\ g_i &= g_{i-1} + \eta e_{i-1} h_{i-1} u_i \\ h_i &= h_{i-1} + \eta e_{i-1} \left(\frac{y_{i-1}}{h_{i-1}} \right) \end{aligned} \quad (36)$$

As illustrated in the figure below, we can use the learning algorithm in equation (36) to estimate the system parameters (f, g, h) in real-time. In other words, the system parameters will be updated at every time instant based on the previous values, and the current or previous data.

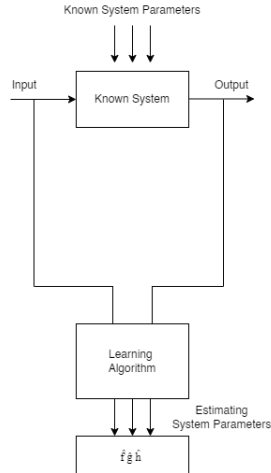


Figure 1: Illustrates How the Learning Algorithm Works.

A MATLAB code was implemented based on the derived learning algorithm in equation (36). As shown in Appendix 6.1, the MATLAB code starts by simulating random data and follows that by estimating the system parameters, in which the parameters were set as follows:

$$f = 0.5$$

$$g = 1.0$$

$$h = 2.0$$

Furthermore, the MATLAB code estimates the system parameter by iterating through 100 iterations. Based on Figure 2, it was concluded that the estimates are not stable, because when comparing the values obtained at different iterations, it was observed that estimates are changing significantly. Moreover, the estimates of f , g , and h obtained at the 100th iteration were the following:

$$f = 0.57268$$

$$g = 1.1701$$

$$h = 1.4245$$

Comparing the estimates of f and g with the exact values, error percentages of 14.5% and 17.0% were obtained for f and g , respectively, which are slightly large. In addition, the estimated h is far from the exact value with an error percentage of 28.8%.

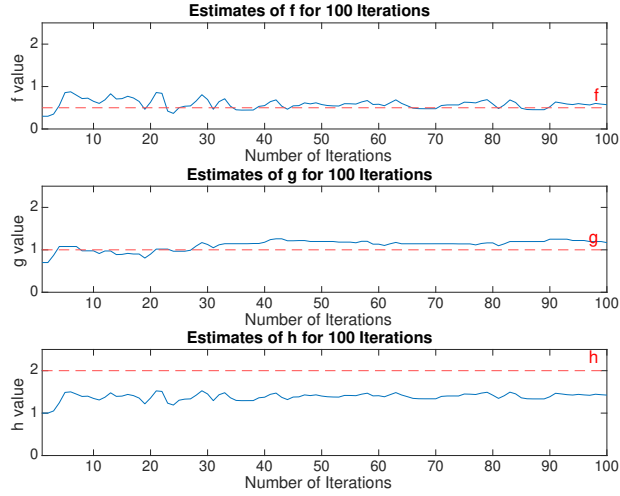


Figure 2: Comparing the Estimates (blue line) of f , g , and h with the Exact Value (red dashed-line) Using 100 Iterations.

Due to the instability observed in the expected values, other algorithms or optimization methods should be used to solve the learning algorithm and estimate the parameters.

2.2 Parameters Estimation

As it was used in the previous part, the error percentage will be applied as a metric to evaluate the performance of each method by comparing the exact parameters to their estimated values.

Furthermore, the error percentage can be calculated by using the following formula:

$$\% \text{ error} = \frac{|\text{estimated value} - \text{exact value}|}{\text{exact value}} \times 100\% \quad (37)$$

In addition, to estimate the unknown parameters using the other methods, the least squares and the recursive least squares (RLS), equation (2) will be reformulated into an equation with two unknowns (f and β) instead of three unknowns (f , g , and h).

To reformulate equation (2), equation (1) will be replaced in equation (2), and the following equations will demonstrate the remaining process thoroughly.

$$y_i = h(fx_{i-1} + gu_i) \quad (38)$$

$$= f \underbrace{hx_{i-1}}_{y_{i-1}} + hgu_i \quad (39)$$

$$= fy_{i-1} + \underbrace{hg}_{\beta=hg} u_i \quad (40)$$

$$y_i = fy_{i-1} + \beta u_i \quad (41)$$

Now, equation (41) has two unknowns, namely f and β , where $\beta = hg$. Therefore, the estimation methods should be applied to estimate f and β that solve equation (41).

2.2.1 Least Squares Method

Equation (41) will result in $n-1$ equations for $n \in \mathbb{N}$ representing the number of iterations.

$$\begin{aligned} y_2 &= f y_1 + \beta u_2 \\ &\vdots \\ y_n &= f y_{n-1} + \beta u_n \end{aligned} \quad (42)$$

Hence, the system of equations in equation (42) can be written as the matrix equation $Ax = b$.

$$\overbrace{\begin{bmatrix} y_1 & u_2 \\ \vdots & \vdots \\ y_{n-1} & u_n \end{bmatrix}}^A \overbrace{\begin{bmatrix} f \\ \beta \end{bmatrix}}^x = \overbrace{\begin{bmatrix} y_2 \\ \vdots \\ y_n \end{bmatrix}}^b \quad (43)$$

$(n-1) \times 2 \qquad 2 \times 1 \qquad (n-1) \times 1$

The system of equations in equation (43) is an over-determined system. Based on [3], a system of equations is considered over-determined when the number of equations is greater than the number of unknowns.

For instance, for $n = 4$, the following system of equations will be obtained. It is obvious that the system of equations has 3 equations, but 2 unknowns only.

$$\overbrace{\begin{bmatrix} y_1 & u_2 \\ y_2 & u_3 \\ y_3 & u_4 \end{bmatrix}}^A \overbrace{\begin{bmatrix} f \\ \beta \end{bmatrix}}^x = \overbrace{\begin{bmatrix} y_2 \\ y_4 \\ y_4 \end{bmatrix}}^b \quad (44)$$

$3 \times 2 \qquad 2 \times 1 \qquad 3 \times 1$

Therefore, for the error vector $\mathbf{e} = A\mathbf{x} - \mathbf{b}$, the least squares method can be used to find the vector \mathbf{x} that minimizes the norm of the error vector defined in the following equation

$$\|A\mathbf{x} - \mathbf{b}\|^2 = (A\mathbf{x} - \mathbf{b})^T (A\mathbf{x} - \mathbf{b}) \quad (45)$$

such that

$$\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b} \quad (46)$$

As shown in Appendix 6.2.1, a MATLAB function that takes two inputs was written. The function generates the data by iterating through a given number of iterations and estimates the unknown parameters based on a given number of data points.

Furthermore, a MATLAB program code is provided in Appendix 6.2.2. The code is written to go through 101 iterations four times and compare how estimating the unknown parameters using a higher or lower number of data points differs.

Data Points	f value	f Error %	β value	β Error %
25	0.48949	2.103	2.1523	7.6165
50	0.49054	1.892	2.0758	3.7911
75	0.49331	1.3388	2.0507	2.5351
100	0.50439	0.87887	1.9905	0.47517

Table 1: Estimates of f and β Using Different Number of Data Points

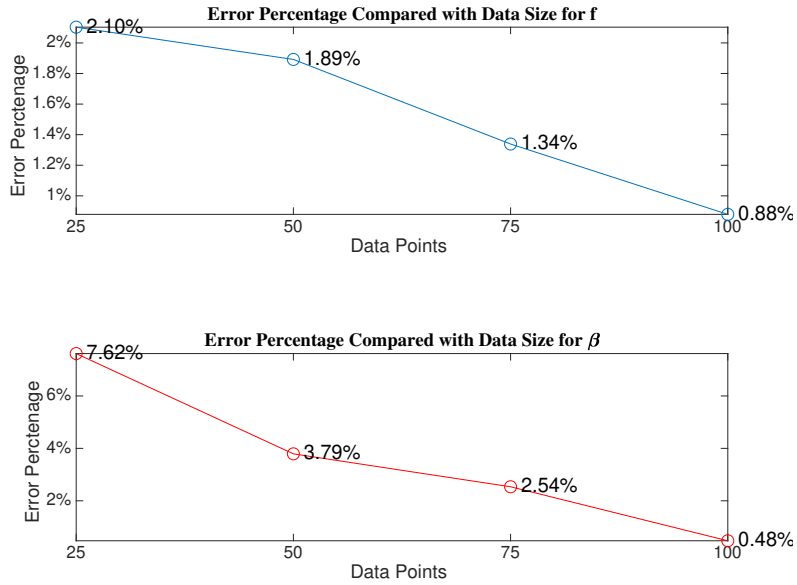


Figure 3: Error Percentages of f and β Estimates for Different Number of Data Points

Based on Table 1 and Figure 3, several findings can be realized. First, it was observed that usually higher number of data points result in a lower error percentage. Second, compared with Figure 2, lower error percentages were achieved by the least squares method. Third, more accurate estimations were obtained using the least squares method than the estimations reported in Figure 2.

Furthermore, the MATLAB code generates four different simulations of data and estimates the unknown parameters for each simulation of data using a given number of data points. Due to that, another run of the code may result in an error percentage achieved by a lower number of data points that is lower than an error percentage achieved by a higher number of data points. However, in all cases, it was observed that the estimations obtained are more accurate than those in Figure 2.

2.2.2 Recursive Least Squares (RLS) Method

As the name of the method may suggest, the recursive least squares (RLS) method is a recursive method that updates the \mathbf{x} vector that minimizes the norm of the error vector shown in equation (45) whenever a new data point is added.

As [3] explained, the idea of the RLS method is that it updates the previous estimate by adding a correction term based on the newly added data point resulting in a new estimate.

To explain the RLS method mathematically, for $k \in \mathbb{N}$, denote A_k as the A matrix with k elements and \mathbf{b}_k as the \mathbf{b} vector with k elements. Then A_k , \mathbf{b}_k , and \mathbf{x}_k can be defined as follows:

$$A_k = \begin{bmatrix} y_1 & u_2 \\ \vdots & \vdots \\ y_k & u_{k+1} \end{bmatrix}, \quad \mathbf{x}_k = \begin{bmatrix} f \\ \beta \end{bmatrix}, \quad \mathbf{b}_k = \begin{bmatrix} y_2 \\ \vdots \\ y_{k+1} \end{bmatrix} \quad (47)$$

If a new data point $(k+1)$ is added, then A_{k+1} and \mathbf{b}_{k+1} can be defined as follows:

$$A_{k+1} = \begin{bmatrix} A_k \\ a_{k+1}^T \end{bmatrix}, \quad \mathbf{b}_{k+1} = \begin{bmatrix} \mathbf{b}_k \\ y_{k+2} \end{bmatrix} \quad (48)$$

where a_{k+1} and y_{k+2} represent the new data point, and a_{k+1} can be expressed as the following equation:

$$a_{k+1} = \begin{bmatrix} y_{k+1} \\ u_{k+2} \end{bmatrix} \quad (49)$$

Furthermore, the inverse matrix P_k can be defined as the following equation:

$$P_k = \left[A_k^T A_k \right]^{-1} \quad (50)$$

And once a new data point (k+1) is added, the inverse matrix updates as expressed in the equation below:

$$P_{k+1} = P_k - P_k a_{k+1} \left[a_{k+1}^T P_k a_{k+1} + 1 \right]^{-1} a_{k+1}^T P_k \quad (51)$$

Obviously, the \mathbf{x} vector also gets updated, and \mathbf{x}_{k+1} can be written as the following equation:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \underbrace{P_{k+1} a_{k+1}}_{\text{gain matrix}} (y_{k+2} - a_{k+1}^T \mathbf{x}_k) \quad (52)$$

The gain matrix $P_{k+1} a_{k+1}$ can be written as expressed in the following equation:

$$\begin{aligned} K_{k+1} &= P_{k+1} a_{k+1} \\ &= P_k a_{k+1} \left[a_{k+1}^T P_k a_{k+1} + 1 \right]^{-1} \end{aligned} \quad (53)$$

Finally, equation (51) can be reformulated resulting in the following equation:

$$P_{k+1} = P_k - K_{k+1} a_{k+1}^T P_k \quad (54)$$

A MATLAB code that can be found in Appendix 6.3.2 is written to iterate through 101 iterations, assumes that $\mathbf{x}_1 = 0$, and estimates \mathbf{x}_k for $k \in [2, 101]$. The code is based on the “RLSE-Online” [3] function attached in Appendix 6.3.1. Furthermore, Table 2 shows below the estimated values f and β for \mathbf{x}_{26} , \mathbf{x}_{51} , \mathbf{x}_{76} , and \mathbf{x}_{101} .

k	f value	f Error %	β value	β Error %
26	0.48986	2.028	2.1498	7.4903
51	0.49348	1.3035	2.0786	3.931
76	0.49526	0.94742	2.0521	2.6032
101	0.49624	0.75272	2.0373	1.8651

Table 2: Estimates of f and β for Different Values of k

Below, Figure 4 compares the error percentages for the estimated parameters for x_{26} , x_{51} , x_{76} , and x_{101} , and shows how the error percentages decreased as the value of k increased. Moreover, comparing the error percentages in Table 2 with the values reported in Table 1, it can be observed that the error percentages are almost identical in most cases. This observation is due to the fact that both methods are based on the least squares method; however, the RLS method updates the estimates in real-time whenever a new data point arrives.

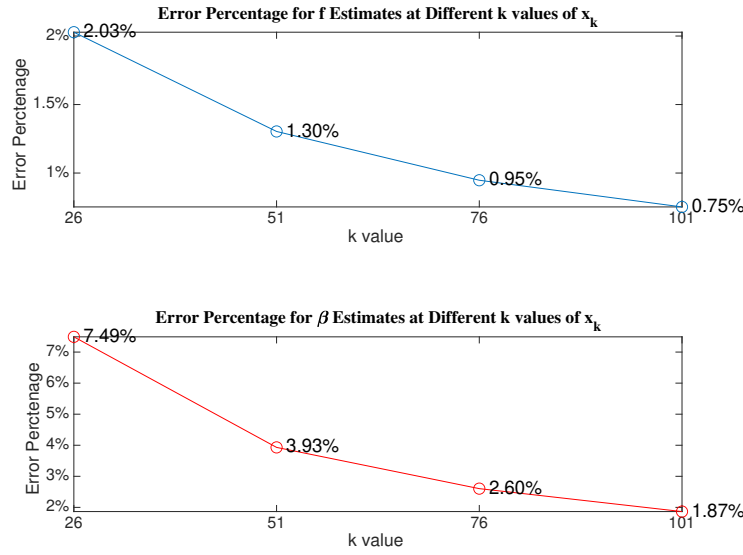


Figure 4: Error Percentages of f and β Estimates for Different Number of k Values

Furthermore, Figure 5 clearly displays the process of how the RLS method estimates the f and β as the number of iterations increases and compares the estimates to the exact values of the parameters.

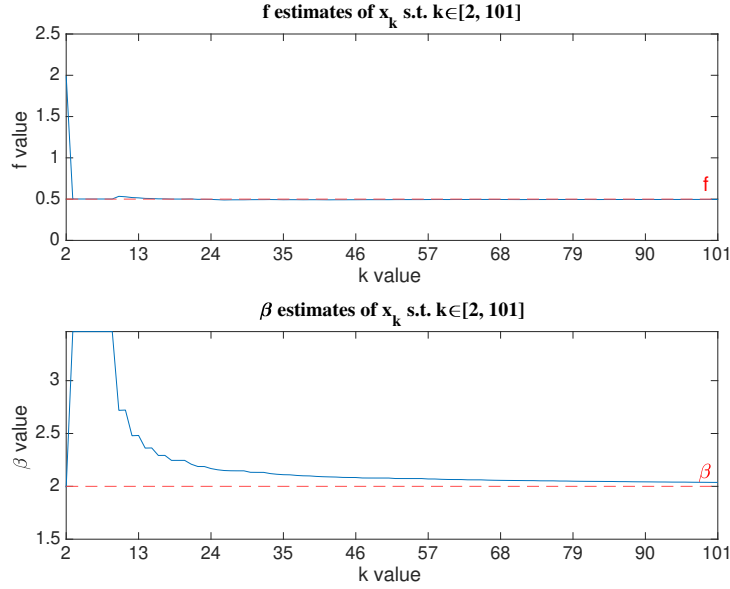


Figure 5: Comparing the Estimates (blue line) of f and β with the Exact Value (red dashed-line) Using 101 Iterations.

3 KRLS Algorithm

3.1 Learning Algorithm

The previous section introduced the recursive least squares (RLS) method and explained how the RLS method can be used to estimate unknown parameters. This section of the paper will introduce the “kernelized version” of the RLS and will show how “non-linearizing” the RLS can be helpful in addressing nonlinear problems.

In this section of the paper, a kernel function will be used to produce a nonlinear version of the RLS introduced in Section 2.2.2, namely the kernel recursive least squares (KRLS) algorithm. Usually, solutions obtained by the kernel methods are of the following form

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^{\ell} \alpha_i k(\mathbf{x}_i, \mathbf{x}) \quad (55)$$

where $\{\mathbf{x}_i\}_{i=1}^{\ell}$ are the training data points, $k(\cdot, \cdot)$ is the kernel function, and

α_i are the coefficients. There are different examples of kernel functions, and two well-known and frequently used kernel functions are the polynomial kernel and the Gaussian kernel. The kernel function that will be used throughout this paper is the Gaussian kernel, and it can be expressed as shown in equation (56).

$$k(\mathbf{x}, \mathbf{x}') = \exp \frac{-(\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}')}{2\sigma^2} \quad (56)$$

As [4] explained, some kernel methods obtain regularization by using sparsification, which is the process that results in the disappearance of many of the coefficients (α_i) in the solution defined in equation (55). There are several approaches for sparsification, but not all approaches are suitable for real-time online algorithms. The sparsification approach that the KRLS uses is an online constructive sparsification. The online sparsification approach compares the samples, and the samples that cannot be represented by a linear combination of the previously added samples are added into the kernel representation. In more details, after mapping the inputs into the feature space or the Hilbert space \mathcal{H} , that sparsification approach determines if the feature vector of a specific sample is approximately dependent on the other samples or not. If the feature vector is not approximately dependent on the previous samples, then it will be added to a dictionary \mathcal{D} , and its coefficient will be added to \hat{f} .

The solution of the kernel methods is previously expressed in equation (55); however, mapping the input vectors into the Hilbert space produces the following expression

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^{\ell} \alpha_i k(\phi(\mathbf{x}_i), \phi(\mathbf{x})) \quad (57)$$

In addition, the authors of [4] demonstrated more clearly how the sparsification process can be implemented mathematically. For input-output pairs $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots\}$, such that $x_i \in \mathcal{X}$ and $y_i \in \mathbb{R}$. At time step t , previous $t - 1$ training samples $\{\mathbf{x}_i\}_{i=1}^{t-1}$ have been observed, and a subset of the training samples is admitted into the dictionary $\mathcal{D}_{t-1} = \{\mathbf{x}_j\}_{j=1}^{m_{t-1}}$, where the feature vectors of the admitted training samples $\{\phi(\tilde{\mathbf{x}}_j)\}_{j=1}^{m_{t-1}}$ are linearly independent. If we considered an additional training sample \mathbf{x}_t , then it can be decided whether the feature vector $\phi(\mathbf{x}_t)$ is linearly independent or not by testing if it is approximately linearly dependent on the dictionary vectors or not. If the feature vector

is not *approximately linearly dependent*, then it will be added to the dictionary. As a result, the training samples $\mathbf{x}_1, \dots, \mathbf{x}_t$ can be approximated as a linear combination of the vectors in the dictionary \mathcal{D}_t .

Furthermore, the authors of [4] introduced the approximate linear dependence (ALD) condition and they stated that if the following condition

$$\delta_t \stackrel{\text{def}}{=} \min_{\mathbf{a}} \left\| \sum_{j=1}^{m_{t-1}} a_j \phi(\tilde{\mathbf{x}}_j) - \phi(\mathbf{x}_t) \right\| \leq \nu, \quad (58)$$

where ν is the accuracy parameter that determines the level of sparsity and $a_j = \tilde{\mathbf{K}}_{j-1}^{-1} \tilde{\mathbf{k}}_{j-1}(\mathbf{x}_j)$ such that $\tilde{\mathbf{K}}$ denotes the dictionary kernel matrix and $\tilde{\mathbf{k}}_{j-1}(\mathbf{x}_t) = \tilde{\mathbf{K}}_{j-1}^T \mathbf{a}_j$, was satisfied by finding coefficients $\mathbf{a} = (a_1, \dots, a_{m_{t-1}})^T$ that satisfies the condition, then adding the training sample \mathbf{x}_t to the dictionary will be avoided.

Generally, at each time step t , either one of two scenarios is faced:

1. If the feature vector $\phi(\mathbf{x}_t)$ is ALD on \mathcal{D}_{t-1} such that $\delta_t \leq \nu$, then \mathbf{x} is not added to the dictionary and $\mathbf{k}_{tt} = k(\mathbf{x}_t, \mathbf{x}_t)$. As a result, $\mathcal{D}_t = \mathcal{D}_{t-1}$, $m_t = m_{t-1}$, and $\tilde{\mathbf{K}}_t = \tilde{\mathbf{K}}_{t-1}$, in which \mathcal{D} denotes the dictionary, m denotes the dictionary size.
2. If the feature vector $\phi(\mathbf{x}_t)$ is not ALD on \mathcal{D}_{t-1} such that $\delta_t > \nu$, then \mathbf{x} is added to the dictionary. Consequently, $\mathcal{D}_t = \mathcal{D}_{t-1} \cup \{\mathbf{x}_t\}$, $m_t = m_{t-1} + 1$, and $\tilde{\mathbf{K}}_t$ grows as expressed in the following formula.

$$\tilde{\mathbf{K}}_t = \begin{bmatrix} \mathbf{K}_{t-1} & \tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t) \\ \tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t)^T & k_{tt} \end{bmatrix} \quad (59)$$

where $\tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t) = \tilde{\mathbf{K}}_{t-1}^T \mathbf{a}_t$, and $\mathbf{k}_{tt} = k(\mathbf{x}_t, \mathbf{x}_t)$

Furthermore, when comparing the dictionary kernel matrix $\tilde{\mathbf{K}}$ to the gain matrix K in the RLS, the dictionary kernel matrix is not exactly the same as the gain matrix. However, the dictionary kernel matrix plays a similar role.

3.2 Example

The authors of [4] experimentally tested the KRLS algorithm in time-series prediction (TSP). They applied the algorithm to the laser time-series data set

from the Santa Fe time-series competition[5]*. The authors used a training data set that consisted of 1000 samples as shown in Figure 6, and they predicted the next 100 samples of the time-series, in which those 100 samples represented the test data set. Due to the three intensity collapses that occur in the training data set shown in Figure 6 below, the chaotic laser time-series data set is considered to be difficult to predict.

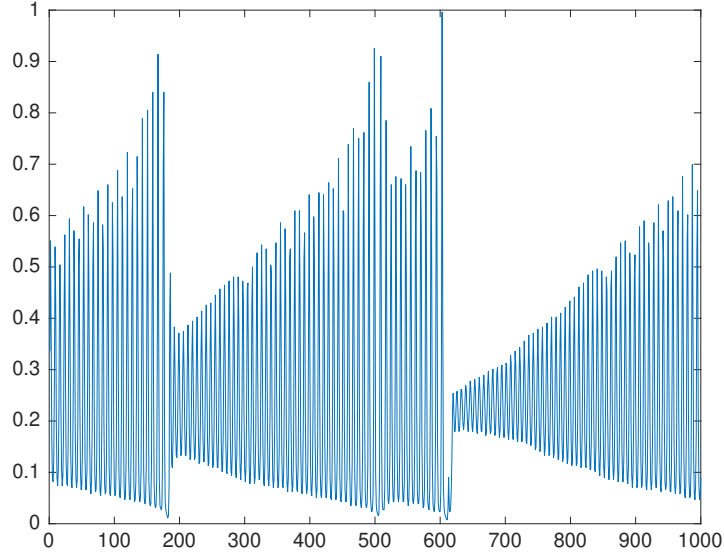


Figure 6: Laser Santa Fe Competition Time-Series Training Data Set

To compare their results to the competition results[6], the authors of [4] calculated the normalized mean squared error (NMSE) they achieved by the KRLS. [6] reported that out of 14 competitors, only two competitors obtained significant prediction accuracy, in which they achieved NMSE of 0.028 and 0.080, respectively. Furthermore, the winning entry obtained an NMSE of 0.028 by using a ‘highly specialized neural network architecture adapted by a temporal version of the back propagation algorithm’[4]. On the other hand, the authors of [4] achieved an NMSE of 0.026 using the KRLS algorithm, which is slightly better than that achieved by the winner in the Santa Fe competition.

This paper also implemented the KRLS algorithm on the same data set, and compared the achieved results to those reported in [4] and [6]. The MATLAB

*<http://www-psych.stanford.edu/~andreas/Time-Series/SantaFe.html>

codes[‡] in Appendix 6.4 used the KRLS algorithm to predict the 100 data points of the time-series. Furthermore, the code utilized the same parameters that the authors of [4] used in their approach. The following parameters were utilized:

$$\begin{aligned}\sigma &= 0.9 \\ \nu &= 0.01 \\ n &= 40\end{aligned}$$

in which they represent the kernel width, the ALD threshold parameter, and the auto-regressive window size, respectively.

Moreover, the MATLAB code addressed the time-series prediction task in two ways. First, normally, the MATLAB code used training set composed of 1000 data points shown in Figure 6 to predict the subsequent 100 data points, and achieved a NMSE of 0.069428 as shown in Figure 7. Comparing the achieved NMSE with the previously reported NMSEs, it can be observed that the NMSE achieved by the attached MATLAB code is worse than the NMSE obtained by the authors of [4] and the first competitor. However, the obtained NMSE is slightly better than the NMSE achieved by the second competitor.

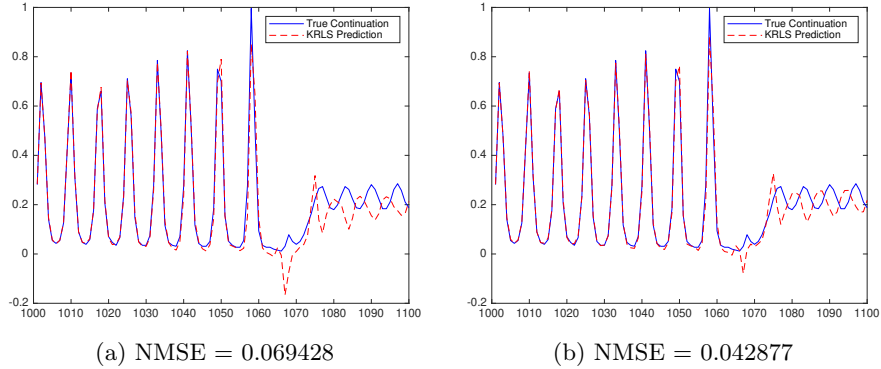


Figure 7: Comparing the 100 Predicted Data Points Using KRLS (red dashed-line) to the Exact Data Points (blue line).

Second, the code implemented another method to predict the time-series, which is the iterative prediction method. The iterative prediction method starts by using the training data as an input data set to predict an output, as usual. Next,

[‡]Following resources were used to implement the codes:

<http://web.mit.edu/wingated/www/resources.html>
<https://gist.github.com/caub/9462102>

the method augments the previously predicted output to the input data set and predicts another output, and this process keeps repeating to generate multiple outputs that make up the predicted time-series. The iterative prediction method predicted the time-series and obtained an NMSE of 0.042877 shown in Figure 7 above. Although the obtained NMSE improved over NMSE acquired by the first method, it is still worst than the NMSE achieved by the authors of [4] and the competition winner. Although the three NMSEs were obtained by the same algorithm, different values of the NMSE were obtained. This is clearly caused by the approach used to predict the time-series.

4 KAARMA Algorithm

4.1 Learning Algorithm

Previously, Section 2 derived, explained, and analyzed the time-series model known as the autoregressive moving average (ARMA) model. In this case, a kernel function is utilized and the ARMA model is trained in the Hilbert space \mathcal{H} to produce a kernel adaptive filtering (KAF) algorithm called the kernel adaptive autoregressive moving average (KAARMA) algorithm. In this section of the paper, the nonlinear version of the ARMA model, namely the kernel adaptive ARMA (KAARMA) algorithm, will be addressed.

As in the case of the kernel methods, solutions of a kernel adaptive filtering (KAF) algorithm can be expressed in equation (55). Compared with the previously addressed algorithms that are composed of inputs and outputs, the KAARMA algorithm has an additional state vector. Furthermore, the authors of [2] proposed a gradient-descent adaptive procedure that uses an input sequence and the observed outputs to learn the unknowns of a nonlinear system.

To demonstrate the Kernel Adaptive ARMA (KAARMA) algorithm mathematically, a dynamical system should be introduced by the following equations

$$\mathbf{x}_i = \mathbf{f}(\mathbf{x}_{i-1}, \mathbf{u}_i) \quad (60)$$

$$\mathbf{y}_i = \mathbf{h}(\mathbf{x}_i) \quad (61)$$

where $\mathbf{f}(\cdot, \cdot)$ is a continuous nonlinear state function and $\mathbf{h}(\cdot)$ is a continuous

nonlinear output function, and they can be expressed as follows

$$\begin{aligned}\mathbf{f}(\mathbf{x}_{i-1}, \mathbf{u}_i) &\triangleq [f^{(1)}(\mathbf{x}_{i-1}, \mathbf{u}_i), \dots, f^{(n_x)}(\mathbf{x}_{i-1}, \mathbf{u}_i)]^T \\ &= [\mathbf{x}_i^{(1)}, \dots, \mathbf{x}_i^{(n_x)}]^T\end{aligned}\quad (62)$$

$$\begin{aligned}\mathbf{h}(\mathbf{x}_i) &\triangleq [h^{(1)}(\mathbf{x}_i), \dots, h^{(n_y)}(\mathbf{x}_i)]^T \\ &= [\mathbf{y}_i^{(1)}, \dots, \mathbf{y}_i^{(n_y)}]^T\end{aligned}\quad (63)$$

in which $\mathbf{u}_i \in \mathbb{R}^{n_u}$ represents the input vector, $\mathbf{x}_i \in \mathbb{R}^{n_x}$ represents the state vector, $\mathbf{y}_i \in \mathbb{R}^{n_y}$ represents the output vector, (n_u, n_x, n_y) represents the dimensions for each vector, and the subscript within brackets $^{(k)}$ represents the k -th component of a vector.

In addition, the authors of [2] presented a “new hidden state” vector that can be expressed using equation (60) and equation (61) as shown in the following equation

$$\mathbf{s}_i \triangleq \begin{bmatrix} \mathbf{x}_i \\ \mathbf{y}_i \end{bmatrix} = \begin{bmatrix} \mathbf{f}(\mathbf{x}_{i-1}, \mathbf{u}_i) \\ \mathbf{h}(\mathbf{f}(\mathbf{x}_{i-1}, \mathbf{u}_i)) \end{bmatrix}\quad (64)$$

Using the expression previously stated in equation (64), y_i can be reformulated to the equation below

$$y_i = \underbrace{\begin{bmatrix} \mathbf{0} & \mathbf{I}_{n_y} \end{bmatrix}}_{\mathbb{I}} \underbrace{\begin{bmatrix} x_i \\ y_i \end{bmatrix}}_{\mathbf{s}_i} = \mathbb{I} \mathbf{s}_i\quad (65)$$

where \mathbb{I} is the fixed selector matrix, \mathbf{I}_{n_y} is a $n_y \times n_y$ identity matrix, $\mathbf{0}$ is an $n_y \times n_x$ zero matrix, and $\mathbf{s}_i \in \mathbb{R}^{n_s}$.

Moreover, to allow the KAARMA algorithm to learn the unknowns of a non-linear system, specifically the continuous state function $\mathbf{f}(\cdot, \cdot)$ and the output function $\mathbf{h}(\mathbf{f}(\cdot, \cdot))$, the authors of [2] applied the theory of the reproducing kernel Hilbert space \mathcal{H} . For instance, to apply the theory of RKHS to the system in equation (60) and equation (61), the inputs \mathbf{x}_{i-1} and \mathbf{u}_i will be mapped separately into the Hilbert space \mathcal{H} reproducing $\varphi(\mathbf{x}_{i-1}) \in \mathcal{H}_x$ and $\phi(\mathbf{u}_i) \in \mathcal{H}_u$, respectively. Next, the tensor product of the feature vectors will be taken to produce the tensor-product feature expressed in the following equation

$$\psi(\mathbf{x}_{i-1}, \mathbf{u}_i) \triangleq \varphi(\mathbf{x}_{i-1}) \otimes \phi(\mathbf{u}_i) \in \mathcal{H}_{su}\quad (66)$$

Furthermore, the authors [2] did not provide an analytical example of “how to” map the input vectors into the Hilbert space \mathcal{H} . However, the authors of [7] provided a concrete example that shows how the input vectors can be mapped into the Hilbert space \mathcal{H} resulting in a feature vector. However, the kernel used in their example was the polynomial kernel.

The authors of [2] also introduced the weights $\mathbf{\Omega}$ in which the equation of the weights $\mathbf{\Omega}_i$ at time i can be expressed as the linear combination of the previous features as stated in the following equation

$$\mathbf{\Omega}_i = \mathbf{\Psi}_i \mathbf{A}_i \quad (67)$$

where

$$\mathbf{\Psi}_i \triangleq [\psi(\mathbf{x}_{-1}, \mathbf{u}_0), \dots, \psi(\mathbf{x}_{m-2}, \mathbf{u}_{m-1})] \in \mathbb{R}^{n_\psi \times m}$$

is a vector of the m previous tensor-product features, and

$$\mathbf{A}_i \triangleq [\alpha_{i,1}, \dots, \alpha_{i,n_s}] \in \mathbb{R}^{m \times n_s}$$

are the corresponding coefficients, and m is the dictionary size in which $\mathbf{\Psi}_i$ represents the dictionary also.

Additionally, the authors of [2] introduced the cost function ϵ_i at time i , and it can be expressed as the following equation

$$\epsilon_i = \frac{1}{2} \mathbf{e}_i^T \mathbf{e}_i \quad (68)$$

where $\mathbf{e}_i = \mathbf{d}_i - \mathbf{y}_i \in \mathbb{R}^{n_y \times 1}$ is the error vector with \mathbf{d}_i representing the desired output.

By using the equation of the cost function expressed in equation (68), the error gradient with respect to $\mathbf{\Omega}_i$ can be written as the following equation

$$\frac{\partial \epsilon_i}{\partial \mathbf{\Omega}_i} = \frac{\partial \mathbf{e}_i^T \mathbf{e}_i}{2 \partial \mathbf{\Omega}_i} = \mathbf{e}_i^T \frac{\partial \mathbf{y}_i}{\partial \mathbf{\Omega}_i} \quad (69)$$

4.2 Comparison with KRLS

As the previous section of the paper showed, the KAARMA and the KRLS algorithms share some properties; however, the algorithms have various differences between them. This section of the paper will go through three major differences that can be observed when comparing the two algorithms.

As the authors of [2] stated, the state vector used in the KRLS algorithm can be expressed by

$$\mathbf{x}_i = \lambda^{-1/2} \mathbf{x}_{i-1}$$

such that $0 < \lambda \leq 1$ is a constant forgetting factor. On the other hand, as it was expressed in equation (60), the state vector utilized in the KAARMA algorithm can be written as

$$\mathbf{x}_i = \mathbf{f}(\mathbf{x}_{i-1}, \mathbf{u}_i)$$

where \mathbf{u}_i is the input vector and \mathbf{x}_{i-1} is the previous state vector.

Both, the KRLS and the KAARMA algorithms include a dictionary with size m . However, the dictionary \mathcal{D} of the KRLS algorithm works in a different way than the dictionary Ψ_i of the KAARMA algorithm. For instance, by observing the pseudo-codes below, it can be realized that the KRLS algorithm initializes $m = 1$, and then the ALD test decides whether m should be increased or not. If the ALD condition was not satisfied, a training sample x_t will be added to the dictionary, and m increases by 1. However, if ALD condition was not satisfied, then the dictionary size will not change. On the other hand, the KAARMA algorithm initializes $m = 1$: dictionary size; in other words, m can be any arbitrary value. Furthermore, m represents the numbers of past tensor-product features contained in Ψ_i .

When comparing the two algorithms with respect to the process of mapping the inputs into the Hilbert space \mathcal{H} , it can be observed that each algorithm has a different mapping process that produces a feature vector or a feature matrix. Both algorithms map the input vectors to the Hilbert space \mathcal{H} producing a feature vector $\varphi(\cdot)$. However, in the case of the KRLS algorithm, producing a feature vector means that the mapping process is done and the ALD test can carry on to check if the feature vectors satisfy the ALD condition or not, as shown in the pseudo-code of algorithm (1). On the other hand, in the case of the KAARMA algorithm, after mapping the two inputs into two separate Hilbert spaces \mathcal{H} producing two feature vectors $\varphi(\cdot)$ and $\phi(\cdot)$, a tensor product of the two feature vectors is taken to produce a tensor-product feature $\psi(\cdot, \cdot) = \varphi(\cdot) \otimes \phi(\cdot)$ that makes up the feature matrix Ψ , as shown in the pseudo-code of algorithm (2).

Algorithm 1 Kernel RLS Algorithm

Parameter: ν
Initialization:
 $\tilde{\mathbf{K}}_1 = [k_{11}]$
 $\tilde{\mathbf{K}}_1^{-1} = [\frac{1}{k_{11}}]$
 $\tilde{\alpha}_1 = (\frac{y_1}{k_{11}})$
 $\mathbf{P}_1 = [1]$
 $m = 1$
for $t = 2, 3, \dots$ **do**
 1. **Get new sample:** (x_t, y_t)
 2. **Compute** $\tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t)$
 3. **ALD test:**
 $\mathbf{a}_t = \tilde{\mathbf{K}}_{t-1} \tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t)$
 $\delta_t = k_{tt} - \tilde{\mathbf{k}}_{t-1}(\mathbf{x}_t)^T \mathbf{a}_t$
 if $\delta_t > \nu$ **then** % add \mathbf{x}_t to dictionary
 $\mathcal{D}_t = \mathcal{D}_{t-1} \cup \{\tilde{\mathbf{x}}_t\}$
 Compute $\tilde{\mathbf{K}}_{t-1}$
 Compute \mathbf{P}_t
 Compute $\tilde{\alpha}_t$
 $m := m + 1$
 else % dictionary unchanged
 $\mathcal{D}_t = \mathcal{D}_{t-1}$
 $\mathbf{q}_t = \frac{\mathbf{P}_{t-1} \mathbf{a}_t}{1 + \mathbf{a}_t^T \mathbf{P}_{t-1} \mathbf{a}_t}$
 Compute \mathbf{P}_t
 Compute $\tilde{\alpha}_t$
 end if
end for
Output: $\mathcal{D}_t, \tilde{\alpha}_t$

Algorithm 2 Kernel Adaptive ARMA Algorithm

Initialization:
 n_u : input dimension
 n_s : state dimension
 n_y : output dimension
 n_s : state kernel parameter
 n_u : input kernel parameter
 η : learning rate
Randomly initialize input $\mathbf{u}_0 \in \mathbb{R}^{1 \times n_u}$
Randomly initialize states \mathbf{s}_{-1} and $\mathbf{s}_0 \in \mathbb{R}^{1 \times n_s}$
Randomly initialize coefficient matrix $\mathbf{A} \in \mathbb{R}^{1 \times n_s}$
 $\Psi = [\psi(\mathbf{s}_{-1}, \mathbf{u}_0)]$: feature matrix
 $\mathbf{S} = [\mathbf{s}_{-1}]$: state dictionary
 $m=1$: dictionary size
 $\mathbb{I} = [\mathbf{0} \quad \mathbf{I}_{n_y}] \in \mathbb{R}^{n_y \times n_s}$: measurement matrix
Computation:
for time $t = 1, \dots, n$ **do**
 Initialization
 $\Psi' = []$: feature matrix update
 $\mathbf{S}' = []$: state matrix update
 $\mathbf{V}_1^{(k)} = \mathbf{I}_{n_s}^{(k)} \in \mathbb{R}^{n_s \times 1}$, for $k = 1, \dots, n_s$
 Update State-Transition Gradient Matrix
 for time $i = 1, \dots, t$ **do**
 Generate Next State
 $\mathbf{s}_i = \Omega^T \psi(\mathbf{s}_{i-1}, \mathbf{u}_i)$
 Update State Gradient
 $\mathbf{D}_i = [(\mathbf{S}^{(1)} - \mathbf{s}_{i-1}), \dots, (\mathbf{S}^{(m)} - \mathbf{s}_{i-1})]$
 $\mathbf{K}_i = \text{diag}(\Psi^T \psi(\mathbf{s}_{i-1}, \mathbf{u}_i))$
 $\Lambda_i = 2a_s \mathbf{A}^T \mathbf{K}_i \mathbf{D}_i^T$
 $\mathbf{V}_i^{(k)} = [\Lambda_i \mathbf{V}_{i-1}^{(k)}, \mathbf{I}_{n_s}^k]$
 Update Feature and State Matrices
 $\Psi' = [\Psi', \psi(\mathbf{s}_{i-1}, \mathbf{u}_i)]$
 $\mathbf{S}' = [\mathbf{S}', \mathbf{s}_{i-1}]$
 end for
 Prediction
 $\mathbf{y}_t = \mathbb{I} \mathbf{s}_t$
 Update Weights in the RKHS
 $\mathbf{e}_t = \mathbf{d}_t - \mathbf{y}_t$
 $\mathbf{A}^{(k)} = \begin{bmatrix} \mathbf{A}_i^{(k)} \\ \eta (\mathbf{I} \mathbf{V}_i^{(k)})^T \mathbf{e}_i \end{bmatrix}$
 $\Psi = [\Psi, \Psi']$
 $\mathbf{S} = [\mathbf{S}, \mathbf{S}']$
 $\Omega = \Psi \mathbf{A}$
 $m = m + t$
end for

5 Conclusion

To sum up, this paper presented a linear time series model called the autoregressive moving average (ARMA) model. Also, the paper derived a learning algorithm for the linear model and showed that it is solvable by estimating the unknown parameters using two methods, namely the least squares method and the recursive least squares (RLS) method. The paper also showed how a kernel function can be used to produce a kernel method as the kernel recursive least squares (KRLS) algorithm, which is the linear version of the RLS. The KRLS algorithm was used for time-series prediction (TSP), and the obtained results were compared with previously reported results. At the end, the showed how a kernel function can be utilized to produce a kernel adaptive filtering algorithm as the kernel adaptive autoregressive moving average (KAARMA) algorithm. The paper also compared between the KRLS algorithm and the KAARMA algorithm as it pointed to few major differences between the two algorithms.

6 References

- [1] V. N. Vapnik, *The nature of statistical learning theory*. Springer-Verlag New York, Inc., 1995.
- [2] K. Li and J. C. Príncipe, “The kernel adaptive autoregressive-moving-average algorithm,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 27, no. 2, pp. 334–346, 2016.
- [3] W. Y. Yang, W. Cao, T.-S. Chung, and J. Morris, *System of Linear Equations*. John Wiley & Sons, Ltd, 2005, ch. 2, pp. 71–115. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/0471705195.ch2>
- [4] Y. Engel, S. Mannor, and R. Meir, “The kernel recursive least-squares algorithm,” *IEEE Transactions on Signal Processing*, vol. 52, no. 8, pp. 2275–2285, 2004.
- [5] A. Weigend and N. Gerschenfeld, *Time Series Prediction: Forecasting the Future and Understanding the Past*, 06 1994, vol. 10.
- [6] A. Weigend and N. Gershenfeld, “Results of the time series prediction competition at the santa fe institute,” in *IEEE International Conference on Neural Networks*, 1993, pp. 1786–1793 vol.3.
- [7] J. C. Príncipe, W. Liu, and S. Haykin, *Kernel adaptive filtering: a comprehensive introduction*. John Wiley & Sons, 2011.

7 Appendix §

7.1 MATLAB Code for the Learning Algorithm

```
1 clc; clear; clf; close all;
2
3 %% LearningAlgorithmCode.m Explanation:
4
5 % Is the m file name suggest. This is a code that I
   wrote to implement the
6 % learning algorithm explained in the research paper.
   The code includes two
7 % phases. A phase of simulating data, and then a phase
   of estimating the
8 % system parameters f,g,h.
9
10 %% Generating Data Phase:
11
12 % To Control Random Number Generator
13 rng(10)
14
15 iter = 100;
16
17 x = zeros(iter,1);
18 y = zeros(iter,1);
19 s = zeros(2, iter);
20
21 x(1) = 2;
22 y(1) = 1;
23 s(:,1) = [x(1) ; y(1)];
24
25 f_exact = 0.5;
26 g_exact = 1;
27 h_exact = 2;
28
```

§All MATLAB codes are available in the following Github repository:
<https://github.com/7mxd/Senior-Research-Project-2022-2023>

```

29 % Using u as a random response instead of the
    commented impulse response
30
31 u = zeros(iter, 1);
32 u = randi([0,1], [1,iter]);
33
34 % for i = 1:iter
35 %     if i == 1
36 %         %u(i) = 1;
37 %     else
38 %         %u(i) = 0;
39 %     end
40 % end
41
42 d = zeros(iter, 1);
43
44 for i = 2:iter
45     x(i) = f_exact*x(i-1) + g_exact*u(i);
46     y(i) = h_exact*x(i);
47     s(:, i) = [x(i) ; y(i)];
48 end
49
50 for i = 1:iter
51     d(i) = y(i);
52 end
53
54 %% Estimating the System Parameters:
55 % Note: I used the explicit scheme.
56
57 n = 0.1; % learning rate
58 f = zeros(iter, 1);
59 g = zeros(iter, 1);
60 h = zeros(iter, 1);
61 f(1) = 0.3;
62 g(1) = 0.7;
63 h(1) = 1;

```

```

64
65
66 e = zeros(iter, 1);
67 yhat = zeros(iter, 1);
68 yhat(1) = 1;
69
70
71 for i = 2:iter % # of iterations
72     f(i) = f(i-1) + n*e(i-1)*yhat(i-1); % yhat(i-1) =
        I*s(:, i-1)
73     g(i) = g(i-1) + n*e(i-1)*h(i-1)*u(i);
74     h(i) = h(i-1) + n*e(i-1)*(yhat(i-1)/h(i-1));
75     yhat(i) = f(i)*yhat(i-1) + h(i)*g(i)*u(i);
76     e(i) = d(i) - yhat(i);
77 end
78
79 %% Plotting and Calculating Errors
80
81 Error_f = zeros(iter, 1);
82 Error_g = zeros(iter, 1);
83 Error_h = zeros(iter, 1);
84
85 for i = 1:iter
86     Error_f(i) = ((abs(f(i) - f_exact))/f_exact)*100;
87     Error_g(i) = ((abs(g(i) - g_exact))/g_exact)*100;
88     Error_h(i) = ((abs(h(i) - h_exact))/h_exact)*100;
89 end
90
91 Error = [Error_f, Error_g, Error_h];
92
93 figure('Name', 'Estimates of f,g,h')
94 % Figure 1 (Estimates of f,g,h for 100 Iterations)
95
96 subplot(3,1,1)
97 plot([1:100]', f)
98 yline(0.5, '--r', 'f')

```

```

99 axis([-inf inf 0 2.5])
100 xlabel('Number of Iterations')
101 ylabel('f Value')
102 title('Estimates of f for 100 Iterations')
103
104 subplot(3,1,2)
105 plot([1:100]', g)
106 yline(1, '--r', 'g')
107 axis([-inf inf 0 2.5])
108 xlabel('Number of Iterations')
109 ylabel('g Value')
110 title('Estimates of g for 100 Iterations')
111
112 subplot(3,1,3)
113 plot([1:100]', h)
114 yline(2, '--r', 'h')
115 axis([-inf inf 0 2.5])
116 xlabel('Number of Iterations')
117 ylabel('h Value')
118 title('Estimates of h for 100 Iterations')
119
120 table([1:100]', f, Error_f, g, Error_g, h, Error_h, '
    VariableNames', ...
121     {'iteration', 'f est.', 'f Error%', 'g est.', 'g
    Error%', 'h est.', ...
122     'h Error%'}))
123
124 % To export .eps figure
125 % print -depsc fghestimatesLA

```

7.2 MATLAB Codes for the Least Squares Method

7.2.1 Least Squares Method Function

```

1 %% Algorithm_LSE.m Explanation:
2

```



```

3  % A function that takes two inputs and returns three
   outputs
4
5  % Inputs:
6  % - iter : Number of Iterations for Generating Data
7  % - rows : Number of Data Points Used for Estimating x
8
9  % Note than x is called here as (z) to avoid some
   confusions when running
10 % the code
11
12 % Outputs:
13 % Az = b ; (x is replaced by z)
14 % Outputs are A matrix, z vector, and b vector.
15
16 %% Code Implementation
17
18 function [z,b,A] = Algorithm_LSE(iter, rows)
19
20
21 % Generating Data Phase:
22 x = zeros(iter,1);
23 y = zeros(iter,1);
24 s = zeros(2, iter);
25
26 x(1) = 2;
27 y(1) = 1;
28 s(:,1) = [x(1) ; y(1)];
29
30 f = 0.5;
31 g = 1;
32 h = 2;
33
34 u = zeros(iter, 1);
35 u = randi([0,1], [iter,1]);
36 % for i = 1:iter

```

```

37 %     if i == 1
38 %         %u(i) = 1;
39 %     else
40 %         %u(i) = 0;
41 %     end
42 % end
43
44 d = zeros(iter, 1);
45
46 for i = 2:iter
47     x(i) = f*x(i-1) + g*u(i);
48     y(i) = h*x(i);
49     s(:, i) = [x(i) ; y(i)];
50 end
51
52 for i = 1:iter
53     d(i) = y(i);
54 end
55
56 % Solving  b = Az
57
58
59 % for i=1:iter-1
60 %     b(i,1) = y(i+1);
61 % end
62 %
63 % for i=1:iter-1
64 %     A(i, 1) = y(i);
65 %     A(i, 2) = u(i+1);
66 % end
67
68 for i=1:rows
69     b(i,1) = y(i+1);
70 end
71
72 for i=1:rows

```

```

73     A(i,1) = y(i);
74     A(i,2) = u(i+1);
75 end
76
77 z = A\b;
78 end

```

7.2.2 Least Squares Method Estimations Analysis

```

1  clc; clear; clf;
2
3  % We have the following:
4                                     f = 0.5;
5                                     g = 1;
6                                     h = 2;
7  % Where beta is f*g as follows:
8                                     Beta = g*h;
9
10 % We are going through 100 iteration to generate the
    data
11
12 % Therefore, to check the impact of data size, I will
    keep increasing the
13 % number of rows to estimate f and B.
14
15 % I will test the following sizes
16                                     datasize = [25, 50, 75,
17                                               100];
18
19 % Set seed to get fixed results:
20 rng(1);
21 [z1,b1,A1] = Algorithm_LSE(101, 25);
22 ERR_PERC(1,1) = ((abs(z1(1) - f))/f)*100;
23 ERR_PERC(1,2) = ((abs(z1(2) - Beta))/Beta)*100;
24
25 [z2,b2,A2] = Algorithm_LSE(101, 50);

```

```

26 ERR_PERC(2,1) = ((abs(z2(1) - f))/f)*100;
27 ERR_PERC(2,2) = ((abs(z2(2) - Beta))/Beta)*100;
28
29 [z3,b3,A3] = Algorithm_LSE(101, 75);
30 ERR_PERC(3,1) = ((abs(z3(1) - f))/f)*100;
31 ERR_PERC(3,2) = ((abs(z3(2) - Beta))/Beta)*100;
32
33 [z4,b4,A4] = Algorithm_LSE(101, 100);
34 ERR_PERC(4,1) = ((abs(z4(1) - f))/f)*100;
35 ERR_PERC(4,2) = ((abs(z4(2) - Beta))/Beta)*100;
36
37 f_values = [z1(1), z2(1), z3(1), z4(1)];
38 beta_values = [z1(2), z2(2), z3(2), z4(2)];
39
40 table(datasize', f_values', ERR_PERC(:,1), beta_values
41       ', ERR_PERC(:,2), ...
42       'VariableNames', {'Data Points', 'f value', 'f
43       Error %', 'beta value', 'beta Error %'})
44
45 subplot(5,1,[1,2])
46 plot(datasize, ERR_PERC(:, 1), '-o')
47 xticks(datasize)
48 ytickformat('percentage')
49 axis([-inf inf -inf inf])
50 xlabel("Data Points")
51 ylabel("Error Perctenage")
52 for i=1:numel(datasize)
53     text(datasize(i) + 1.2, ERR_PERC(i,1), [num2str(
54         ERR_PERC(i,1), '%0.2f'), '%'])
55 end
56 title("Error Percentage Compared with Data Size for f
57       ", 'fontName', 'Times New Roman')
58
59 subplot(5,1,[4,5])
60 plot(datasize, ERR_PERC(:, 2), 'r-o')

```

```

58 xticks(datasize)
59 ytickformat('percentage')
60 axis([-inf inf -inf inf])
61 xlabel("Data Points")
62 ylabel("Error Perctenage")
63 for i=1:numel(datasize)
64     text(datasize(i) + 1.2, ERR_PERC(i,2), [num2str(
        ERR_PERC(i,2), '%0.2f'), '%'])
65 end
66 title("Error Percentage Compared with Data Size for \
    beta", 'fontName', 'Times New Roman')
67
68 % To export .eps figure
69 print -depsc LSEstimations

```

7.3 MATLAB Codes for the RLS Method

7.3.1 RLSE_Online.m Function

```

1  %% RLSE on-line processing
2  % Implemeneted by following instructions given by (
    APPLIED NUMERICAL
3  % METHODS USING MATLAB), pages 76-79
4
5  %% Explanation:
6
7  % at_k1 represents the added data point of A matrix
8  % b_k1 represents the added data point of b vector
9  % K represents the gain matrix
10 % P represents the inverse matrix "[A_k' A_k]^-1"
11
12 %% Defining RLSE_Online Function:
13
14 function [x, K, P] = RLSE_Online(aT_k1, b_k1, x, P)
15     K = P*aT_k1' / (aT_k1*P*aT_k1' + 1); % Equation
        (2.1.17) Page 78

```

```

16     x = x + K * (b_k1 - aT_k1*x); % Equation (2.1.16)
    Page 78
17     P = P - K*aT_k1*P; % Equation (2.1.18) Page 78
18 end

```

7.3.2 RLS Algorithm Analysis

```

1 %% Algorithm_RLSE Analysis:
2
3 clear; clc; clf; close all;
4
5 rng(1);
6
7 % Note that we have the following parameters:
8         %f = 0.5;
9         Beta = 2;
10
11 % and I will test the following numbers of data points
12         ;
13         k_data = [26, 51, 76,
14                   101];
15
16 %% Program Implementation:
17
18 % Initializing Data:
19 iter = 101; % Number of Iterations
20 x = zeros(iter,1); % `x` in learning algorithm
21 y = zeros(iter,1); % `y` in learning algorithm
22 x(1) = 2;
23 y(1) = 1;
24 u = randi([0,1], [iter, 1]);
25 f = 0.5;
26 g = 1;
27 h = 2;
28
29 z_o = [0.5 2]'; % True Values of the Parameters (f,
30               beta)

```

```

28 zsize = length(z_o);
29 z = zeros(zsize, 1);
30 P = 100*eye(zsize, zsize);
31
32 for k=1:(iter - 1)
33     x(k+1) = f*x(k) + g*u(k+1);
34     y(k+1) = h*x(k+1);
35
36     A(k, :) = [y(k) u(k+1)];
37     b(k, :) = y(k+1);
38     [z,K,P] = RLSE_Online(A(k,:), b(k,:), z, P); %
        Updating z vector (f, beta) estimates
39     online(:,k) = z;
40
41     ERR_PERC(k,1) = ((abs(online(1,k) - f))/f)*100;
42     ERR_PERC(k,2) = ((abs(online(2,k) - Beta))/Beta)
        *100;
43 end
44
45 %% Table & Plots:
46
47 table1 = table([2:101]', online(1,:)', ERR_PERC(:,1),
    online(2,:)', ERR_PERC(:,2), ...
48     'VariableNames', {'k', 'f value', 'f Error %', '
        beta value', 'beta Error %'})
49
50 figure('Name', 'Estimates of f and beta for different
    k values')
51 % Figure 1 (Estimates of f and beta for different data
    points)
52 subplot(2,1,1)
53 plot([2:101]', online(1,:)',
54     yline(0.5, '--r', 'f')
55     xlim([2 101])
56     xticks([2:11:102])
57     ylim([0 2.5])

```

```

58 xlabel('k value')
59 ylabel('f value')
60 title('f estimates of x_{k} s.t. k\in[2, 101]', '
      fontName', 'Times New Roman')
61
62 subplot(2,1,2)
63 plot([2:101]', online(2,:))
64 yline(2, '--r', '\beta')
65 xlim([2 101])
66 xticks([2:11:102])
67 ylim([1.5 inf])
68 xlabel('k value')
69 ylabel('\beta value')
70 title('\beta estimates of x_{k} s.t. k\in[2, 101]', '
      fontName', 'Times New Roman')
71
72 % To export .eps figure
73 print -depsc RLSEstimations1
74
75 figure('Name', "Error Percentage Compared with k
      values for f and beta")
76 % Figure 2 (Error Percentage Compared with Data Size
      for f and beta)
77 ERR = [ERR_PERC(25,:), ERR_PERC(50,:), ERR_PERC
      (75,:), ERR_PERC(100,:)'];
78 subplot(5,1,[1,2])
79 plot(k_data, ERR(1, :), '-o')
80 xticks(k_data)
81 ytickformat('percentage')
82 axis([-inf inf -inf inf])
83 xlabel("k value")
84 ylabel("Error Perctenage")
85 for i=1:numel(k_data)
86     text(k_data(i) + 1.2, ERR(1,i), [num2str(ERR(1,i),
      '%0.2f'), '%'])
87 end

```



```

88 title("Error Percentage for f Estimates at Different k
    values of x_{k}", 'fontName', 'Times New Roman')
89
90 subplot(5,1,[4,5])
91 plot(k_data, ERR(2, :), 'r-o')
92 xticks(k_data)
93 ytickformat('percentage')
94 axis([-inf inf -inf inf])
95 xlabel("k value")
96 ylabel("Error Perctenage")
97 for i=1:numel(k_data)
98     text(k_data(i) + 1.2, ERR(2,i), [num2str(ERR(2,i),
        '%0.2f'), '%'])
99 end
100 title("Error Percentage for \beta Estimates at
    Different k values of x_{k}", 'fontName', 'Times New
    Roman')
101
102 % To export .eps figure
103 print -depsc RLSEstimations2

```

7.4 MATLAB Codes for the KRLS Algorithm

7.4.1 KRLS MATLAB Functions

```

1 %% rbf4nn Function Explanation:
2
3 % Inputs:
4 % X -> a matrix that contains all samples as columns
5 % Y -> another matrix that contains all samples as
    columns
6 % sigsq -> sigma (squared) ; where sigma represents
    the kernal width
7
8 % Outputs:
9 % K_ij = exp(-1 / (2*sigsq) * (X_i' * Y_j))
10 % where K_ij is the Gaussian kernel function equation

```

```

11
12 function K = rbf4nn(X, Y, sigsq)
13 K = X' * Y;
14
15 Xsq = X .* X;
16 Xsum = sum(Xsq, 1);
17
18 Ysq = Y .* Y;
19 Ysum = sum(Ysq, 1);
20
21 K = K - Xsum' * ones(1, length(Ysum)) / 2;
22 K = K - ones(length(Xsum), 1) * Ysum / 2;
23 K = K ./ sigsq;
24 K = exp(K);
25 end

1 %% dict_init Function Explanation:
2
3 % Inputs:
4 % kfunc -> Handle to the kernel function.
5 % kparam -> The "kernel specific" parameter that will
   be passed to the kernel function on every call.
6 % thresh -> the almost-linearly-independent threshold
7 % state -> the initial data point (column vector)
8
9 % Outputs:
10 % dp -> A dictionary data structure. Used to be called
   in "dict.m".
11
12 function dp = dict_init(kfunc, kparam, thresh, state)
13 dp.kfunc = kfunc;
14 dp.kparam = kparam;
15 dp.thresh = thresh;
16
17 dp.Dict = state;
18
19 dp.K = feval(kfunc, state, state, kparam); % evaluate

```

```

    the kernel function on the three other inputs
20 dp.Kinv = 1 / dp.K; % inverted K
21
22 dp.addedFlag = 1;
23
24 return;
25 end

```

```

1 %% dict Function Explanation:
2
3 % Inputs:
4 % dp -> Data structure represting the problem setup,
   returned from
5 % dict_init.m or subsequent calls to this function (
   dict.m)
6 % state -> A new data point
7
8 % Outputs:
9 % dp -> An updated data structure.
10
11 function dp = dict(dp, state)
12 m = size(dp.Dict, 2); % Number of Enteries ; note that
   size(' ',2) returns the columns number
13
14 ktt = feval(dp.kfunc, state, state, dp.kparam);
15 ktwid = feval(dp.kfunc, dp.Dict, state, dp.kparam);
16
17 at = dp.Kinv * ktwid;
18 dt = ktt - ktwid' * at;
19
20 dp.addedFlag = 0;
21
22 if (dt > dp.thresh) % if NOT almost-linearly-dependent
   , add to dictionary
23     dp.Dict = [dp.Dict, state];
24     dp.K = [dp.K, ktwid; ktwid', ktt];
25     dp.Kinv = (1 / dt) * [dt * (dp.Kinv) + at * at', -

```

```

        at; -at', 1];
26     dp.addedFlag = 1;
27 end
28
29 dp.at = at;
30 dp.dt = dt;
31 dp.ktwid = ktwid;
32 kp.ktt = ktt;
33
34 return;
35 end

```

```

1  %% krls_init Function Explanation:
2
3  % Inputs:
4  % kfunc -> Handle to the kernel function.
5  % kparam -> The "kernel specific" parameter that will
        be passed to the kernel function on every call.
6  % thresh -> the almost-linearly-independent threshold
7  % state -> the initial data point (column vector)
8  % target -> the initial target (scalar)
9
10 % Outputs:
11 % kp -> A data structure that includes everything
        needed for subsequent KRLS calls (in krls.m)
12
13 function kp = krls_init(kfunc, kparam, thresh, state,
        target)
14 kp.kfunc = kfunc;
15 kp.kparam = kparam;
16 kp.thresh = thresh;
17
18 kp.dp = dict_init(kfunc, kparam, thresh, state);
19
20 kp.P = 1;
21 kp.Alpha = target / kp.dp.K;
22

```

```

23 return;
24 end

1 %% krls Function Explanation:
2
3 % Inputs:
4 % kp -> the data structure returned from "krls_init.m"
   or subsequent calls to this function (krls.m)
5 % state -> the next state vector (data point) (column
   vector)
6 % target -> the target value (scalar)
7
8 % Outputs:
9 % kp -> A data structure that includes everything
   needed for subsequent KRLS calls (in krls.m)
10
11 function kp = krls(kp, state, target)
12 m = size(kp.dp.Dict, 2); % Number of Entries ; note
   that size(' ',2) returns the columns number
13
14 kp.dp = dict(kp.dp, state);
15
16 at = kp.dp.at;
17 dt = kp.dp.dt;
18 ktwid = kp.dp.ktwid;
19 Kinv = kp.dp.Kinv;
20
21 if(kp.dp.addedFlag)
22     kp.P = [kp.P, zeros(m, 1); zeros(1, m), 1];
23     inno = (target - ktwid' * kp.Alpha) / dt;
24     kp.Alpha = [kp.Alpha - at * inno; inno];
25     kp.addedFlag = 1;
26 else
27     tmp = kp.P * at;
28     qt = tmp / (1 + at' * tmp);
29     kp.P = kp.P - qt * tmp';
30     kp.Alpha = kp.Alpha + Kinv * qt * (target - ktwid'

```

```

        * kp.Alpha);
31     kp.addedFlag = 0;
32 end
33
34 return;
35 end

```

```

1  %% krls_query Function Explanation:
2
3  % Inputs:
4  % kp -> the data structure returned from "krls_init.m"
   % or subsequent calls to the "krls.m" function
5  % state -> the vector you want to query
6
7  % Outputs:
8  % val -> the value of regression function at state
9
10 function val = krls_query(kp, state)
11 kernvals = feval(kp.kfunc, state, kp.dp.Dict, kp.
   kparam);
12
13 val = kernvals * kp.Alpha;
14
15 return;
16 end

```

7.4.2 KRLS Santa Fe TSP

```

1  clc; clear; clf; close all;
2
3  %% krls_santafe.m Explanation:
4
5  % This code is written to use the previous KRLS
   function and apply them to "The Santa Fe Time
   Series Competition Data Set A: Laser generated data
   ". Based on 1000 data points of the laser time-

```

```

        series, KRLS will be used to predict the next 100
        data points.
6
7 %% Load Data and Prepare Inputs
8
9 tic; % Time Started for NSME
10
11 % Read Content from Webpage:
12 weboptions('Timeout', 15);
13 data = [webread('https://web.archive.org/web
           /20160427182805if_/http://www-psych.stanford.edu
           :80/~andreas/Time-Series/SantaFe/A.dat') ...
14         webread('https://web.archive.org/web
           /20160427182805if_/http://www-psych.stanford.
           edu:80/~andreas/Time-Series/SantaFe/A.cont')];
15 y = textscan(data, '%f'); % Read Formatted Data (%f)
           from a text file (data)
16 y = y{1}(1:1100) / 256; % y{1} to access cell (access
           data)
17 figure('Name', 'Santa Fe Laset Time Series Data')
18 plot(y(1:1000))
19
20
21 %% Set an Auto-Regressive Matrix for the inputs
22 X = zeros(length(y), 40); % 1100 rows and 40 columms
23 for i = 1:size(X, 2) % 1:40
24     X(:, i) = [NaN(i,1); y(1:end-i)]; % 'NaN' on the
           Upper Triangle of the Matrix
25 end
26
27 %% Kernel Recursive Least Squares Regression
28
29 % Parameters:
30 % Here we test them directly. Normally we must
           allocate room for

```

```

31 % cross-validation in order to find optimal
    parameters in a grid. Once
32 % the best parameters are assessed, they an be
    evaluated on the
33 % training + validation sets
34
35 kernel_func = @rbf4nn; % Gaussian Kernel
36 kparam = 0.9; % Variance of the Gaussian
37 ald_thresh = 0.01; % almost-linearly-dependance
    threshold
38 n = 40; % Auto-Regressive Window Size
39
40 ltest = 100; % Test Set Size
41 Ytest = y(end - ltest + 1 : end); % Test Set
42 X_ = X(n + 1 : end - ltest, 1:n); % 960 x 40
43 Y_ = y(n + 1 : end - ltest); % 960 x 1
44
45 % We could try to train with missing values, starting
    at 2 instead of n+1
46 ltraining = size(X_, 1); % size of training data (960)
47 lvalidation = 0;
48
49 kp = krls_init(kernel_func, kparam, ald_thresh, X_
    (1,:) ', Y_(1));
50
51 for i = 2:ltraining
52     kp = krls(kp, X_(i,:) ', Y_(i));
53 end
54
55 v = [Y_(ltraining) X_(ltraining, 1:n-1)];
56
57 prediction = zeros(ltest, 1);
58 for j = 1:ltest
59     prediction(j) = krls_query(kp, v');
60     v = [prediction(j) v(1:n-1)];
61 end

```



```

62
63 testNSME = goodnessOfFit(Ytest, prediction, 'NMSE');
64 disp("NSME = " + testNSME)
65
66 figure('Name', 'KRLS predicting 100 steps into the
        future with NMSE = 0.069428')
67 plot(1000 + (1:length(Ytest)), Ytest, 'b-')
68 hold on
69 plot(1000 + (1:length(Ytest)), prediction, 'r--')
70 hold off
71 legend('True Continuation', 'KRLS Prediction')
72
73 toc % Time Ended for NMSE
74
75 %% Reinforcement Training
76
77 % To improve the performance of a machine learning
    algorithm, a technique
78 % called "Reinforcement Training" is used to
    iteratively add new training
79 % samples into the existing dataset and re-train the
    model on the
80 % "augmented" dataset.
81
82 tic % Time Started for Reinforcement Training
83
84 ireinforce = 40; % We will have 10 values of NMSE
85 nmse = repmat(testNSME, 1, ireinforce); % Replicate
    the 'testNMSE' on 1 x 10 matrix
86 prediction = repmat(prediction, 1, ireinforce); %
    Replicate the 'prediction' on 1 x 10 matrix
87
88 for i = 2:ireinforce
89     % Make new inputs by injecting the fitted values
90     fitted = zeros(ltraining, 1); % fitted values are
        composed of 100 x 1 column vector

```

```

91     for j = i:ltraining
92         fitted(j) = krls_query(kp, X_(j, :));
93     end
94     X_(2:end, 2:end) = X_(1:end-1, 1:end-1);
95     X_(i:end, 1) = fitted(1:end-i+1);
96     for j = i:ltraining
97         kp = krls(kp, X_(j, :)', Y_(j));
98     end
99
100     v = [Y_(ltraining) X_(ltraining, 1:n-1)];
101     prediction = zeros(ltest, 1);
102     for j = 1:ltest
103         prediction(j) = krls_query(kp, v');
104         v = [prediction(j) v(1:n-1)];
105     end
106
107     nmse(i) = goodnessOfFit(Ytest, prediction, 'NMSE'
108         );
109     predictions(:, i) = prediction; % Prediction
110                                     Vector that will contain prediction vectors for
111                                     each iteration
112
113 end
114
115 % Plot Best Reinforcement
116
117 [minNMSE, iNMSE] = min(nmse); % 'iNMSE' is the index
118                                of the minimum NMSE
119 disp("MIN NSME = " + minNMSE)
120
121 figure('Name', 'KRLS predicting 100 steps into the
122     future with NMSE = 0.042877')
123 plot(1000 + (1:length(Ytest)), Ytest, 'b-')
124 hold on
125 plot(1000 + (1:length(Ytest)), predictions(:, iNMSE),
126     'r--')
127 hold off

```

```
121 legend('True Continuation', 'KRLS Prediction')
122
123 toc % Time Ended for Reinforcement Training
```