

Subspace Identification of Hammerstein Systems Using Least Squares Support Vector Machines

Ivan Goethals, Kristiaan Pelckmans, Johan A. K. Suykens, and Bart De Moor

Abstract—This paper presents a method for the identification of multiple-input–multiple-output (MIMO) Hammerstein systems for the goal of prediction. The method extends the numerical algorithms for subspace state space system identification (N4SID), mainly by rewriting the oblique projection in the N4SID algorithm as a set of componentwise least squares support vector machines (LS-SVMs) regression problems. The linear model and static nonlinearities follow from a low-rank approximation of a matrix obtained from this regression problem.

Index Terms—Hammerstein models, least squares support vector machines, subspace identification.

I. INTRODUCTION

THROUGHOUT the last few decades, the field of linear modeling has been explored to the level that most linear identification problems can be solved efficiently with fairly standard and well known tools. Extensions to nonlinear systems are often desirable but in general much harder from a practical as well as a theoretical perspective. In many situations, Hammerstein systems are seen to provide a good tradeoff between the complexity of general nonlinear systems and interpretability of linear dynamical systems (see, e.g., [1]). They have been used e.g., for modeling biological processes [13], [32], chemical processes [7], and in signal processing applications [21]. Hammerstein models have also been shown to be useful for control problems (as, e.g., in [14]).

Identification of Hammerstein systems has been explored from different perspectives. Published approaches mainly differ in the way the static nonlinearity is represented and in the type of optimization problem that is finally obtained. Known approaches include the expansion of the nonlinearity as a sum

of (orthogonal or nonorthogonal) basis functions [16], [17], [15], the use of a finite number of cubic spline functions as presented in [6], piecewise linear functions [28] and neural networks [12]. Regardless of the parameterization scheme that is chosen, the final cost function will involve cross-products between parameters describing the static nonlinearity and those describing the linear dynamical system. Employing a maximum likelihood criterion results in a nonconvex optimization problem where global convergence is not guaranteed [20]. Hence, in order to find a good optimum for these techniques, a proper initialization is often necessary [5].

Different approaches were proposed in the literature to overcome this difficulty. These result in convex methods which generate models of the same, or almost the same quality as their nonconvex counterparts. Unfortunately, convexity is either obtained by placing heavy restrictions on the input sequence (e.g., whiteness) and the nonlinearity under consideration [2] or by using a technique known as overparameterization [3], [1]. In the latter, one replaces every cross-product of unknowns by new independent parameters resulting in a convex but overparameterized method. In a second stage the obtained solution is projected onto the Hammerstein model class using a singular value decomposition. A classical problem with the overparameterization approach is the increased variance of the estimates due to the increased number of unknowns in the first stage.

In [8] and [9], it was seen that by combining ideas from the overparameterization approach with concepts of least squares support vector machines (LS-SVMs), a Hammerstein autoregressive with exogenous inputs (ARX) identification algorithm was obtained which outperforms existing overparameterization approaches, mostly due to the effect of regularization. LS-SVMs [24], [23] are reformulations to the standard support vector machines (SVM). SVMs [29], [11], [19] and related methods constitute a powerful methodology for solving problems in linear and nonlinear classification, function approximation and density estimation and also stimulated new results in kernel based methods in general. They have been introduced on the interplay between learning theory, statistics, machine learning, neural networks and optimization theory.

A drawback with the method introduced in [8] is that the ARX model class is a rather restricted model class and is for instance not suitable to describe systems involving output noise. To this extent, identification algorithms based on state-space models are in many cases preferable. In this paper, we study the extension of the linear N4SID subspace identification algorithm to Hammerstein systems. It will be shown that by using the concept of componentwise LS-SVM regression, the state reconstruction step in classical identification algorithms can readily be extended to

Manuscript received May 28, 2004; revised March 14, 2005. Recommended by Guest Editor A. Vicino. This work was supported by Research Council KUL: GOA-Mefisto 666, GOA AMBioRICS, several Ph.D./postdoctoral and fellow grants; Flemish Government: FWO: Ph.D./postdoctoral grants, Projects, G.0240.99, G.0407.02, G.0197.02, G.0141.03, G.0491.03, G.0120.03, G.0452.04, G.0499.04, G.0211.05, G.0080.01 research communities (ICCoS, ANMMM, MLDM); AWI: Bil. Int. Collaboration Hungary/Poland; IWT: Ph.D. Grants, GBOU (McKnow), Belgian Federal Science Policy Office: IUAP P5/22; PODO-II; EU: FP5-Quprodis; ERNSI; Eureka 2063-IMPACT; Eureka 2419-FlITE; Contract Research/agreements: ISMC/IPCOS, Data4s, TML, Elia, LMS, Mastercard.

I. Goethals and J. Suykens are with the Department of Electrical Engineering ESAT-SCD, the Katholieke Universiteit Leuven (K. U. Leuven), B-3001 Leuven, Belgium, and also with the Fund for Scientific Research-Flanders (FWO-Vlaanderen) (e-mail: ivan.goethals@esat.kuleuven.be; johan.suykens@esat.kuleuven.be).

K. Pelckmans and B. De Moor are with the Department of Electrical Engineering ESAT-SCD, the Katholieke Universiteit Leuven (K. U. Leuven), B-3001 Leuven, Belgium (e-mail: kristiaan.pelckmans@esat.kuleuven.be; bart.demoor@esat.kuleuven.ac.be).

Digital Object Identifier 10.1109/TAC.2005.856647

Hammerstein systems. The linear system and static nonlinearity are recovered in a second step.

The outline of this paper is as follows. In Section II, the N4SID subspace algorithm for linear systems is reviewed briefly. Section III extends the N4SID algorithm toward a nonlinear setting using a variation on the theme of LS-SVMs. Section IV presents some illustrative examples for single-input–single-output (SISO) and multiple-input–multiple-output (MIMO) systems and relates the presented algorithm to existing approaches. A brief introduction into LS-SVM regression and component-wise LS-SVM regression is provided in Appendices I and II.

As a general rule in this paper, lowercase symbols will be used to denote column vectors. Uppercase symbols are used for matrices. Elements of matrices and vectors are selected using Matlab standards, e.g., $A(i, j)$ denotes the ij th entry of a matrix A , and $A(:, i)$ symbolizes the i th column of the same matrix. Estimates for a parameter x will be denoted by \hat{x} . The symbol \triangleq is used for definitions.

II. N4SID ALGORITHM FOR LINEAR SUBSPACE IDENTIFICATION

The subspace algorithm considered in this paper is the so-called N4SID algorithm, which is part of the set of combined deterministic-stochastic subspace algorithms as presented in [26] and [27]. We consider systems of the form

$$\begin{cases} x_{t+1} = Ax_t + Bu_t + \nu_t \\ y_t = Cx_t + Du_t + v_t \end{cases} \quad (1)$$

with $u_t \in \mathbb{R}^m$ and $y_t \in \mathbb{R}^l$ the input and output at time t , $x_t \in \mathbb{R}^n$ the state, and $\nu_t \in \mathbb{R}^n$ and $v_t \in \mathbb{R}^l$ zero mean white Gaussian noise vector sequences with covariance matrix

$$E \left\{ \begin{bmatrix} \nu_p \\ v_q \end{bmatrix} \begin{bmatrix} \nu_q^T & v_q^T \end{bmatrix} \right\} = \begin{bmatrix} Q & S \\ S^T & R \end{bmatrix} \delta_{pq}.$$

Given observed sequences $\{(u_t, y_t)\}_{t=0}^N$, N4SID identification algorithms are concerned with finding an estimate for the model order n of the system (1), estimates for the matrices $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{l \times n}$, $D \in \mathbb{R}^{l \times m}$ up to a similarity transformation, and the noise covariance matrices $Q \in \mathbb{R}^{n \times n}$, $S \in \mathbb{R}^{n \times l}$, and $R \in \mathbb{R}^{l \times l}$.

Block Hankel matrices play an important role in these algorithms. The input block Hankel matrices are defined as

$$U_{0|2i-1} \triangleq \begin{bmatrix} u_0 & u_1 & u_2 & \dots & u_{j-1} \\ u_1 & u_2 & u_3 & \dots & u_j \\ \vdots & \vdots & \vdots & & \vdots \\ u_{i-1} & u_i & u_{i+1} & \dots & u_{i+j-2} \\ u_i & u_{i+1} & u_{i+2} & \dots & u_{i+j-1} \\ u_{i+1} & u_{i+2} & u_{i+3} & \dots & u_{i+j} \\ \vdots & \vdots & \vdots & & \vdots \\ u_{2i-1} & u_{2i} & u_{2i+1} & \dots & u_{2i+j-2} \end{bmatrix} \\ \triangleq \begin{bmatrix} U_{0|i-1} \\ U_{i|2i-1} \end{bmatrix} \triangleq \begin{bmatrix} U_p \\ U_f \end{bmatrix} \triangleq \begin{bmatrix} U_{0|i} \\ U_{i+1|2i-1} \end{bmatrix} \triangleq \begin{bmatrix} U_p^+ \\ U_f^- \end{bmatrix}$$

with i and j user defined indexes such that $2i + j - 1 = N$. The output block Hankel matrices $Y_p, Y_f \in \mathbb{R}^{l \times j}$, $Y_p^+ \in \mathbb{R}^{(i+1)l \times j}$ and $Y_f^- \in \mathbb{R}^{(i-1)l \times j}$ are defined in a similar way. Finally

$$W_p \triangleq \begin{bmatrix} U_p \\ Y_p \end{bmatrix} \quad W_p^+ \triangleq \begin{bmatrix} U_p^+ \\ Y_p^+ \end{bmatrix}$$

are introduced as the past input–output block Hankel matrices and

$$\Gamma_i \triangleq \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{i-1} \end{bmatrix}$$

as the so-called extended observability matrix of order i . Defining $A^*/_{B^*}C^*$ as the oblique projection of the row space of a matrix A^* into the row space of a matrix C^* along the row space of a matrix B^* : $A^*/_{B^*}C^* = L_{C^*}C^*$ whereby

$$A^* = L_{B^*}B^* + L_{C^*}C^* + L_{B^*\perp, C^*\perp} \begin{bmatrix} B^* \\ C^* \end{bmatrix}^\perp$$

the main reasoning behind N4SID subspace algorithms follows from the fact that under the assumptions that

- 1) the process noise ν_t and measurement noise v_t are uncorrelated with the input u_t ;
- 2) the input u_t is persistently exciting of order $2i$, i.e., the input block Hankel matrix $U_{0|2i-1}$ is of full rank;
- 3) the sample size goes to infinity: $j \rightarrow \infty$;
- 4) the process noise ν_t and the measurement noise are not identically zero.

The following relation holds:

$$Y_{f/U_f}W_p = \Gamma_i \tilde{X}_i$$

with $Y_{f/U_f}W_p$ the so-called oblique projection of the future outputs Y_f onto the past data W_p along the future inputs U_f , which can be written explicitly as [27]

$$Y_{f/U_f}W_p = (Y_f - Y_f U_f^\dagger U_f) (W_p - W_p U_f^\dagger U_f)^\dagger W_p$$

where \tilde{X}_i can be shown to correspond to an estimate for the state in (1), resulting from a bank of nonsteady state Kalman filters [26]. Hence, the order of the system and a realization of the state can be obtained from a singular value decomposition of the oblique projection. Once the state is known, extraction of A, B, C and D is straightforward. Without going into further theoretical details of the N4SID algorithm (interested readers are referred to [25]–[27]), we summarize here a practical N4SID algorithm that will be used toward the Hammerstein model extension.

- 1) Calculate the oblique projections of the future outputs along the future inputs onto the past:

$$\begin{cases} \mathcal{O}_i = Y_{f/U_f}W_p \\ \mathcal{O}_{i+1} = Y_f^-/U_f^-W_p^+ \end{cases} \quad (2)$$

where $\mathcal{O}_i \in \mathbb{R}^{l \times j}$ and $\mathcal{O}_{i+1} \in \mathbb{R}^{(i+1)l \times j}$. This projection can be implemented using a least squares algorithm as follows:

$$\begin{aligned} (\hat{L}_u, \hat{L}_y) &= \underset{L_u, L_y}{\operatorname{argmin}} \left\| [L_u \quad L_y] \begin{bmatrix} U_p \\ U_f \\ Y_p \end{bmatrix} - Y_f \right\|_F^2 \\ (\hat{L}_u^-, \hat{L}_y^-) &= \underset{L_u^-, L_y^-}{\operatorname{argmin}} \left\| [L_u^- \quad L_y^-] \begin{bmatrix} U_p^+ \\ U_f^- \\ Y_p^+ \end{bmatrix} - Y_f^- \right\|_F^2. \end{aligned}$$

Estimates for \mathcal{O}_i and \mathcal{O}_{i+1} are then obtained as follows:

$$\begin{aligned} \mathcal{O}_i &= \hat{L}_u(:, 1 : im)U_p + \hat{L}_y Y_p \\ \mathcal{O}_{i+1} &= \hat{L}_u^-(:, 1 : (i+1)m)U_p^+ + \hat{L}_y^- Y_p^+. \end{aligned}$$

- 2) Calculate the SVD of the oblique projection \mathcal{O}_i , determine the order by inspecting the singular values and partition the SVD accordingly to obtain U_1 and S_1

$$\mathcal{O}_i = USV^T = [U_1 \quad U_2] \begin{bmatrix} S_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix}.$$

- 3) Determine the extended observability matrices Γ_i and Γ_{i-1} from

$$\Gamma_i = U_1 S_1^{1/2} \quad \Gamma_{i-1} = \Gamma_i(1 : l(i-1), :). \quad (3)$$

- 4) Determine estimates for the state sequences from the equations

$$\begin{cases} \hat{X}_i = \Gamma_i^\dagger \mathcal{O}_i \\ \hat{X}_{i+1} = \Gamma_{i-1}^\dagger \mathcal{O}_{i+1}. \end{cases}$$

- 5) Extract estimates for A, B, C and D from

$$\begin{bmatrix} \hat{X}_{i+1} \\ Y_{i|i} \end{bmatrix} = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} \hat{X}_i \\ U_{i|i} \end{bmatrix} + \begin{bmatrix} \rho_\nu \\ \rho_v \end{bmatrix}$$

by minimizing the least-squares residuals ρ_ν and ρ_v .

- 6) Determine estimates for Q, S and R from

$$\begin{bmatrix} \hat{Q} \\ \hat{S}^T \\ \hat{R} \end{bmatrix} = \frac{1}{j} \begin{bmatrix} \rho_\nu \\ \rho_v \end{bmatrix} \begin{bmatrix} \rho_\nu^T & \rho_v^T \end{bmatrix}.$$

The extension of this approach toward the identification of a Hammerstein system mainly concentrates on steps 1) and 5) where one uses the technique of componentwise LS-SVMs [18] instead.

III. EXTENDING THE N4SID ALGORITHM TOWARDS IDENTIFICATION OF HAMMERSTEIN MODELS

In this section, the linear N4SID algorithm will be extended to the identification of Hammerstein systems making use of the concept of overparameterization in an LS-SVM framework.

Equation (1) is transformed into a Hammerstein system by introducing a static nonlinearity $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$ which is applied to the inputs

$$\begin{cases} x_{t+1} = Ax_t + Bf(u_t) + \nu_t & \forall t = 0, \dots, N-1 \\ y_t = Cx_t + Df(u_t) + v_t & \forall t = 0, \dots, N-1. \end{cases} \quad (4)$$

Inputs and outputs $\{(u_t, y_t)\}_{t=0}^{N-1}$, are assumed to be available. The sequences of process and measurement noise $\{\nu_t\}_{t=0}^{N-1}$ and $\{v_t\}_{t=0}^{N-1}$ follow the same statistics as outlined in Section II. We define the matrix operator Φ as an operator on a block Hankel matrix and a nonlinear function ρ on \mathbb{R}^m which applies $\rho(\cdot)$ to every block matrix Z_i in Z and stacks the results in the original Hankel configuration

$$\begin{aligned} \Phi_{\mathcal{F}} \left(\begin{bmatrix} Z_1 & Z_2 & \dots & Z_p \\ Z_2 & Z_3 & \dots & Z_{p+1} \\ \vdots & \vdots & \ddots & \vdots \\ Z_q & Z_{q+1} & \dots & Z_{p+q-1} \end{bmatrix} \right) \\ = \begin{bmatrix} \mathcal{F}(Z_1) & \mathcal{F}(Z_2) & \dots & \mathcal{F}(Z_p) \\ \mathcal{F}(Z_2) & \mathcal{F}(Z_3) & \dots & \mathcal{F}(Z_{p+1}) \\ \vdots & \vdots & \ddots & \vdots \\ \mathcal{F}(Z_q) & \mathcal{F}(Z_{q+1}) & \dots & \mathcal{F}(Z_{p+q-1}) \end{bmatrix} \in \mathbb{R}^{qm \times p}. \end{aligned}$$

A. Overparameterization for the Oblique Projection \mathcal{O}_i

The oblique projection $\mathcal{O}_i = Y_f / U_f W_p$ can be calculated from estimates for L_u, L_y and f obtained by minimizing the residuals E of the following equation [27]:

$$Y_f = [L_u \quad L_y] \begin{bmatrix} \Phi_f(U_{0|2i-1}) \\ Y_p \end{bmatrix} + E \quad (5)$$

in a least-squares sense. This can be rewritten as

$$\begin{aligned} Y_f(s, t) &= L_y(s, :) Y_p(:, t) \\ &+ \sum_{h=1}^{2i} L_u(s, (h-1)m+1 : hm) f(u_{h+t-2}) + E(s, t) \end{aligned} \quad (6)$$

for $s = 1, \dots, il$ and $t = 1, \dots, j$. Once estimates for \hat{L}_u, \hat{L}_y and \hat{f} occurring in (5) and (6) are obtained, the oblique projection is calculated as

$$\begin{aligned} \mathcal{O}_i(s, t) &= \hat{L}_y(s, :) Y_p(:, t) \\ &+ \sum_{h=1}^i \hat{L}_u(s, (h-1)m+1 : hm) \hat{f}(u_{h+t-2}) \end{aligned} \quad (7)$$

for $s = 1, \dots, il$ and $t = 1, \dots, j$. Note that in (6) and (7), products between parameter matrices L_u and L_y and the static nonlinearity f appear which are hard to incorporate in an optimization problem. In order to deal with the resulting nonconvexity, we apply the concept of overparameterization (see Appendix II) by introducing a set of functions $g_{h,s} : \mathbb{R}^m \rightarrow \mathbb{R}$ such that [1]

$$g_{h,s} \triangleq c_{h,s}^T f \quad \text{s.t.} \quad c_{h,s}^T = L_u(s, (h-1)m+1 : hm) \quad (8)$$

for $h = 1, \dots, 2i$ and $s = 1, \dots, il$. With these new functions we obtain a generalization to (6) and (7)

$$Y_f(s, t) = L_y(s, :)Y_p(:, t) + \sum_{h=1}^{2i} g_{h,s}(u_{h+t-2}) + E(s, t) \quad (9)$$

$$\mathcal{O}_i(s, t) = \hat{L}_y(s, :)Y_p(:, t) + \sum_{h=1}^i \hat{g}_{h,s}(u_{h+t-2}) \quad (10)$$

for $s = 1, \dots, il$ and $t = 1, \dots, j$. Note that (9) is now linear in the functions $g_{h,s} : \mathbb{R}^m \rightarrow \mathbb{R}$. The central idea behind the algorithm presented in this paper is that \hat{L}_y and estimates for the functions $g_{h,s}$ in (9)–(10) can be determined from data using the concept of componentwise LS-SVM regression as presented in Appendix II.

Let the kernel function be defined as $K : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ such that $K(u_p, u_q) = \varphi(u_p)^T \varphi_k(u_q)$ for all $p, q = 0, \dots, N-1$ and the kernel matrix $\Omega \in \mathbb{R}^{N \times N}$ such that $\Omega(i, j) = K(u_{i-1}, u_{j-1})$ for all $i = 1, \dots, N$, $j = 1, \dots, N$. Substituting $g_{h,s}$ for the primal model $w_{h,s}^T \varphi$ (see Appendix II) in (9) results in

$$Y_f(s, t) = L_y(s, :)Y_p(:, t) + \sum_{h=1}^{2i} w_{h,s}^T \varphi(u_{h+t-2}) + E(s, t) \quad \forall s = 1, \dots, li, \quad t = 1, \dots, j. \quad (11)$$

As argued in Appendix II, the expansion of a nonlinear function as the sum of a set of nonlinear functions is not unique, e.g.,

$$(w_1^T \varphi_1(u)) + (w_2^T \varphi_2(u)) = (w_1^T \varphi_1(u) + \delta) + (w_2^T \varphi_2(u) - \delta)$$

for all $\delta \in \mathbb{R}$. It was seen that this problem can be avoided by including centering constraints of the form

$$\sum_{t=0}^{N-1} f(u_t) = 0. \quad (12)$$

This constraint can always be applied since for any constant δ_u , and any function $\underline{f} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ such that $f = \underline{f} + \delta_u$ there exists a state transformation $\xi_t = \Psi(x_t)$ with $\Psi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and a constant δ_y such that (4) is transformed as follows:

$$\begin{cases} \xi_{t+1} = A\xi_t + B\underline{f}(u_t) + \nu_t \\ y_t - \delta_y = C\xi_t + D\underline{f}(u_t) + v_t \end{cases} \quad (13)$$

with $\xi_t \in \mathbb{R}^n$ and $\delta_y \in \mathbb{R}^l$ defined as

$$\begin{cases} \xi_t = \Psi(x_t) = x_t - (I - A)^{-1} B \delta_u \\ \delta_y = (C(I - A)^{-1} B + D) \delta_u. \end{cases}$$

Hence, the constraint (12) can be applied provided that a new parameter δ_y is added to the model, transforming (11) into

$$Y_f(s, t) + [1_i \otimes \delta_y](s) = L_y(s, :)Y_p(:, t) + 1_i \otimes \delta_y + \sum_{h=1}^{2i} w_{h,s}^T \varphi(u_{h+t-2}) + E(s, t) \quad \forall s = 1, \dots, li, t = 1, \dots, j$$

where \otimes denotes the matrix kronecker product. Through the equality $w_{h,s}^T \varphi = g_{h,s} = c_{h,s}^T f$ for all $h = 1, \dots, 2i, s = 1, \dots, il$, the constraint (12) amounts to

$$\sum_{t=0}^{N-1} w_{h,s}^T \varphi(u_t) = 0 \quad \forall h, s.$$

The LS-SVM primal problem is then formulated as a constrained optimization problem

$$\begin{aligned} \min_{w_{h,s}, L_y, E, \delta_y} \mathcal{J}(w_{h,s}, L_y, E, \delta_y) \\ = \frac{1}{2} \sum_{s=1}^{il} \sum_{h=1}^{2i} w_{h,s}^T w_{h,s} + \frac{\gamma}{2} \sum_{s=1}^{il} \sum_{t=1}^j E(s, t)^2 \\ \text{s.t.} \begin{cases} Y_f(s, t) + [1_i \otimes \delta_y](s) \\ = L_y(s, :)Y_p(:, t) + 1_i \otimes \delta_y \\ + \sum_{h=1}^{2i} w_{h,s}^T \varphi(u_{h+t-2}) + E(s, t) \\ \forall s = 1, \dots, il, t = 1, \dots, j \\ \sum_{t=0}^{N-1} w_{h,s}^T \varphi(u_t) = 0 \\ \forall h = 1, \dots, 2i, \quad s = 1, \dots, li. \end{cases} \end{aligned} \quad (14)$$

Lemma 3.1: Given the primal problem (14), estimates for L_y and δ_y follow from the dual system:

$$\begin{bmatrix} 0 & 0 & 1^T & 0 \\ 0 & 0 & Y_p & 0 \\ 1 & Y_p^T & \mathcal{K}_p + \mathcal{K}_f + \gamma^{-1}I & \mathcal{S} \\ 0 & 0 & \mathcal{S}^T & \mathcal{T} \end{bmatrix} \begin{bmatrix} \bar{d} \\ \frac{L_y^T}{\mathcal{A}} \\ \frac{\mathcal{A}}{\mathcal{B}} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \frac{Y_f^T}{\mathcal{B}} \\ 0 \end{bmatrix}$$

where $\bar{d} = (1_i \otimes I_l - L_y(1_i \otimes I_l))\delta_y$, 1_j is a column vector of length j with elements 1, $\mathcal{T} = I_{2i} \times 1_N^T \Omega^1 1_N$, $\mathcal{S}_q = \sum_{t=1}^N \Omega(t, q)$ and

$$\begin{aligned} \mathcal{A} &= \begin{bmatrix} \alpha_{1,1} & \alpha_{2,1} & \dots & \alpha_{li,1} \\ \alpha_{1,2} & \alpha_{2,2} & \dots & \alpha_{li,2} \\ \vdots & \vdots & & \vdots \\ \alpha_{1,j} & \alpha_{2,j} & \dots & \alpha_{li,j} \end{bmatrix} \\ \mathcal{B} &= \begin{bmatrix} \beta_{1,1} & \beta_{1,2} & \dots & \beta_{1,li} \\ \beta_{2,1} & \beta_{2,2} & \dots & \beta_{2,li} \\ \vdots & \vdots & & \vdots \\ \beta_{2i,1} & \beta_{2i,2} & \dots & \beta_{2i,li} \end{bmatrix} \\ \mathcal{S} &= \begin{bmatrix} \mathcal{S}_1 & \mathcal{S}_2 & \dots & \mathcal{S}_{2i} \\ \mathcal{S}_2 & \mathcal{S}_3 & \dots & \mathcal{S}_{2i+1} \\ \vdots & \vdots & & \vdots \\ \mathcal{S}_j & \mathcal{S}_{j+1} & \dots & \mathcal{S}_N \end{bmatrix}. \end{aligned}$$

The matrices $\mathcal{K}_p \in \mathbb{R}^{j \times j}$ and $\mathcal{K}_f \in \mathbb{R}^{j \times j}$ have elements

$$\begin{aligned} \mathcal{K}_p(p, q) &= \sum_{h=1}^i K(u_{h+p-2}, u_{h+q-2}) \\ \mathcal{K}_f(p, q) &= \sum_{h=i+1}^{2i} K(u_{h+p-2}, u_{h+q-2}) \end{aligned}$$

for all $p, q = 1, \dots, j$. Estimates for the $g_{h,s}$ in (9) are given as:

$$\hat{g}_{h,s} : \mathbb{R}^m \rightarrow \mathbb{R} : u^* \rightarrow \sum_{t=1}^j \alpha_{s,t} K(u_{h+t-2}, u^*) + \beta_{h,s} \sum_{t=0}^{N-1} K(u_t, u^*) \quad \forall h, s. \quad (15)$$

Proof: This directly follows from the Lagrangian:

$$\begin{aligned} \mathcal{L}(w, d, L_y, E; \alpha, \beta, d_s) = & \mathcal{J}(w, E) \\ & - \sum_{h=1}^{2i} \sum_{s=1}^{li} \beta_{h,s} \left\{ \sum_{t=0}^{N-1} w_{h,s}^T \varphi(u_t) \right\} \\ & - \sum_{s=1}^{li} \sum_{t=1}^j \alpha_{s,t} \left\{ L_y(s, :) Y_p(:, t) \right. \\ & \left. + \sum_{h=1}^{2i} w_{h,s}^T \varphi(u_{h+t-2}) + d_s + E(s, t) - Y_f(s, t) \right\} \end{aligned}$$

with $\bar{d} = [d_1^T \dots d_l^T]^T$ by taking the conditions for optimality $(\partial \mathcal{L})/(\partial w_{h,s}) = 0, (\partial \mathcal{L})/(\partial L_y(s, :)) = 0, (\partial \mathcal{L})/(\partial E(s, t)) = 0, (\partial \mathcal{L})/(\partial d_s) = 0, (\partial \mathcal{L})/(\partial \alpha_{s,t}) = 0, (\partial \mathcal{L})/(\partial \beta_{h,s,k}) = 0, (\partial \mathcal{L})/(\partial d_s) = 0$ and after elimination of the primal variables $w_{h,s}$ and E . ■

Combining the results from Lemma (3.1) with (10), we have

$$\begin{aligned} \mathcal{O}_i = & \sum_{h=1}^i \left(\begin{bmatrix} \Phi_{\hat{g}_{h,1}} U_{h|h} \\ \Phi_{\hat{g}_{h,2}} U_{h|h} \\ \vdots \\ \Phi_{\hat{g}_{h,2li}} U_{h|h} \end{bmatrix} \right) + \hat{L}_y \left(Y_p - 1_{li} 1_{li}^T \otimes \hat{\delta}_y \right) \\ = & \mathcal{A}^T \mathcal{K}_p + \mathcal{B}_p^T \mathcal{S}_p^T + \hat{L}_y \left(Y_p - (1_i 1_j^T) \otimes \hat{\delta}_y \right) \quad (16) \end{aligned}$$

with $\mathcal{B}_p = \mathcal{B}(1 : i, :)$ and $\mathcal{S}_p = \mathcal{S}(:, 1 : i)$.

B. Calculating the Oblique Projection \mathcal{O}_{i+1}

The calculation of \mathcal{O}_{i+1} is entirely equivalent to that of \mathcal{O}_i . Without further proof, we state that \mathcal{O}_{i+1} is obtained as

$$\begin{aligned} \mathcal{O}_{i+1} = & (\mathcal{A}^-)^T (\mathcal{K}_p^+)^T + (\mathcal{B}_p^-)^T (\mathcal{S}_p^-)^T \\ & + L_y^- \left(Y_p^+ - 1_{(i+1)} 1_j^T \otimes \delta_y \right) \quad (17) \end{aligned}$$

with $\mathcal{K}_p^+(p, q) = \sum_{h=1}^{i+1} K(u_{h+p-2}, u_{h+q-2})$ for all $p, q = 1, \dots, j$, and

$$\begin{aligned} \mathcal{B}_p^- &= \mathcal{B}^-(1 : i+1, :) \\ \mathcal{S}_p^- &= \mathcal{S}^-(:, 1 : i+1). \end{aligned}$$

$\mathcal{A}^-, \mathcal{B}^-$ and L_y^- follow from:

$$\begin{bmatrix} 0 & 0 & 1^T & 0 \\ 0 & 0 & Y_p^+ & 0 \\ 1 & (Y_p^+)^T & \mathcal{K}_{pf} & \mathcal{S} \\ 0 & 0 & \mathcal{S}^T & \mathcal{T} \end{bmatrix} \begin{bmatrix} \bar{d}^- \\ (L_y^-)^T \\ \mathcal{A}^- \\ \mathcal{B}^- \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ (Y_f^-)^T \\ 0 \end{bmatrix}$$

with $\mathcal{K}_{pf} = \mathcal{K}_p + \mathcal{K}_f + \gamma^{-1} I$ and

$$\bar{d}^- = ((1_{(i-1)} \otimes I_l) - L_y^- (1_{(i+1)} \otimes I_l)) \delta_y.$$

C. Obtaining Estimates for the States

The state sequences \tilde{X}_i and \tilde{X}_{i+1} can now be determined from \mathcal{O}_i and \mathcal{O}_{i+1} in line with what is done in the linear case discussed in Section II. These state sequences will be used in a second step of the algorithm to obtain estimates for the system matrices and the nonlinearity f . Note that in the linear case, it is well known that the obtained state sequences \tilde{X}_i and \tilde{X}_{i+1} can be considered as the result of a bank of nonsteady-state Kalman filters working in parallel on the columns of the block-Hankel matrix W_p [27]. In the Hammerstein case, and if f were known, this relation would still hold provided that W_p is replaced by $[\Phi_f(U_p)^T Y_p^T]^T$. However, an estimate \hat{f} for f based on a finite amount of data will in general be subject to approximation errors [29]. As the classical results for the bank of linear Kalman filters are not applicable if the inputs $\hat{f}(u_t)$ to the linear model are not exact the obtained states \tilde{X}_i and \tilde{X}_{i+1} can no longer be seen as the result of a bank of Kalman filters working on $[\Phi_f(U_p)^T Y_p^T]^T$. Despite the loss of this property, it will be illustrated in the examples that the proposed method outperforms existing Hammerstein approaches such as approaches based on nonlinear autoregressive with exogenous inputs (NARX) models and N4SID identification algorithms with an expansion in Hermite polynomials.

D. Extraction of the System Matrices and the Static Nonlinearity F

The linear model and static nonlinearity are estimated from

$$\begin{aligned} & (\hat{A}, \hat{B}, \hat{C}, \hat{D}, \hat{f}) \\ = & \underset{A, B, C, D, f}{\operatorname{argmin}} \left\| \begin{bmatrix} \tilde{X}_{i+1} \\ Y_{i|i} - \delta_y \end{bmatrix} - \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} \tilde{X}_i \\ \Phi_f(U_{i|i}) \end{bmatrix} \right\|_F^2. \quad (18) \end{aligned}$$

It will be shown in this subsection that this least-squares problem can again be written as an LS-SVM regression problem. Denoting

$$\mathcal{X}_{i+1} = \begin{bmatrix} \tilde{X}_{i+1} \\ Y_{i|i} - \delta_y \end{bmatrix} \quad \Theta_{AC} = \begin{bmatrix} A \\ C \end{bmatrix} \quad \Theta_{BD} = \begin{bmatrix} B \\ D \end{bmatrix} \quad (19)$$

and replacing $\Theta_{BD}(s, :)$ by $\omega_s^T \varphi$, where again an expansion of a product of scalars and nonlinear functions is written as a linear combination of nonlinear functions, we have

$$\mathcal{X}_{i+1} = \Theta_{AC} \tilde{X}_i + \begin{bmatrix} \omega_1^T \\ \omega_2^T \\ \vdots \\ \omega_{n+t}^T \end{bmatrix} \Phi_\varphi(U_{i|i}) + E$$

with E the residuals of (18). The resulting LS-SVM primal problem can be written as

$$\min_{\omega_s, E, \Theta_{AC}} \mathcal{J}(\omega, E) = \frac{1}{2} \sum_{s=1}^{n+l} \omega_s^T \omega_s + \frac{\gamma_{BD}}{2} \sum_{s=1}^{n+l} \sum_{t=1}^j E(s, t)^2$$

$$\text{s.t.} \begin{cases} \mathcal{X}_{i+1}(s, t) = \Theta_{AC}(s, :) \tilde{X}_i(:, t) + \omega_s^T \varphi(u_{i+t-1}) & (a) \\ \forall s = 1, \dots, li, \quad t = 1, \dots, j \\ \sum_{t=0}^{N-1} \omega_s^T \varphi(u_t) = 0 \quad \forall s = 1, \dots, li & (b) \end{cases}$$

where γ_{BD} denotes a regularization constant which can be different from the γ used in Subsection III-A.

Lemma 3.2: Estimates for A and C in Θ_{AC} are obtained from the following dual problem:

$$\left[\begin{array}{c|c|c} 0 & \tilde{X}_i & 0 \\ \tilde{X}_i^T & \mathcal{K}_{BD} + \gamma_{BD}^{-1} I & \mathcal{S}_{BD} \\ 0 & \mathcal{S}_{BD}^T & \mathcal{T}_{BD} \end{array} \right] \begin{bmatrix} \Theta_{AC}^T \\ \mathcal{A}_{BD} \\ \mathcal{B}_{BD} \end{bmatrix} = \begin{bmatrix} 0 \\ \mathcal{X}_{i+1}^T \\ 0 \end{bmatrix} \quad (20)$$

whereby $\omega_s = \sum_{t=1}^j \alpha_{s,t} \varphi(u_{i+t-1}) + \beta_s \sum_{t=0}^{N-1} \varphi(u_t)$ for all $s = 1, \dots, n+l$, $\mathcal{K}_{BD}(p, q) = K(u_{i+p-1}, u_{i+q-1})$ for all $p, q = 1, \dots, j$, $\mathcal{B}_{BD} = [\beta_1 \ \beta_2 \ \dots \ \beta_{n+l}]$, $\mathcal{T}_{BD} = 1_N^T K 1_N$, and

$$\mathcal{A}_{BD} = \begin{bmatrix} \alpha_{1,1} & \alpha_{2,1} & \dots & \alpha_{n+l,1} \\ \alpha_{1,2} & \alpha_{2,2} & \dots & \alpha_{n+l,2} \\ \vdots & \vdots & & \vdots \\ \alpha_{1,j} & \alpha_{2,j} & \dots & \alpha_{n+l,j} \end{bmatrix} \quad \mathcal{S}_{BD} = \begin{bmatrix} \mathcal{S}_{i+1} \\ \mathcal{S}_{i+2} \\ \vdots \\ \mathcal{S}_{i+j} \end{bmatrix}.$$

Proof: This follows directly from the Lagrangian:

$$\mathcal{L} = \mathcal{J}(\omega, E) - \sum_{s=1}^{n+l} \beta_s \left\{ \sum_{t=0}^{N-1} \omega_s \varphi(u_t) \right\} - \sum_{s=1}^{n+l} \alpha_{s,t} \left\{ \mathcal{X}_{i+1}(s, t) - \Theta_{AC}(s, :) \tilde{X}_i(:, t) - \omega_s^T \varphi(u_{i+t-1}) \right\}$$

by taking the conditions for optimality $(\partial \mathcal{L})/(\partial \omega_s) = 0$, $(\partial \mathcal{L})/(\partial E) = 0$, $(\partial \mathcal{L})/(\partial \Theta_{AC}) = 0$, $(\partial \mathcal{L})/(\partial \alpha_{s,t}) = 0$, $(\partial \mathcal{L})/(\partial \beta_s) = 0$, and after elimination of the primal variables ω_s and E . ■

By combining the results from Lemma 3.2 with (18) and (19), we have

$$\Theta_{BD} [f(u_0) \ f(u_1) \ \dots \ f(u_{N-1})] = \mathcal{A}_{BD}^T \Omega(i+1 : i+j, :) + \mathcal{B}_{BD}^T \sum_{t=1}^N \Omega(t, :). \quad (21)$$

Hence, estimates for B, D in Θ_{BD} and the nonlinearity f can be obtained from a rank m approximation of the right hand side of (21), for instance using a singular value decomposition. This is a typical step in overparameterization approaches [1] and amounts to projecting the results for the overparameterized model as used in the estimation onto the class of Hammerstein models.

E. Practical Implementation

Following the discussion in the previous sections, the final algorithm for Hammerstein N4SID subspace identification can be summarized as follows.

- 1) Find estimates for the oblique projections \mathcal{O}_i and \mathcal{O}_{i+1} from (16) and (17).
- 2) Find estimates for the state following the procedure outlined in Subsection III-C.
- 3) Obtain estimates for A, C, \mathcal{A}_{BD} and \mathcal{B}_{BD} following the procedure outlined in Subsection III-D.
- 4) Obtain estimates for B, D and f from a rank- m approximation of (21).

It should be noted at this point that given the fact that regularization is inherently present in the proposed identification technique, lack of persistency of excitation will not lead to any numerical problems. However, in order to ensure that all aspects of the linear system are properly identified, persistency of excitation of $f(u)$ of at least order $2i$ is desired (see also Section II). Persistency of excitation of $f(u)$ can for some nonlinear functions f be expressed as a condition on the original inputs u but the relation is certainly not always straightforward (see for instance [30] for a discussion on this issue).

Furthermore, it is important to remark that the estimate of the static nonlinearity will only be reliable in regions where the input density is sufficiently high.

IV. ILLUSTRATIVE EXAMPLES

In this section, the presented algorithm is compared to the Hammerstein ARX approach presented in [8] and classical subspace Hammerstein identification algorithms involving overparameterization in orthogonal basis functions. Two properties of the Hammerstein N4SID subspace approach will thereby be highlighted.

- The greater flexibility that comes with the use of state-space models over more classical Hammerstein ARX approaches.
- The superior performance of the introduced algorithm over existing overparameterization approaches for Hammerstein subspace identification.

A. Comparison Between Hammerstein N4SID and Hammerstein ARX

Consider the following system which belongs to the Hammerstein class of models:

$$A(z)(y + \nu) = B(z)f(u) + e \quad (22)$$

with A and B polynomials in the forward shift operator z where $B(z) = z^6 + 0.8z^5 + 0.3z^4 + 0.4z^3$, $A(z) = (z - 0.98e^{\pm i})(z - 0.98e^{\pm 1.6i})(z - 0.97e^{\pm 0.4i})$. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ such that $f(u) = \text{sinc}(u)u^2$ be the static nonlinearity. A dataset was generated from this system where $u_t \sim \mathcal{N}(0, 2)$ is a white Gaussian noise sequence for $t = 0, \dots, N-1$ with $N = 1000$ and $\{e_t\}_{t=0}^{N-1}$ is a sequence of Gaussian white noise with a level of 10% of the level of the nonlinearity $f(u)$.

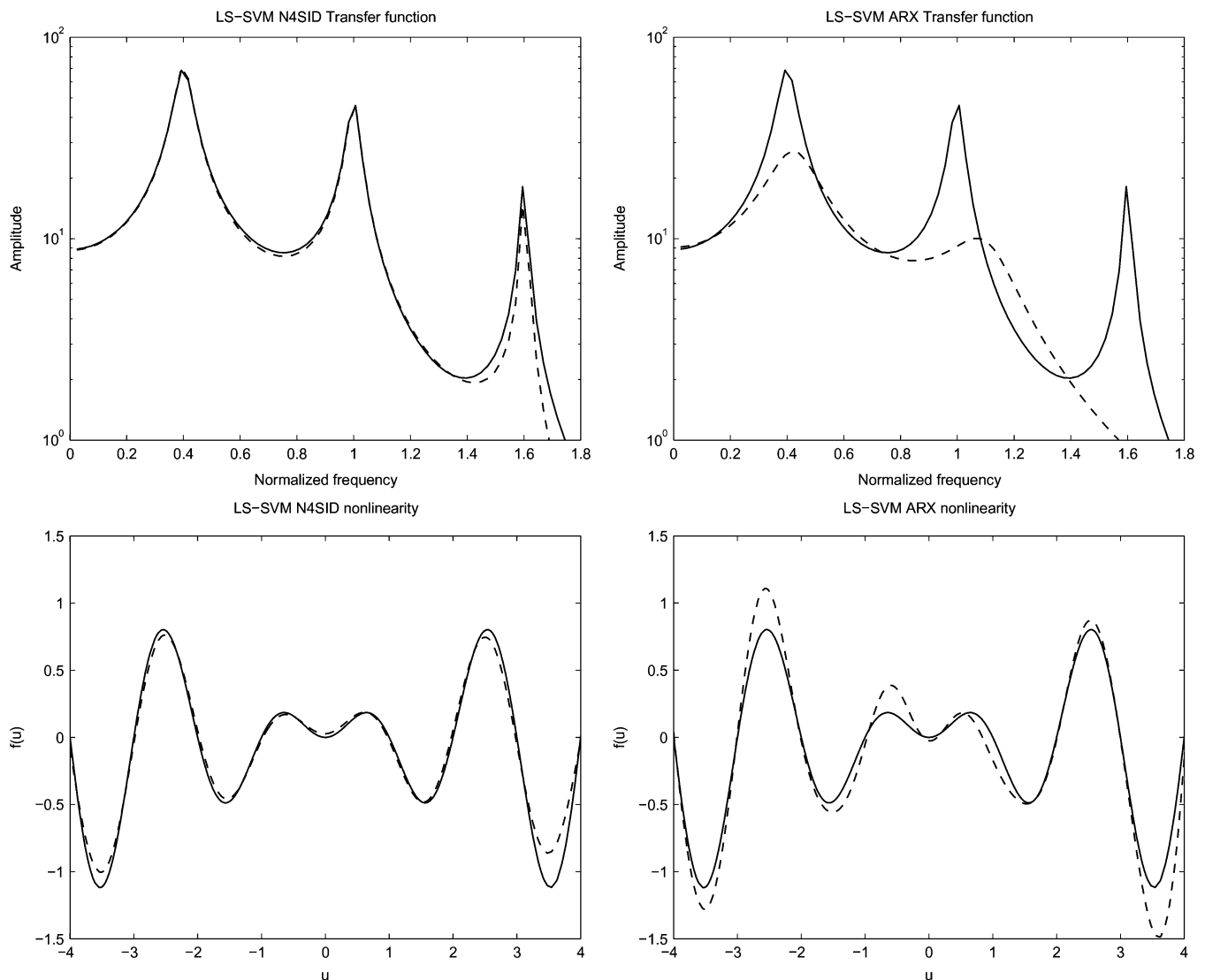


Fig. 1. True transfer function (solid) and estimated ones (dashed) for the LS-SVM N4SID subspace algorithm (top left) and the LS-SVM ARX algorithm (top right), as estimated from a sequence of 1000 input/output measurements on a simulated system, with the addition of 10% output noise. The true nonlinearities (solid) and estimated ones (dashed) are displayed below the transfer functions, for the N4SID case (lower left), and the ARX-case (lower right).

The measurement noise terms ν_t were chosen to be zero mean Gaussian white noise such that a signal to noise ratio of 10 was obtained at the output signal. The Hammerstein N4SID subspace identification algorithm as derived in Section III was used to extract the linear model and the static nonlinearity f from the dataset described above. The number of block-rows i in the Block Hankel matrices was set to 10 which is a common choice in subspace identification algorithms [27]. An advantage of the N4SID algorithm is that the model order, 6 in this case, follows automatically from the spectrum of the SVD. The hyper-parameters in the LS-SVM N4SID algorithm were selected as $\sigma = 0.1, \gamma = 1000, \gamma_{BD} = 10$ by validation on an independent validation set. The resulting linear system and static nonlinearity are displayed in Fig. 1.

As a comparison, the results of the LS-SVM ARX estimator [8] are also displayed in Fig. 1. For the ARX-estimator, the number of poles and zeros were assumed to be fixed a priori. Two hyper-parameters (the regularization constant and the bandwidth of the RBF kernel) which need to be set in this

method were chosen in accordance with the choices reported in [8]. Note that although the Hammerstein ARX method performed very well for this example in the absence of output noise (see the examples in [8]), its performance deteriorates in the presence of output noise as evidenced by the poor fit in Fig. 1.

This highlights one of the main advantages of the use of subspace identification methods [27] over more classical ARX procedures, namely that they allow for the successful estimation of a much wider class of linear systems. Note on the other hand that if the true system fits well into the more restricted ARX framework, use of the latter is to be preferred [8].

B. Comparison With Classical Subspace Overparameterization Approaches

As mentioned before, a classical approach to Hammerstein system identification is to expand the static nonlinearity in a set of orthogonal or nonorthogonal basis-functions [17]. The same idea can be applied to subspace algorithms [15]. Once a set

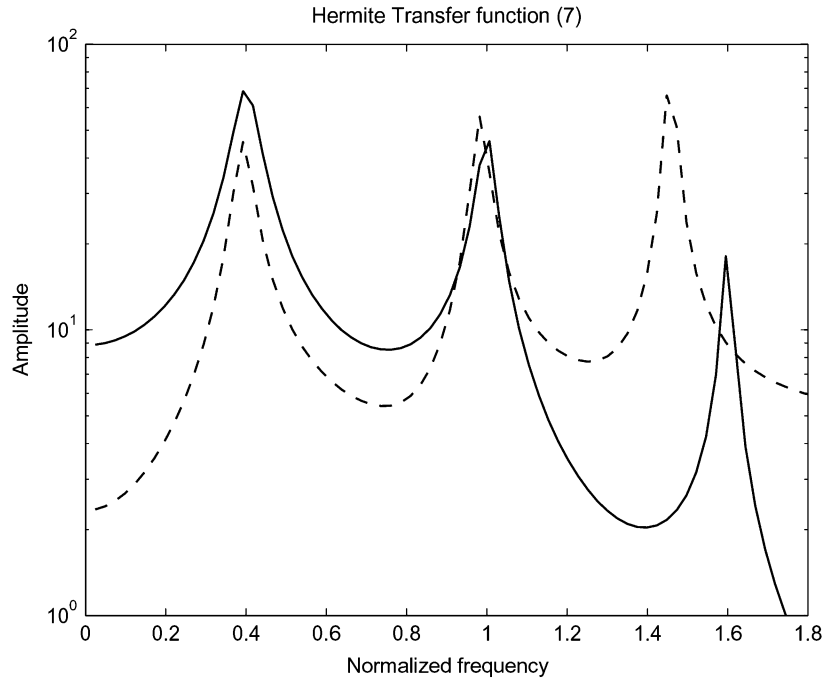


Fig. 2. True transfer function (solid) and the estimated one (dashed) for the Hermite N4SID subspace algorithm as estimated from a sequence of 1000 input/output measurements on a simulated system, with the addition of 10% output noise.

of basis-functions is considered, the one-dimensional input is transformed into a higher-dimensional input vector which contains the coefficients of the expansion of $f(u)$ in its basis. The classical N4SID subspace algorithm as outlined in Section II is thereafter applied. The linear system and static nonlinearities can be obtained from the obtained matrices B and D (see [15] for a detailed procedure).

This example will adopt the common choice of the Hermite polynomials as a basis. The best results on the dataset with output noise were obtained when selecting seven Hermite polynomials with orders ranging from 0 to 6. The obtained linear system corresponding to this choice of basis functions is displayed in Fig. 2. Note the rather poor performance of this method, compared to the LS-SVM N4SID algorithm. This can largely be attributed to the fact that the performance of subspace algorithms degrades as the number of inputs increases, certainly if these inputs are highly correlated [4]. This as a result of a bad conditioning of the matrices U_p and U_f as the number of rows increases and these rows get more correlated. For the zero-order Hermite polynomial (which is a constant) this is certainly the case but also when leaving out this polynomial, condition numbers of 10^5 and higher are encountered. This problem does not occur in the N4SID LS-SVM algorithm as the latter features an inherently available regularization framework. An additional advantage is the flexibility one gets by plugging in an appropriate kernel and the fact that if localized kernels are used, no specific choices have to be made for their locations. The locations follow directly from the formulation of costfunctions as (14).

V. CONCLUSION

In this paper, a method for the identification of Hammerstein systems was presented based on the well-known N4SID

subspace identification algorithm. The basic framework of the N4SID algorithm is largely left untouched, except for the ordinary least squares steps which are replaced by a set of componentwise LS-SVM regressions. The proposed algorithm was observed to be able to extract the linear system and the static nonlinearity from data, even in the presence of output noise.

APPENDIX I LS-SVM FUNCTION ESTIMATION

Let $\{(x_i, y_i)\}_{i=0}^N \subset \mathbb{R}^d \times \mathbb{R}$ be a set of independently and identically distributed (i.i.d.) input/output training data with input samples x_i and output samples y_i . Consider the static regression model $y_i = f(x_i) + e_i$ where where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is an unknown real-valued smooth function and $\{e_i\}_{i=1}^N$ are i.i.d. (uncorrelated) random errors with $E[e_i] = 0, E[e_i^2] = \sigma_e^2 < \infty$. Originating from the research on classification algorithms, Support Vector Machines (SVM's) and other kernel methods have been used for the purpose of estimating the nonlinear f . The following model is assumed:

$$f(x) = w^T \varphi(x) + b$$

where $\varphi(x) : \mathbb{R}^d \rightarrow \mathbb{R}^{n_H}$ denotes a potentially infinite ($n_H = \infty$) dimensional feature map which doesn't have to be known explicitly. In the following paragraph we will see how the feature map can be induced in an implicit way by the adoption of a proper kernel function. The regularized cost function of the LS-SVM is given as

$$\begin{aligned} \min_{w, b, e} \mathcal{J}(w, e) &= \frac{1}{2} w^T w + \frac{\gamma}{2} \sum_{i=1}^n e_i^2 \\ \text{s.t. } y_i &= w^T \varphi(x_i) + b + e_i, \quad i = 1, \dots, N. \end{aligned}$$

The relative importance between the smoothness of the solution and the data fitting is governed by the scalar $\gamma \in \mathbb{R}_0^+$ referred to as the regularization constant. The optimization performed corresponds to ridge regression (see, e.g., [10]) in feature space. In order to solve the constrained optimization problem, the Lagrangian $\mathcal{L}(w, b, e; \alpha) = \mathcal{J}(w, e) - \sum_{i=1}^N \alpha_i \{w^T \varphi(x_i) + b + e_i - y_i\}$ is constructed, with $\alpha_i \in \mathbb{R}$ the Lagrange multipliers. After application of the conditions for optimality: $(\partial \mathcal{L})/(\partial w) = 0, (\partial \mathcal{L})/(\partial b) = 0, (\partial \mathcal{L})/(\partial e_i) = 0, (\partial \mathcal{L})/(\partial \alpha_i) = 0$, the following set of linear equations is obtained:

$$\begin{bmatrix} 0 & \mathbf{1}_N^T \\ \mathbf{1}_N & \Omega + \gamma^{-1} I_N \end{bmatrix} \begin{bmatrix} b \\ \hat{\alpha} \end{bmatrix} = \begin{bmatrix} 0 \\ y \end{bmatrix} \quad (23)$$

where $y = [y_1 \dots y_N]^T, \mathbf{1}_N = [1 \dots 1]^T, \alpha = [\alpha_1 \dots \alpha_N]^T, \Omega_{i,j} = K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j), \forall i, j = 1, \dots, N$, with K a positive-definite Mercer kernel function. Note that in order to solve the set of (23), the feature map φ does never have to be defined explicitly. Only its inner product, a positive definite kernel, is needed. This is called the kernel trick [29], [19]. For the choice of the kernel $K(\cdot, \cdot)$, see, e.g., [19]. Typical examples are the use of a linear kernel $K(x_i, x_j) = x_i^T x_j$, a polynomial kernel $K(x_i, x_j) = (\tau + x_i^T x_j)^d$ of degree d or the RBF kernel $K(x_i, x_j) = \exp(-|x_i - x_j|_2^2 / \sigma^2)$ where σ denotes the bandwidth of the kernel. The resulting LS-SVM model can be evaluated at a new point $x_* \in \mathbb{R}^d$ as

$$\hat{f}(x_*) = \sum_{i=1}^N \hat{\alpha}_i K(x_*, x_i) + \hat{b}$$

where $(\hat{b}, \hat{\alpha})$ is the solution to (23).

LS-SVMs are reformulations to the original SVM's employed for tasks in classification [24], regression [23] and provides primal-dual optimization formulations to the algorithm of kernel principal component analysis (KPCA), kernel partial least squares (KPLS), kernel canonical correlation analysis (KCCA), and others [23]. By the use of the least squares criterion and the use of equality instead of inequality constraints, the estimation typically boils down to the solution of a set of linear equations or eigenvalue problems instead of the optimization of quadratic programming problems [22]–[24]. In [8], the task of identifying a Hammerstein model using an LS-SVM based approach of the nonlinearity and an ARX system was considered.

APPENDIX II

HAMMERSTEIN FIR MODELS USING LS-SVMS

The extension of LS-SVMs toward the estimation of additive models was studied in [18]. It was applied toward the identification of Hammerstein ARX models in [8]. We review briefly the basic steps as they will reoccur in the presented technique. Let $\{(u_t, y_t)\}_{t=0}^{N-1} \subset \mathbb{R} \times \mathbb{R}$ be a (SISO) sequence of observations. Consider a Hammerstein FIR model of order $L \in \mathbb{N}_0$.

$$y_t = \sum_{l=0}^{L-1} c_l f(u_{t-l}) + e_t = \sum_{l=0}^{L-1} c_l (w_l^T \varphi(u_{t-l}) + b) + e_t \quad (24)$$

for all $t = L - 1, \dots, N - 1$ where $c = (c_1, \dots, c_L)^T \in \mathbb{R}^L$ is a vector. Whenever both w, b as well as c are unknown, the simultaneous estimation of those parameters is known to be a hard problem. Following [1], an overparameterization technique can be adopted. Consider in the first stage the identification of the parameters w_l, b_l and c of the slightly broader model

$$y_t = \sum_{l=0}^{L-1} f_l(u_{t-l}) + e_t = \sum_{l=0}^{L-1} (w_l^T \varphi(u_{t-l}) + b_l) + e_t \quad (25)$$

for all $t = L - 1, \dots, N - 1$ where $w_l \in \mathbb{R}^{N_h}$ and $b_l \in \mathbb{R}$ for all $l = 0, \dots, L - 1$. A necessary and sufficient condition for the restriction of (25) to the Hammerstein class (24) can be written as the rank constraint

$$[f_0(u) \dots f_{L-1}(u)] = f(u)c^T \quad \forall u \in \mathbb{R}. \quad (26)$$

It becomes clear that the right hand side occurring in the term (25) has a nonunique representation as one can always add (and subtract) a row-vector $\delta \in \mathbb{R}^{1 \times L}$ to the nonlinear function $f_L: \mathbb{R} \rightarrow \mathbb{R}^{1 \times L}$ such that $f_L(u) = [f_0(u) \dots f_{L-1}(u)] \in \mathbb{R}^{1 \times L}$ and $\sum_{l=1}^L \delta(l) = 0$. This follows from the following relation:

$$(f_L(u) + \delta)1_L = f_L(u)1_L \quad \forall u \in \mathbb{R}.$$

However, this operation does not preserve the constraint (26) if $\delta \neq 0_L$. As a bias term b can be found such that $y_t = f(u_{t-L}, \dots, u_t) + b$ and $E[f_l(u)] = 0$, the nonlinear functions can be centered around zero without loss of generality. Then a necessary linear condition for (26) becomes

$$E[f_0(u) \dots f_{L-1}(u)] = E[f(u)]c^T = 0_L^T \quad (27)$$

or using the empirical counterpart

$$\sum_{t=0}^{N-1} f_l(u_t) = \sum_{t=0}^{N-1} w_l^T \varphi(u_t) = 0 \quad \forall l = 0, \dots, L - 1 \quad (28)$$

which are referred to as the centering constraint.

The overparameterization procedure amounts to first obtaining estimates of the model class (25) subject to the centering constraints (28) and afterwards projecting the result onto the Hammerstein class by calculating a rank one approximation of the estimate using an SVD. The primal-dual derivation can be summarized as follows [8].

Lemma 2.1: Consider the primal estimation problem

$$\begin{aligned} \underset{w, e, b}{\operatorname{argmin}} &= \frac{\gamma}{2} \sum_{t=L-1}^{N-1} e_t^2 + \frac{1}{2} \sum_{l=0}^{L-1} w_l^T w_l \\ \text{s.t.} & \begin{cases} \sum_{l=0}^{L-1} w_l^T \varphi(u_{t-l}) + b + e_t - y_t \\ \forall t = L - 1, \dots, N - 1 \\ \sum_{t=0}^{N-1} w_l^T \varphi(u_t) = 0 \quad \forall l = 0, \dots, L - 1. \end{cases} \quad (29) \end{aligned}$$

Let $\Omega_l \in \mathbb{R}^{N-L \times N-L}$ be a positive-definite matrix defined as $\Omega_{l,i,j} = \varphi(u_{i-l})^T \varphi(u_{j-l})$ for all $i, j = L - 1, \dots, N - 1$ and $l = 0, \dots, L - 1$. Let $\Omega = \sum_{l=0}^{L-1} \Omega_l$ be the kernel matrix and let $\mathcal{S} \in \mathbb{R}^{(N-L) \times L}$ such that $\mathcal{S}(:, l) = \Omega_{(l-1)} 1_{N-L}$ for

all $l = 0, \dots, L-1$. The solution is uniquely characterized by the following dual problem:

$$\begin{bmatrix} 0 & 0_L^T & 1_{N-L}^T \\ 0_L & 0_{L \times L} & \mathcal{S}^T \\ 1_{N-L} & \mathcal{S} & \Omega + \gamma^{-1} I_{N-L} \end{bmatrix} \begin{bmatrix} b \\ \beta \\ \alpha \end{bmatrix} = \begin{bmatrix} 0 \\ 0_L \\ y \end{bmatrix} \quad (30)$$

where $\alpha \in \mathbb{R}^{N-L}$ and $\beta = (\beta_1, \dots, \beta_L)^T \in \mathbb{R}^L$ denote the Lagrange multipliers to the constraints in (29). The estimate can be evaluated in a new datapoint $u \in \mathbb{R}$ as

$$\hat{f}_l(u) = \sum_{t=L-1}^{N-1} \hat{\alpha}_t K(u_t, u) + \hat{\beta}_l \sum_{t=L-1}^{N-1} K(u_t, u)$$

and $\hat{f}(u_t) = \sum_{l=0}^{L-1} \hat{f}_l(u_t - l + 1) + \hat{b}$ where $\hat{\alpha}_t, \hat{\beta}$ and \hat{b} are the solution to (30).

REFERENCES

- [1] E. W. Bai, "An optimal two-stage identification algorithm for Hammerstein-Wiener nonlinear systems," *Automatica*, vol. 4, no. 3, pp. 333–338, 1998.
- [2] —, "A blind approach to Hammerstein model identification," *IEEE Trans. Signal Process.*, vol. 50, no. 7, pp. 1610–1619, Jul. 2002.
- [3] F. H. I. Chang and R. Luus, "A noniterative method for identification using the Hammerstein model," *IEEE Trans. Autom. Control*, vol. AC-16, no. 5, pp. 464–468, Oct. 1971.
- [4] A. Chiuso and G. Picci, "On the ill-conditioning of subspace identification with inputs," *Automatica*, vol. 40, no. 4, pp. 575–589, 2004.
- [5] P. Crama and J. Schoukens, "Initial estimates of Wiener and Hammerstein systems using multisine excitation," *IEEE Trans. Measure. Instrum.*, vol. 50, no. 6, pp. 1791–1795, Jun. 2001.
- [6] E. J. Dempsey and D. T. Westwick, "Identification of Hammerstein models with cubic spline nonlinearities," *IEEE Trans. Biomed. Eng.*, vol. 51, no. 2, pp. 237–245, Feb. 2004.
- [7] E. Eskinat, S. H. Johnson, and W. L. Luyben, "Use of Hammerstein models in identification of nonlinear systems," *AIChE J.*, vol. 37, no. 2, pp. 255–268, 1991.
- [8] I. Goethals, K. Pelckmans, J. A. K. Suykens, and B. De Moor, "Identification of MIMO Hammerstein models using least squares support vector machines," ESAT-SISTA, Leuven, Belgium, Tech. Rep. 04-45, 2004.
- [9] —, "NARX identification of Hammerstein models using least squares support vector machines," in *Proc. 6th IFAC Symp. Nonlinear Control Systems (NOLCOS 2004)*, Stuttgart, Germany, Sep. 2004, pp. 507–512.
- [10] G. H. Golub and C. F. Van Loan, *Matrix Computations*. Baltimore, MD: John Hopkins Univ. Press, 1989.
- [11] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Heidelberg, Germany: Springer-Verlag, 2001.
- [12] A. Janczak, "Neural network approach for identification of Hammerstein systems," *Int. J. Control*, vol. 76, no. 17, pp. 1749–1766, 2003.
- [13] M. J. Korenberg, "Recent advances in the identification of nonlinear systems: Minimum-variance approximation by hammerstein models," in *Proc. IEEE EMBS*, vol. 13, 1991, pp. 2258–2259.
- [14] Z. H. Lang, "Controller design oriented model identification method for Hammerstein systems," *Automatica*, vol. 29, pp. 767–771, 1993.
- [15] T. McKelvey and C. Hanner, "On identification of Hammerstein systems using excitation with a finite number of levels," in *Proc. 13th Int. Symp. System Identification (SYSID2003)*, 2003, pp. 57–60.
- [16] K. S. Narendra and P. G. Gallman, "An iterative method for the identification of nonlinear systems using the Hammerstein model," *IEEE Trans. Autom. Control*, vol. AC-11, no. 3, pp. 546–550, Jul. 1966.
- [17] M. Pawlak, "On the series expansion approach to the identification of Hammerstein systems," *IEEE Trans. Autom. Control*, vol. 36, no. 6, pp. 736–767, Jun. 1991.
- [18] K. Pelckmans, I. Goethals, J. De Brabanter, J. A. K. Suykens, and B. De Moor, "Componentwise least squares support vector machines," in *Support Vector Machines: Theory and Applications*, L. Wang, Ed. New York: Springer-Verlag, 2005.
- [19] B. Schölkopf and A. Smola, *Learning with Kernels*. Cambridge, MA: MIT Press, 2002.
- [20] J. Sjöberg, Q. Zhang, L. Ljung, A. Benveniste, B. Delyon, P. Glorennec, H. Hjalmarrsson, and A. Juditsky, "Nonlinear black-box modeling in system identification: A unified overview," *Automatica*, vol. 31, no. 12, pp. 1691–1724, 1995.
- [21] J. C. Stapleton and S. C. Bass, "Adaptive noise cancellation for a class of nonlinear dynamic reference channels," *IEEE Trans. Circuits Syst.*, vol. CAS-32, no. 2, pp. 143–150, Feb. 1985.
- [22] J. A. K. Suykens, G. Horvath, S. Basu, C. Micchelli, and J. Vandewalle, Eds., *Advances in Learning Theory: Methods, Models and Applications, Volume 90 of NATO Science Series III: Computer & Systems Sciences*. Amsterdam, The Netherlands: IOS, 2003.
- [23] J. A. K. Suykens, T. Van Gestel, J. De Brabanter, B. De Moor, and J. Vandewalle, *Least Squares Support Vector Machines*. Singapore: World Scientific, 2002.
- [24] J. A. K. Suykens and J. Vandewalle, "Least squares support vector machine classifiers," *Neural Process. Lett.*, vol. 9, pp. 293–300, 1999.
- [25] P. Van Overschee and B. De Moor, "Subspace algorithms for the stochastic identification problem," *Automatica*, vol. 29, no. 3, pp. 649–660, 1993.
- [26] —, "A unifying theorem for three subspace system identification algorithms," *Automatica*, vol. 31, no. 12, pp. 1853–1864, 1995.
- [27] P. Van Overschee and B. De Moor, *Subspace Identification for Linear Systems: Theory, Implementation, Applications*. Norwell, MA: Kluwer, 1996.
- [28] T. H. van Pelt and D. S. Bernstein, "Nonlinear system identification using Hammerstein and nonlinear feedback models with piecewise linear static maps—Part I: Theory," in *Proc. Amer. Control Conf. (ACC2000)*, 2000, pp. 225–229.
- [29] V. N. Vapnik, *Statistical Learning Theory*. New York: Wiley, 1998.
- [30] M. Verhaegen and D. Westwick, "Identifying MIMO Hammerstein systems in the context of subspace model identification methods," *Int. J. Control*, vol. 63, pp. 331–349, 1996.
- [31] G. Wahba, *Spline Models for Observational Data*. Philadelphia, PA: SIAM, 1990.
- [32] D. Westwick and R. Kearney, "Identification of a hammerstein model of the stretch reflex EMG using separable least squares," in *Proc. World Congr. Medical Physics and Biomedical Engineering*, Chicago, IL, 2000.



Ivan Goethals was born in Wilrijk, Belgium, in 1978. He received the M.Sc. degree in nuclear physics from the K. U. Leuven, Leuven, Belgium, in 2000, and the Ph.D. degree for his research with the SCD-SISTA Group of the Department of Electrical Engineering (ESAT) at the same university in 2005.

His main research interests are in the fields of linear and nonlinear system identification.



Kristiaan Pelckmans was born on November 3, 1978, in Merksplas, Belgium. He received the M.Sc. degree in computer science and the Ph.D. degree from K. U. Leuven, Leuven, Belgium, in 2000 and 2005, respectively.

After a projectwork for an implementation of kernel machines and LS-SVMs (LS-SVMlab), he was a researcher at the K. U. Leuven in the Department of Electrical Engineering in the SCD SISTA Laboratory. His research mainly focuses on machine learning and statistical inference using primal-dual

kernel machines.



Johan A. K. Suykens was born in Willebroek, Belgium, on May 18 1966. He received the degree in electro-mechanical engineering and the Ph.D. degree in applied sciences from K. U. Leuven, Leuven, Belgium, in 1989 and 1995, respectively.

In 1996, he was a Visiting Postdoctoral Researcher at the University of California, Berkeley. He has been a Postdoctoral Researcher with the Fund for Scientific Research FWO Flanders and is currently an Associate Professor with K. U. Leuven. His research interests are mainly in the areas of the theory and appli-

cation of neural networks and nonlinear systems. He is author of the books *Artificial Neural Networks for Modeling and Control of Non-linear Systems* (Norwell, MA: Kluwer, 1995) and *Least Squares Support Vector Machines* (Singapore: World Scientific, 2002) and editor of the books *Nonlinear Modeling: Advanced Black-Box Techniques* (Norwell, MA: Kluwer, 1998) and *Advances in Learning Theory: Methods, Models and Applications* (Amsterdam, The Netherlands: IOS, 2003). In 1998, he organized an International Workshop on Nonlinear Modeling with Time-series Prediction Competition.

Dr. Suykens has served as Associate Editor for the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS—PART I (1997–1999) and PART II (since 2004), and since 1998, he has been an Associate Editor for the IEEE TRANSACTIONS ON NEURAL NETWORKS. He received an IEEE Signal Processing Society 1999 Best Paper (Senior) Award and several Best Paper Awards at International Conferences. He is a recipient of the International Neural Networks Society INNS 2000 Young Investigator Award for significant contributions in the field of neural networks. He has served as Director and Organizer of a NATO Advanced Study Institute on Learning Theory and Practice (Leuven 2002) and as a program co-chair for the International Joint Conference on Neural Networks IJCNN 2004.



Bart De Moor received the M.S. degree and a Ph.D. in electrical engineering at the K. U. Leuven, Leuven, Belgium, in 1983 and 1988, respectively.

He was a Visiting Research Associate at Stanford University, Stanford, CA (1988–1990). Currently, he is a Full Professor at the Department of Electrical Engineering of the K. U. Leuven. His research interests are in numerical linear algebra and optimization, system theory, control and identification, quantum information theory, data-mining, information retrieval, and bio-informatics, in which he (co-)authored more

than 400 papers and three books.

Dr. De Moor's work has won him several scientific awards [Leybold-Heraeus Prize (1986), Leslie Fox Prize (1989), Guillemin-Cauer best paper Award of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS (1990), Laureate of the Belgian Royal Academy of Sciences (1992), bi-annual Siemens Award (1994), best paper award of *Automatica* (IFAC, 1996), IEEE Signal Processing Society Best Paper Award (1999)]. From 1991 to 1999, he was the Chief Advisor on Science and Technology of several ministers of the Belgian Federal and the Flanders Regional Governments. He is on the board of three spin-off companies, of the Flemish Interuniversity Institute for Biotechnology, the Study Center for Nuclear Energy, and several other scientific and cultural organizations. Since 2002, he also makes regular television appearances in the Science Show "Hoe?Zo!" on national television in Belgium. Full biographical details can be found at www.esat.kuleuven.be/~demoor.