# AN ORTHOGONAL FORWARD REGRESSION ALGORITHM COMBINED WITH BASIS PURSUIT AND D-OPTIMALITY

*X. Hong[†], M. Brown[‡], S. Chen[§], C. J. Harris[§]*

[†]Department of Cybernetics
University of Reading, Reading, RG6 6AY, UK
[‡] Department of Computing and Mathematics
Manchester Metropolitan University, Manchester, UK
[§] Department of Electronics and Computer Science
University of Southampton, Southampton SO17 1BJ, UK

x.hong@reading.ac.uk

## ABSTRACT

A new forward regression model identification algorithm is introduced. The derived model parameters, in each forward regression step, are initially estimated via orthogonal least squares (OLS) (using the modified Gram-Schmidt procedure), followed by being tuned with a new gradient descent learning algorithm based on the basis pursuit that minimizes the $l^1$ norm of the parameter estimate vector. The model subset selection cost function includes a D-optimality design criterion. Both the parameter tuning procedure, based on basis pursuit, and the model selection criterion, based on the D-optimality that is effective in ensuring model robustness, are integrated with the forward regression, so as to maintain computational efficiency. An illustrative example is included to demonstrate the effectiveness of the new approach.

## 1. INTRODUCTION

A main obstacle in non-linear modelling using associative memory networks or fuzzy logic has been the problem of *the curse of dimensionality* [1]. This factor applies to all lattice based networks or knowledge representations [2, 3, 4, 5]. For these systems it is essential to use some model construction procedures to overcome the obstacle by deriving a model with an appropriate dimension. An orthogonal least squares (OLS) algorithm including parameter regularization technique based on Gram-Schmidt orthogonal decomposition can be used to determine the significant model elements and associated parameter estimates, and the overall model structure [6, 7, 8]. Model selection criteria such as the Akaike information criterion (AIC) [9] are usually incorporated into the procedure to determinate the model construction process. The use of AIC or other information based criteria, if used in forward regression, however only affects the stopping point of the model selection, but does not determine which regressor to be selected.

In recent studies, variants of OLS algorithm were introduced to improve model robustness via experimental design and parameter regularization [10, 11, 12, 13]. Alternatively the model sparsity can be achieved by a novel concept of the basis pursuit or least angle regression [14, 15] that aims to obtain a model by minimizing the $l^1$ norm of the parameters. Both parameter regularization and basis pursuit can be integrated into a Bayesian framework [12, 13, 16]. The advantage of the basis pursuit is that it can achieve much sparser models by forcing more parameters to zero, than models derived from the minimization of the $l^p$ norm, as most $l^p$ norms will produces parameters small, but nonzero, values. Compared to method of the regularization [7, 8], the basis pursuit method, however, will not generally be computationally efficient, because by simply changing from $l^2$ norm to $l^1$ norm in the cost function, this effectively changes a quadratic optimization problem with a simple solution into a more sophisticated problem for which a convex, nonquadratic optimization is generally required [14, 15].

In this paper, a new model identification technique is introduced by using forward regression with basis pursuit and D-optimality design. Based on the previous work [11], we incorporate the concept of basis pursuit to tune the parameter estimates as derived from the orthogonal least squares method. A gradient descent parameter learning method is initially introduced with proven convergence, followed by its application to the parameters tuning in the modified Gram-Schmidt algorithm. It is shown that parameter tuning by basis pursuit, following the initialization of least squares inherent in the Gram-Schmidt procedure will enforce model sparsity, yet fit well in the procedure automated by the D-optimality model selective criterion. The computational ef-

ficiency of the method due to the forward OLS regression maintains.

## 2. PRELIMINARIES

A linear regression model (RBF neural network, B-spline neurofuzzy network) can be formulated as [2, 3]

$$y(t) = \sum_{k=1}^{M} p_k(\mathbf{x}(t))\theta_k + \xi(t) \tag{1}$$

where $t = 1, 2, \cdots, N$, and $N$ is the size of the estimation data set. $y(t)$ is system output variable, $\mathbf{x}(t) = [y(t-1), \cdots, y(t-n_y), u(t-1), \cdots, u(t-n_u)]^T$ is system input vector with assumed known dimension of $(n_y + n_u)$. $u(t)$ is system input variable. $p_k(\bullet)$ is a known nonlinear basis function, such as RBF, or B-spline fuzzy membership functions. $\xi(t)$ is an uncorrelated model residual sequence with zero mean and variance of $\sigma^2$. Eq.(1) can be written in the matrix form as

$$\mathbf{y} = \mathbf{P}\Theta + \Xi \tag{2}$$

where $\mathbf{y} = [y(1), \cdots, y(N)]^T$ is the output vector. $\Theta = [\theta_1, \cdots, \theta_M]^T$ is parameter vector, $\Xi = [\xi(1), \cdots, \xi(N)]^T$ is the residual vector, and $\mathbf{P}$ is the regression matrix $\mathbf{P} = [\mathbf{p}_1, ... \mathbf{p}_M]$, where $\mathbf{p}_k = [p_k(1), \cdots, p_k(N)]^T$, with $p_k(t) = p_k(\mathbf{x}(t))$. An orthogonal decomposition of $\mathbf{P}$ is

$$\mathbf{P} = \mathbf{W}\mathbf{A} \tag{3}$$

where $\mathbf{A} = \{\alpha_{ij}\}$ is an $M \times M$ unit upper triangular matrix and $\mathbf{W}$ is an $N \times M$ matrix with orthogonal columns that satisfy

$$\mathbf{W}^T\mathbf{W} = diag\{\kappa_1, \cdots, \kappa_M\} \tag{4}$$

with

$$\kappa_k = \mathbf{w}_k^T \mathbf{w}_k, \qquad k = 1, \cdots, M \tag{5}$$

so that (2) can be expressed as

$$\mathbf{y} = (\mathbf{P}\mathbf{A}^{-1})(\mathbf{A}\Theta) + \Xi = \mathbf{W}\Gamma + \Xi \tag{6}$$

where $\Gamma = [\gamma_1, \cdots, \gamma_M]^T$ is an auxiliary vector.

### 2.1. The modified Gram-Schmidt algorithm and basis pursuit

Clearly for the orthogonalised system (6), the least squares estimates is given by

$$\gamma_k^{(0)} = \frac{\mathbf{w}_k^T \mathbf{y}}{\mathbf{w}_k^T \mathbf{w}_k}, \qquad k = 1, \cdots, M \tag{7}$$

The original model coefficient vector $\Theta = [\theta_1, \cdots, \theta_M]^T$ can then be calculated from $\mathbf{A}\Theta = \Gamma$ through back substitution.

The modified Gram-Schmidt procedure, described below, can be used to perform the orthogonalization of (3) and parameter estimation (7). Starting from $k = 1$, the columns $\mathbf{p}_j$, $k + 1 \leq j \leq M$ are made orthogonal to the $k$th column at the $k$th stage. The operation is repeated for $1 \leq k \leq M - 1$. Specifically, denoting $\mathbf{p}_j^{(0)} = \mathbf{p}_j$, $1 \leq j \leq M$, then for $k = 1, \cdots, M - 1$

$$\mathbf{w}_k = \mathbf{p}_k^{(k-1)}$$
$$\alpha_{kj} = \frac{\mathbf{w}_k^T \mathbf{p}_j^{(k-1)}}{\mathbf{w}_k^T \mathbf{w}_k}, \qquad k + 1 \leq j \leq M$$
$$\mathbf{p}_j^{(k)} = \mathbf{p}_j^{(k-1)} - \alpha_{kj}\mathbf{w}_k, \qquad k + 1 \leq j \leq M \tag{8}$$

where $\alpha_{kj}$'s are components of the upper triangular matrix $\mathbf{A}$. The last stage of the procedure is simply $\mathbf{w}_M = \mathbf{p}_M^{(M-1)}$. The elements of the auxiliary vector $\Gamma$ are computed by transforming $\mathbf{y}^{(0)} = \mathbf{y}$ in a similar way. For $1 \leq k \leq M$

$$\gamma_k^{(0)} = \frac{\mathbf{w}_k^T \mathbf{y}^{(k-1)}}{\mathbf{w}_k^T \mathbf{w}_k}$$
$$\mathbf{y}^{(k)} = \mathbf{y}^{(k-1)} - \gamma_k^{(0)}\mathbf{w}_k \tag{9}$$

It can be easily verified that $\gamma_k^{(0)}$ as derived from (9) is equivalent to (7). Geometrically the system output vector $\mathbf{y}$, at step $k$, is projected onto a set of orthogonal basis vectors, $\{\mathbf{w}_1, ... \mathbf{w}_k\}$. The model residual is decreased by projecting the system output vector $\mathbf{y}$ onto a new basis $\mathbf{w}_k$ at this step. Effectively, (9) can be regarded as a linear fitting of $\mathbf{y}^{(k-1)}$ by using a single variable $\mathbf{w}^{(k)}$, and to derive the new model residual $\mathbf{y}^{(k)}$, and so on. This observation will be explored further in Section 3.1 for the development of the proposed algorithm in Section 3.2. For better model parameter estimation bias/variance tradeoff, the regularization can be applied [7] with a solution from a quadratic form optimization, and the regularization parameters can be optimized by being treated as hyper-parameters in Bayesian approach [12]. Alternatively the basis pursuit method is simply given by changing the $l^2$ norm into $l^1$ such that

$$V = \frac{1}{2}E[\xi^2(t)] + \boldsymbol{\lambda}^T\|\Gamma\|_1 \tag{10}$$

is minimized for basis pursuit parameter estimates, where $\boldsymbol{\lambda} = [\lambda_1, ..., \lambda_{n_\theta}]^T$, $\|\Gamma\|_1 = [|\gamma_1|, ..., |\gamma_{n_\theta}|]^T$, and $n_\theta \leq M$ denotes the size of parameter vector of $\Gamma$ with nonzero parameters. $\lambda_k \geq 0$, are basis pursuit parameters. Note that only nonzero parameters that are actually included in the model are penalized, because a regressor with zero parameter does not influence model performance. Both parameter regularization and basis pursuit can be integrated into a Bayesian framework [12, 13, 16]. The basis pursuit method tends to produce model with greater sparsity than that of $l^2$ parameter regularization. Because the solution of (10) is a nonquadratic optimization problem, so there is no readily available closed form solution [14].

## 2.2. Model structure selection by D-optimality

A significant advantage due to orthogonalisation is that the contribution of model regressors to the model can be evaluated. The forward OLS estimator involves selecting a set of $n_\theta$ variables $\mathbf{p}_k = [p_k(1), \cdots, p_k(N)]^T$, $k = 1, \cdots, n_\theta$, from $M$ regressors to form a set of orthogonal basis $\mathbf{w}_k$, $k = 1, \cdots, n_\theta$, in a forward regression manner. As the orthogonality property $\mathbf{w}_i^T \mathbf{w}_j = 0$ for $i \neq j$ holds, if (6) is multiplied by itself and then the time average is taken, the following equation is easily derived

$$\frac{1}{N}\mathbf{y}^T\mathbf{y} = \frac{1}{N}\sum_{k=1}^{M}\gamma_k^2\mathbf{w}_k^T\mathbf{w}_k + \frac{1}{N}\Xi^T\Xi \qquad (11)$$

Conventional OLS [6] uses an Error Reduction Ratio $[ERR]$ to select a candidate regressor as the $k$th basis of the subset if it produces the largest value of $[ERR]_k$ from the remaining $(M - k + 1)$ candidates. By setting an appropriate tolerance $\rho$, which can be found by trial and error or via some statistical information criterion such as Akaike's information criterion(AIC) [9] that forms a compromise between the model performance and model complexity. Equivalently, this procedure can be expressed as

$$J^{(k)} = J^{(k-1)} - \frac{1}{N}\gamma_k^2\kappa_k \qquad (12)$$

where $J^{(0)} = \mathbf{y}^T\mathbf{y}$. At the $k$th forward regression stage, a candidate regressor is selected as the $k$th regressor if it produces the smallest $J^{(k)}$. (12) can be modified to form an alternative model selective criterion to enhance model robustness. D-optimality based cost function is one of robustness design criterion in experimental design criteria [10]. The D-optimality criterion is to maximize the determinant of the design matrix defined as $\mathbf{W}_k^T\mathbf{W}_k$, where $\mathbf{W}_k \in \Re^{N \times n_\theta}$ denotes the resultant regression matrix, consisting of $n_\theta$ regressors selected from $M$ regressors in $\mathbf{W}$.

$$\max\{J_D = \det(\mathbf{W}_k^T\mathbf{W}_k) = \prod_{k=1}^{n_\theta}\kappa_k\} \qquad (13)$$

It can be easily verified that the selection of the a subset of $\mathbf{W}_k$ from $\mathbf{W}$ is equivalent to the selection of the a subset of $n_\theta$ regressors from $\mathbf{P}$ [11]. In order to include D-optimality as a model selective criterion for improved model robustness, construct an augmented cost function as

$$\begin{aligned}J &= \frac{1}{N}\Xi^T\Xi + \alpha\log(\frac{1}{J_D}) \\ &= \frac{1}{N}(\mathbf{y}^T\mathbf{y} - \sum_{k=1}^{n_\theta}\gamma_k^2\kappa_k) + \alpha\sum_{k=1}^{n_\theta}\log[\frac{1}{\kappa_k}] \quad (14)\end{aligned}$$

where $\alpha$ is a positive small number. Note that this composite cost function simultaneously minimizes (12) and maximizes (13) [11]. Eq.(14) can be directly incorporated into

the forward OLS algorithm to select the most relevant $k$th regressor at the $k$th forward regression stage, via

$$J^{(k)} = J^{(k-1)} - \frac{1}{N}\gamma_k^2\kappa_k + \alpha\log[\frac{1}{\kappa_k}] \qquad (15)$$

Because $\log(\frac{1}{J_D})$ is an increasing function if $\kappa_k < 1$, which is true for some $k > K$, the selection procedure will terminate if $J^{(k)} \geq J^{(k-1)}$ at the derived model size $n_\theta$ if an proper $\alpha$ is set. The proposed approach can detect a parsimonious model size in an automatic manner.

## 3. MODEL IDENTIFICATION ALGORITHM USING FORWARD REGRESSION WITH BASIS PURSUIT AND D-OPTIMALITY

### 3.1. Parameter estimation by basis pursuit function's gradient descent

Before the introduction of the proposed algorithm, we initially introduce a general concept (algorithm) of parameter estimation by basis pursuit function's gradient descent, followed by the basic idea as how to incorporate this algorithm in the modified Gram-Schmidt orthogonal procedure.
*Theorem 1*(see [16] for the proof): Suppose that the dynamics underlying data set $D_N$ can be described by

$$y(t) = f(\mathbf{x}(t), \Theta) + \xi(t) \qquad (16)$$

where functional $f(\bullet)$ is given as appropriate. If the following parameter learning law is applied

$$\Theta(t+1) = \Theta(t) + \eta\overline{\xi(t)\frac{\partial f}{\partial\Theta}} - \eta\,\boldsymbol{\lambda}^T\,\text{sgn}(\Theta(t)) \quad (17)$$

where the operator $\overline{(\bullet)}$ denotes the time averaging, and $\text{sgn}(\Theta) = [\text{sgn}(\theta_1), ..., \text{sgn}(\theta_M)]^T$, in which,

$$\text{sgn}(u) = \begin{cases} 1 & \text{if } u > 0 \\ 0 & \text{if } u = 0 \\ -1 & \text{if } u < 0 \end{cases}$$

$\eta$ is an arbitrarily small positive number, then

$$\begin{array}{ll}\text{(i)} & \lim_{t\to+\infty}V(t) \to c \qquad\qquad (18)\end{array}$$

$$\begin{array}{ll}\text{(ii)} & \lim_{t\to+\infty}\|\Theta(t) - \Theta(t-k)\| = 0 \text{ for any finite } k\end{array}$$

where the basis pursuit cost function $V(t) = \frac{1}{2}\overline{\xi^2(t)} + \boldsymbol{\lambda}^T\|\Theta\|_1$, and $\|\Theta\|_1 = [|\theta_1|, ..., |\theta_{n_\theta}|]^T$ is constructed based on a subvector of $\Theta$ with nonzero parameters (see also (10)). $c = \min V(t)$ is the lower bound of $V(t)$.

In the proposed algorithm of Subsection 3.2, the above gradient descent of basis pursuit error function is combined with the modified Gram-Schmidt algorithm of Section 2.1 to derive a new model identification procedure. The basic

idea is introduced here. Consider (9), which can be regarded as a linear fitting of $\mathbf{y}^{(k-1)}$ by using a single variable $\mathbf{w}^{(k)}$ with the least squares method. The derived model residual vector $\Xi$ is then set as $\mathbf{y}^{(k)}$. This observation suggests that for each step $k$ in the modified Gram-Schmidt algorithm, the parameter estimates, calculated by (9) can be further tuned by learning algorithm of (17) that optimizes the basis pursuit's function given by (10). Following (9), denote $\mathbf{y}^{(k-1)} = [y^{(k-1)}(1), y^{(k-1)}(2), ..., y^{(k-1)}(N)]^T$ and $\mathbf{w}_k = [w_k(1), ..., w_k(N)]^T$. The tuning process is an extremely simple case based on Theorem 1, as illustrated by the following Theorem.

*Theorem 2* (see [16] for the proof): If the learning law given by (17) is applied to a special case of one dimensional linear system

$$y^{(k-1)}(t) = \gamma_k w_k(t) + \xi(t) \qquad (19)$$

with the parameter estimates $\gamma_k$ initialized as the least square parameter estimate $\gamma_k^{(0)} \neq 0$, given by (9), and if $\lambda_k < \frac{1}{2N}|\mathbf{w}_k^T \mathbf{y}|$, then the final converged parameter estimate $\gamma_k$

$$\text{(i)} \quad |\gamma_k| < |\gamma_k^{(0)}|$$
$$\text{(ii)} \quad \text{sgn}(\gamma_k) = \text{sgn}(\gamma_k^{(0)}) \qquad (20)$$

The significance of Theorem 2 is that by setting the basis pursuit parameters $\lambda_k$ below a certain value, for each step $k$, the overall effect of the tuning process is that the parameters $\gamma_k$ is pulled towards 0. In forward regression, as model size $k$ increases, the parameter estimates $\gamma_k$, as initialized by least squares algorithm with very small magnitudes, followed by basis pursuit gradient tuning, will shrink below some threshold value, and can therefore be obtained as zero, to achieve model sparsity. For a sufficiently small $\lambda_k$, the optimality condition can be derived as

$$\overline{\xi(t)w_k(t)} - \lambda_k \, \text{sgn}(\gamma_k(t)) = 0 \qquad (21)$$

or

$$\gamma_k = \frac{\mathbf{w}_k^T \mathbf{y}^{(k-1)} - N\lambda_k \text{sgn}(\gamma_k)}{\mathbf{w}_k^T \mathbf{w}_k}$$
$$= \gamma_k^{(0)} - \frac{N\lambda_k \text{sgn}(\gamma_k^{(0)})}{\mathbf{w}_k^T \mathbf{w}_k} \qquad (22)$$

### 3.2. The new algorithm using combined modified Gram-Schmidt algorithm, basis pursuit and D-optimality

In this Section a new algorithm is introduced that combines the modified Gram-Schmidt algorithm with the basis pursuit gradient tuning for new parameter estimation. The model selective criteria by D-optimality of Section 2.2 [11] is applied in the proposed algorithm. The algorithm is introduced as follows, in which, the basis pursuit parameters are initially assumed to be predetermined.

*The modified Gram-Schmidt algorithm combining basis pursuit and D-optimality:*

The Gram-Schmidt orthogonalisation scheme can be used to derive a simple and efficient algorithm for selecting subset models. Introducing the definition of $\mathbf{P}^{(k-1)}$ as

$$\mathbf{P}^{(k-1)} = [\mathbf{w}_1, \cdots, \mathbf{w}_{k-1}, \mathbf{p}_k^{(k-1)}, \cdots, \mathbf{p}_M^{(k-1)}] \qquad (23)$$

If some of the columns $\mathbf{p}_k^{(k-1)}, \cdots, \mathbf{p}_M^{(k-1)}$ in $\mathbf{P}^{(k-1)}$ have been interchanged, this will still be referred to as $\mathbf{P}^{(k-1)}$ for notational convenience. The $k$th stage of the forward regression selection procedure is given below

1. For $k \leq j \leq M$, compute

$$\gamma_k^{(j)} = \frac{(\mathbf{p}_j^{(k-1)})^T \mathbf{y}^{(k-1)}}{(\mathbf{p}_j^{(k-1)})^T \mathbf{p}_j^{(k-1)}} \qquad (24)$$

$$J_j^{(k)} = J^{(k-1)} - \frac{1}{N}[\gamma_k^{(j)}]^2 \kappa_k^{(j)} + \alpha \log[\frac{1}{\kappa_k^{(j)}}] \qquad (25)$$

2. Find

$$J^{(k)} = J_{j_k}^{(k)} = \min\{J_j^{(k)}, \quad k \leq j \leq M\} \qquad (26)$$

Then the $j_k$th column of $\mathbf{P}^{(k-1)}$ is interchanged with the $k$th column of $\mathbf{P}^{(k-1)}$, and the $j_k$th column of $\mathbf{A}$ up to the $(k-1)$th row is interchanged with the $k$th column of $\mathbf{A}$. This effectively selects the $j_k$th candidates as the $k$th regressor in the subset model. Then set $\gamma_k^{(0)} = \gamma_k^{(j_k)}$.

3. Perform the orthogonalization as follows

$$\mathbf{w}_k = \mathbf{p}_k^{(k-1)}$$
$$\alpha_{kj} = \frac{\mathbf{w}_k^T \mathbf{p}_j^{(k-1)}}{\mathbf{w}_k^T \mathbf{w}_k}, \quad k+1 \leq j \leq M$$
$$\mathbf{p}_j^{(k)} = \mathbf{p}_j^{(k-1)} - \alpha_{kj} \mathbf{w}_k, \quad k+1 \leq j \leq M \qquad (27)$$

to transform $\mathbf{P}^{(k-1)}$ into $\mathbf{P}^{(k)}$ and derive the $k$th row of $\mathbf{A}$. Update $\kappa_k$.

4. With $\gamma_k^{(0)} \neq 0$ as initialized parameter estimates, the optimal solution of learning law (17) is given by (22), and is rewritten here

$$\gamma_k = \gamma_k^{(0)} - \frac{N\lambda_k \text{sgn}(\gamma_k^{(0)})}{\mathbf{w}_k^T \mathbf{w}_k} \qquad (28)$$

where $\lambda_k < \frac{1}{2N}|\mathbf{w}_k^T \mathbf{y}^{(k-1)}|$.

5. Update $\mathbf{y}^{(k-1)}$ into $\mathbf{y}^{(k)}$ by

$$\mathbf{y}^{(k)} = \mathbf{y}^{(k-1)} - \gamma_k \mathbf{w}_k \qquad (29)$$

and update

$$J^{(k)} = J^{(k-1)} - \frac{1}{N}\gamma_k^2 \kappa_k + \alpha \log[\frac{1}{\kappa_k}] \quad (30)$$

6. The selection is terminated at the $n_\theta$th stage where a subset model containing $n_\theta$ significant regressors by the D-optimality model selective criteria $J^{(k)}$ achieves a minimum.

It is shown by analysis [11] that if the parameter estimates are initialized with very small magnitudes from least squares estimates, the basis pursuit gradient tuning procedure of Step 4, will pull it even more towards zero by applying Theorem 2. The conclusion is that the proposed algorithm can achieve a sparser model than that of without basis pursuit gradient tuning procedure. The identification algorithm introduced above uses a predetermined basis pursuit parameters $\lambda$, which reflects a tradeoff between modelling errors and the $l^1$ norm of parameter vector. By the general principle in data modelling of that a model with generalization is preferred, the choice of $\lambda$ may be derived based on the commonly used method of cross-validation. In the following, a simple method of choosing $\lambda$ is introduced, by the basic principle of cross-validation. i.e. using two data sets, one for training and another for testing. For simplicity a single global basis pursuit $\lambda$ is used, that is, $\lambda_1 = \lambda_2 = ... = \lambda$. The complete modelling procedure of iterating the proposed algorithm, by incrementally increasing $\lambda$ from zero in a controlled manner, is given as follows.

*The iterative procedure of the proposed algorithm including choosing basis pursuit parameters*

1. Initialization. Set an arbitrarily small $\alpha$, applying the modelling procedure of [11] to derive a model with size $n_\theta^{(0)}$. (This is equivalent to the proposed algorithm with $\lambda = 0$.) and set $\lambda = \frac{1}{2N}|\mathbf{w}_{n_\theta^{(0)}}^T \mathbf{y}|$. Set a counter for iteration $j = 1$;

2. Applying the proposed algorithm with the new $\lambda$, to derive a model with the size of $n_\theta^{(j)} < n_\theta^{(j-1)}$. Set a new $\lambda = \frac{1}{2N}|\mathbf{w}_{n_\theta^{(j)}}^T \mathbf{y}|$ for next iteration of this step, while the mean squares errors (MSE) of the test data set is monitored; $j = j + 1$;

3. Step 2 is terminated when the MSE of the test data set achieves a minimum.

It is shown by analysis [11] that as the iteration step $j$ increases, the effect of basis pursuit cost function (shrinking the small parameters to zero) would derive at the smaller size $n_\theta^{(j)}$ compared to previous iteration step. Alternatively, $\lambda$ can be set as a very small value for general improvement in model sparseness.

## 4. ILLUSTRATIVE EXAMPLE

Consider the chaotic two dimensional time series, *Ikeda map* [17], given by

$$\begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \begin{bmatrix} 1 + 0.9[x(t-1)\cos(r) - y(t-1)\sin(r)] \\ 0.9[x(t-1)\sin(r) + y(t-1)\cos(r)] \end{bmatrix}$$

$$\text{with} \qquad r = 0.4 - \frac{6.0}{1 + x^2(t-1) + y^2(t-1)} \qquad (31)$$

1000 data points were generated with an initial condition $x(1) = 0.1, y(1) = 0.1$. Two models were constructed to model $x(t)$ and $y(t)$ respectively. For both models, the input vector is set as $\mathbf{x}(t) = [x(t-1), y(t-1)]^T$. 498 data samples from $t = 1 \sim 500$, were used as estimation set, and 500 data samples $t = 499 \sim 1000$ were used as test data. The Gaussian radial basis function was used to construct full model sets by using all the data in the estimation data set as centers $\mathbf{c}_i$, $i = 1, ...498$, and $p_i(\mathbf{x}(t)) = \exp\{-\frac{\|\mathbf{x}(t)-\mathbf{c}_i\|^2}{\sigma_i^2}\}$, with $\sigma_i = 0.5, \forall i$. For the first model that models $x(t)$, the modelling starts with $\lambda = 0$, and $\alpha = 10^{-8}$ (an arbitrarily small coefficient for D-optimality). The iterative procedure of the proposed algorithm was applied. The model was automatically terminated at a 63 centers networks. The final basis pursuit parameter was derived at $\lambda = 7.7 \times 10^{-8}$. The modelling MSE for the test data set is derived at $3.13 \times 10^{-5}$. Equivalently $99.81\%$ output variance of the test data has been explained by the model. For the second model that models $y(t)$, the modelling starts with $\lambda = 0$, and $\alpha = 10^{-8}$ (an arbitrarily small coefficient for D-optimality). The iterative procedure of the proposed algorithm was applied. The model was automatically terminated at a 66 centers networks. The final basis pursuit parameter was derived at $\lambda = 1.7 \times 10^{-8}$. The modelling MSE for the test data set is derived at $1.36 \times 10^{-5}$. Equivalently $99.94\%$ output variance of the test data has been explained by the model. To illustrate the overall performance of the model in capturing the underlying system dynamics, the modelling results for both estimation and test data set is shown in Figure 1.

## 5. CONCLUSIONS

This paper has introduced a new forward regression model identification algorithm combining the modified Gram-Schmidt algorithm with basis pursuit and D-optimality design.

Figure 1: Modelling results for illustrative example.

## 6. REFERENCES

[1] R. Bellman, *Adaptive Control Processes.* Princeton University Press, 1966.

[2] C. J. Harris, X. Hong, and Q. Gan, *Adaptive Modelling, Estimation and Fusion from Data: A Neuro-fuzzy Approach.* Springer-Verlag, 2002.

[3] M. Brown and C. J. Harris, *Neurofuzzy Adaptive Modelling and Control.* Prentice Hall, Hemel Hempstead, 1994.

[4] J. S. R. Jang, C. Sun, and E. Mizutani, *Neuro-Fuzzy and Soft Computing: A Computational Approach to Learning and Machine Intelligence.* Upper Saddle River, NJ : Prentice Hall, 1997.

[5] T. Takagi and M. Sugeno, "Fuzzy identification of systems and its applications to modelling and control," *IEEE Trans. on Systems, Man, and Cybernetics*, vol. 15, pp. 116–132, 1985.

[6] S. Chen, S. A. Billings, and W. Luo, "Orthogonal least squares methods and their applications to non-linear system identification," *International Journal of Control*, vol. 50, pp. 1873–1896, 1989.

[7] S. Chen, Y. Wu, and B. L. Luk, "Combined genetic algorithm optimization and regularized orthogonal least squares learning for radial basis function networks," *IEEE Trans. on Neural Networks*, vol. 10, pp. 1239–1243, 1999.

[8] M. J. L. Orr, "Regularisation in the selection of radial basis function centers," *Neural Computation*, vol. 7, no. 3, pp. 954–975, 1995.

[9] H. Akaike, "A new look at the statistical model identification," *IEEE Trans. on Automatic Control*, vol. AC-19, pp. 716–723, 1974.

[10] A. C. Atkinson and A. N. Donev, *Optimum Experimental Designs.* Clarendon Press, Oxford, 1992.

[11] X. Hong and C. J. Harris, "Nonlinear model structure design and construction using orthogonal least squares and d-optimality design," *IEEE Transactions on Neural Networks*, vol. 13, no. 5, pp. 1245–1250, 2001.

[12] S. Chen, X. Hong, and C. J. Harris, "Sparse kernel regression modelling using combined locally regularised orthogonal least squares and d-optimality experimental design," *IEEE Trans. on Automatic Control*, vol. 48, no. 6, pp. 1029–1036, 2003.

[13] D. J. C. MacKay, "Bayesian methods for adaptive models," Ph.D. dissertation, California Institute of Technology, USA, 1991.

[14] S. S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM Review*, vol. 43, no. 1, pp. pp129–159, 2001.

[15] B. Efron, I. Johnstone, T. Hastie, and R. Tibshirani, "Least angle regression," *Annals of Statistics*, p. To Appear, 2003.

[16] X. Hong, M. Brown, S. Chen, and C. J. Harris, "Sparse model identification using orthogonal forward regression with basis pursuit and d-optimality," *Submitted to IEE Proc. - Control Theory and Applications*, 2003.

[17] K. Ikeda, "Multiple-valued stationary state and its instability of the transmitted light by a rign cavity system," *Optics Communications*, vol. 30, no. 2, pp. 257–261, 1979.