

l1-norm penalized orthogonal forward regression

Article

Accepted Version

Hong, X., Chen, S., Guo, Y. and Gao, J. (2017) l1-norm penalized orthogonal forward regression. *International Journal of Systems Science*, 48 (10). pp. 2195-2201. ISSN 0020-7721 doi: <https://doi.org/10.1080/00207721.2017.1311383> Available at <http://centaur.reading.ac.uk/72078/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

To link to this article DOI: <http://dx.doi.org/10.1080/00207721.2017.1311383>

Publisher: Taylor & Francis

All outputs in CentAUR are protected by Intellectual Property Rights law, including copyright law. Copyright and IPR is retained by the creators or other copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

www.reading.ac.uk/centaur

CentAUR

Central Archive at the University of Reading

Reading's research outputs online

l^1 -norm Penalized Orthogonal Forward Regression

Xia Hong^{a*}, Sheng Chen^b, Yi Guo^c and Junbin Gao^d

^a*Department of Computer Science, School of Mathematical, Physical and Computational Sciences, University of Reading, Reading, RG6 6AY;* ^b*Electronics and Computer Science, University of Southampton, Southampton SO17 1BJ, UK, and also with King Abdulaziz University, Jeddah 21589, Saudi Arabia;* ^c*CSIRO Mathematics and Information Sciences, North Ryde, NSW 1670, Australia;* ^d*Discipline of Business Analytics, University of Sydney Business School, University of Sydney, Camperdown NSW 2006, Australia*

(v5.0 released May 2016)

A l^1 -norm penalized orthogonal forward regression (l^1 -POFR) algorithm is proposed based on the concept of leave-one-out mean square error (LOOMSE), by defining a new l^1 -norm penalized cost function in the constructed orthogonal space and associating each orthogonal basis with an individually tunable regularization parameter. Due to orthogonality, the LOOMSE can be analytically computed without actually splitting the data set, and moreover a closed form of the optimal regularization parameter is derived by greedily minimizing the LOOMSE incrementally. We also propose a simple formula for adaptively detecting and removing regressors to an inactive set so that the computational cost of the algorithm is significantly reduced. Examples are included to demonstrate the effectiveness of this new l^1 -POFR approach.

Keywords: Cross validation, forward regression, leave-one-out errors, regularization

1. Introduction

One of the main aims in data modeling is good generalization, i.e. the model's capability to approximate accurately the system output for unseen data. Sparse models can be constructed using the l^1 -penalized cost function, e.g., the basis pursuit or least absolute shrinkage and selection operator (LASSO) (Chen et al., 1998; Efron et al., 2004; Tibshirani, 1996). Based on a fixed single l^1 -penalized regularization parameter, the LASSO can be configured as a standard quadratic programming optimization problem. By exploiting piecewise linearity of the problem, the least angle regression (LAR) procedure (Efron et al., 2004) was developed for solving the problem efficiently. Note that the computational efficiency in LASSO is facilitated by a *single* regularization parameter setting. For more complicated constraints, e.g., multiple regularizers, the cross validation by actually splitting data sets as the means of evaluating model generalization comes with considerably large overall computational overheads.

Fundamental to evaluate model generalization capability is the concept of cross-validation (Rao et al., 2008; Stone, 1974), and one commonly used version of cross-validation is the leave-one-out (LOO) cross validation. For the linear-in-the-parameters models, the LOO mean square error (LOOMSE) can be calculated without actually splitting the training data set and estimating the associated models, by making use of Sherman-Morrison-Woodbury theorem (Sherman and Morrison, 1950). Using the LOOMSE as the model term selective criterion, an orthogonal forward regression

*Corresponding author. Email: x.hong@reading.ac.uk

(OFR) procedure was introduced in (Hong et al., 2003). Furthermore, the l^2 -norm based regularization techniques (MacKay, 1991; Orr, 1995) were incorporated into the orthogonal least squares (OLS) algorithm of (Chen et al., 1989) to produce a regularized OLS algorithm that carries out model term selection while reduces the variance of parameter estimate simultaneously (Chen et al., 2003). The optimization of l^1 -norm regularizer with respect to model generalization analytically is however less studied (Ji et al, 2008).

We propose a l^1 -norm penalized OFR (l^1 -POFR) algorithm to carry out the regularizer optimization as well as model term selection and parameter estimation simultaneously in an OFR manner. The algorithm is based on a new l^1 -norm penalized cost function with multiple l^1 regularizers, each of which is associated with an orthogonal basis vector, by orthogonal decomposition of the regression matrix of the selected model terms. We derive a closed form of the optimal regularization parameter by greedily minimizing the LOOMSE incrementally. To save computational costs an inactive set is used along the OFR process by predicting whether any model terms will be unselectable in future regression steps.

2. Preliminaries

Consider the general nonlinear system represented by the nonlinear model (Chen and Billings, 1989):

$$y(k) = f(\mathbf{x}(k)) + v(k), \quad (1)$$

where $\mathbf{x}(k) = [x_1(k) \ x_2(k) \ \cdots \ x_m(k)]^T \in \mathbb{R}^m$ denotes the input vector at sample time index k and $y(k)$ is the system output variable, respectively, while $v(k)$ denotes the system white noise and $f(\bullet)$ is the unknown system mapping.

The unknown system (1) is to be identified based on an observation data set $D_N = \{\mathbf{x}(k), y(k)\}_{k=1}^N$ using a linear-in-the-parameters model of the form:

$$\hat{y}^{(M)}(k) = f^{(M)}(\mathbf{x}(k)) = \sum_{i=1}^M \theta_i \phi_i(\mathbf{x}(k)), \quad (2)$$

where $\hat{y}^{(M)}(k)$ is the model prediction output for $\mathbf{x}(k)$ based on the M -term regression model, and M is the total number of nonlinear regressors, while θ_i are the model weights. While there exist many suitable choices for regressor, without loss of generality, we choose $\phi_i(\mathbf{x})$ to be Gaussian radial basis function (RBF)

$$\phi_i(\mathbf{x}) = e^{-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\tau^2}} \quad (3)$$

in which $\mathbf{c}_i = [c_{1,i} \ c_{2,i} \ \cdots \ c_{m,i}]^T$ is known as the center vector of the i th RBF unit and τ is an RBF width parameter. We assume that each RBF unit is placed on a training data, namely, all the RBF center vectors $\{\mathbf{c}_i\}_{i=1}^M$ are selected from the training data $\{\mathbf{x}(k)\}_{k=1}^N$, and the RBF width τ has been predetermined, for example, using cross validation.

Let us denote $e^{(M)}(k) = y(k) - \hat{y}^{(M)}(k)$ as the M -term modeling error for the input data $\mathbf{x}(k)$. Over the training data set D_N , further denote $\mathbf{y} = [y(1) \ y(2) \ \cdots \ y(N)]^T$, $\mathbf{e}^{(M)} = [e^{(M)}(1) \ e^{(M)}(2) \ \cdots \ e^{(M)}(N)]^T$, and $\Phi_M = [\phi_1 \ \phi_2 \ \cdots \ \phi_M]$ with $\phi_n = [\phi_n(\mathbf{x}(1)) \ \phi_n(\mathbf{x}(2)) \ \cdots \ \phi_n(\mathbf{x}(N))]^T$, $1 \leq n \leq M$. We have the M -term model in the matrix form of

$$\mathbf{y} = \Phi_M \boldsymbol{\theta}_M + \mathbf{e}^{(M)}, \quad (4)$$

where $\boldsymbol{\theta}_M = [\theta_1 \ \theta_2 \ \cdots \ \theta_M]^\top$. Let an orthogonal decomposition of the regression matrix $\boldsymbol{\Phi}_M$ be

$$\boldsymbol{\Phi}_M = \mathbf{W}_M \mathbf{A}_M, \quad (5)$$

where

$$\mathbf{A}_M = \begin{bmatrix} 1 & a_{1,2} & \cdots & a_{1,M} \\ 0 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & a_{M-1,M} \\ 0 & \cdots & 0 & 1 \end{bmatrix} \quad (6)$$

and

$$\mathbf{W}_M = [\mathbf{w}_1 \ \mathbf{w}_2 \ \cdots \ \mathbf{w}_M] \quad (7)$$

with columns satisfying $\mathbf{w}_i^\top \mathbf{w}_j = 0$, if $i \neq j$. The regression model (4) can alternatively be expressed as

$$\mathbf{y} = \mathbf{W}_M \mathbf{g}_M + \mathbf{e}^{(M)}, \quad (8)$$

where $\mathbf{g}_M = [g_1 \ g_2 \ \cdots \ g_M]^\top$ satisfies the triangular system $\mathbf{A}_M \boldsymbol{\theta}_M = \mathbf{g}_M$, which can be used to determine the original model parameter vector $\boldsymbol{\theta}_M$, given \mathbf{A}_M and \mathbf{g}_M . The space spanned by the original model bases $\boldsymbol{\phi}_n$, $1 \leq n \leq M$, is the same space spanned by the orthogonal model bases \mathbf{w}_n , $1 \leq n \leq M$. Also since only the k th row of $\boldsymbol{\Phi}_M$ depends on $\mathbf{x}(k)$, only the k th row of \mathbf{W}_M depends on $\mathbf{x}(k)$.

Further consider the following weighted l^1 -norm penalized OLS criterion for the model (8)

$$L_e(\boldsymbol{\Lambda}_M, \mathbf{g}_M) = \|\mathbf{y} - \mathbf{W}_M \mathbf{g}_M\|^2 + \sum_{i=1}^M \lambda_i |g_i|, \quad (9)$$

where $\boldsymbol{\Lambda}_M = \text{diag}\{\lambda_1, \lambda_2, \cdots, \lambda_M\}$, which contains the regularization parameters $\lambda_i \geq \varepsilon$, $1 \leq i \leq M$, and $\varepsilon > 0$ is a predetermined lower bound for the regularization parameters. Given $\boldsymbol{\Lambda}_M$, the solution for \mathbf{g}_M can be obtained by setting the subderivative vector of L_e to zero, i.e. $\frac{\partial L_e}{\partial \mathbf{g}_M} = \mathbf{0}$, yielding

$$g_i^{(\text{olasso})} = \left(|g_i^{(\text{LS})}| - \frac{\lambda_i/2}{\mathbf{w}_i^\top \mathbf{w}_i} \right)_+ \text{sign}(g_i^{(\text{LS})}) \quad (10)$$

for $1 \leq i \leq M$, with the usual least squares solution given by $g_i^{(\text{LS})} = \frac{\mathbf{w}_i^\top \mathbf{y}}{\mathbf{w}_i^\top \mathbf{w}_i}$, and the operator $(\)_+$

$$z_+ = \begin{cases} z, & \text{if } z > 0, \\ 0, & \text{if } z \leq 0. \end{cases} \quad (11)$$

Unlike the LASSO (Chen et al., 1998; Tibshirani, 1996), our objective $L_e(\boldsymbol{\Lambda}_M, \mathbf{g}_M)$ is constructed on the orthogonal space and the l^1 -norm parameter constraints are associated with the orthogonal bases \mathbf{w}_i , $1 \leq i \leq M$. Since the cost function (9) contains sparsity inducing l^1 norm, some parameters $g_i^{(\text{olasso})}$ will be returned as zeros, producing a sparse model in the orthogonal space spanned

by the columns of \mathbf{W}_M , which corresponds to a sparse model in the original space spanned by the columns of Φ_M .

3. Regularization parameter optimization and model construction with LOOMSE

Each OFR stage involves the joint regularization parameter optimization, model term selection and parameter estimation. The regularization parameters with respect to their associated candidate regressors are optimized using the approximate LOOMSE formula that is derived in Section 3.2, and the regressor with the smallest LOOMSE is selected. This OFR procedure is inherently suboptimal as it is based on greedy incremental optimization.

3.1. Model representation and LOOMSE in n -th stage OFR

Consider the OFR modeling process that has produced the $(n-1)$ -term model. The model output vector of this $(n-1)$ -term model is given by

$$\hat{\mathbf{y}}^{(n-1)} = \sum_{i=1}^{n-1} g_i^{(\text{lasso})} \mathbf{w}_i, \quad (12)$$

and we denote the corresponding modeling error vector by $\mathbf{e}^{(n-1)} = \mathbf{y} - \hat{\mathbf{y}}^{(n-1)}$.

Consider the n th OFR stage where n columns of regressors are constructed as $\mathbf{W}_n = [\mathbf{w}_1 \ \mathbf{w}_2 \ \cdots \ \mathbf{w}_n]$, with $\mathbf{w}_i = [w_i(1) \ w_i(2) \ \cdots \ w_i(N)]^T$, $i = 1, \dots, n$. Clearly, the n th OFR stage can be represented by

$$\mathbf{e}^{(n-1)} = g_n \mathbf{w}_n + \mathbf{e}^{(n)}. \quad (13)$$

The model form (13) illustrates the fact that the n th OFR stage is simply to fit a one-variable model using the current model residual produced after the $(n-1)$ th stage as the desired system output. Since $\mathbf{w}_n^T \hat{\mathbf{y}}^{(n-1)} = 0$, it is easy to verify that $g_n^{(\text{LS})} = \frac{\mathbf{w}_n^T \mathbf{y}}{\mathbf{w}_n^T \mathbf{w}_n} = \frac{\mathbf{w}_n^T \mathbf{e}^{(n-1)}}{\mathbf{w}_n^T \mathbf{w}_n}$.

The selection of one regressor from the candidate regressors involves initially generating candidate \mathbf{w}_n by making each candidate regressor to be orthogonal to the $(n-1)$ orthogonal basis vectors, \mathbf{w}_i for $1 \leq i \leq n-1$ obtained in the previous $(n-1)$ OFR stages, followed by evaluating their contributions. Consider the case of $2|\mathbf{w}_n^T \mathbf{e}^{(n-1)}| > \varepsilon$. Applying (10) to (13), we note that clearly as λ_n decreases away from $2|\mathbf{w}_n^T \mathbf{e}^{(n-1)}|$ towards ε , $g_n^{(\text{lasso})}$ increases its magnitude at a linear rate to λ_n , from zero to an upper bound $|g_n^{(\text{B})}|$ with

$$g_n^{(\text{B})} = \left(|g_n^{(\text{LS})}| - \frac{\varepsilon}{2\mathbf{w}_n^T \mathbf{w}_n} \right)_+ \text{sign}(g_n^{(\text{LS})}). \quad (14)$$

For any candidate regressor, it is vital that we evaluate its potential model generalization performance using the most suitable value of λ_n . The optimization of the LOOMSE with respect to λ_n is detailed in Section 3.2, based on the idea of the LOO cross validation outlined below.

Suppose that we sequentially set aside each data point in the estimation set D_N in turn and estimate a model using the remaining $(N-1)$ data points. The prediction error is calculated on the data point that has not been used in estimation. That is, for $k = 1, 2, \dots, N$, the models are estimated based on $D_N \setminus (\mathbf{x}(k), y(k))$, respectively, and the outputs are denoted as $\hat{\mathbf{y}}^{(n-1, -k)}(k, \lambda_n)$.

Then, the LOO prediction error based on the k th data sample is calculated as

$$e^{(n,-k)}(k, \lambda_n) = y(k) - \hat{y}^{(n-1,-k)}(k, \lambda_n). \quad (15)$$

The LOOMSE is defined as the average of all these prediction errors, given by $J(\lambda_n) = E \left[(e^{(n,-k)}(k, \lambda_n))^2 \right]$. Thus the optimal regularization parameter for the n th stage is given by

$$\lambda_n^{\text{opt}} = \arg \min_{\lambda_n} \left\{ J(\lambda_n) = \frac{1}{N} \sum_{k=1}^N (e^{(n,-k)}(k, \lambda_n))^2 \right\}. \quad (16)$$

Evaluation of $J(\lambda_n)$ by directly splitting the data set requires extensive computational efforts. We show in Section 3.2 that $J(\lambda_n)$ can be approximately calculated without actually sequentially splitting the estimation data set. Furthermore, we also show that the optimal value λ_n^{opt} can be obtained in a closed-form expression in the orthogonal modeling space.

3.2. Optimal regularization parameter estimate

We notice from (10) that $g_n^{(\text{olasso})} = 0$ if $2|\mathbf{w}_n^T \mathbf{e}^{(n-1)}| < \lambda_n$, and thus a sufficient condition that a given \mathbf{w}_n may be excluded from the candidate pool without explicitly determining λ_n is $2|\mathbf{w}_n^T \mathbf{e}^{(n-1)}| < \varepsilon$, which is the regularizer's lower bound, a preset value indicating the correlation of the candidate regressor. Hence, in the following we assume that $2|\mathbf{w}_n^T \mathbf{e}^{(n-1)}| > \varepsilon$, and we have

$$\mathbf{g}_n^{(\text{olasso})} = \mathbf{H}_n^{-1} \left(\mathbf{W}_n^T \mathbf{y} - \Lambda_n \text{sign}(\mathbf{g}_n^{(\text{LS})})/2 \right), \quad (17)$$

where $\mathbf{g}_n^{(\text{olasso})} = [g_1^{(\text{olasso})} \ g_2^{(\text{olasso})} \ \dots \ g_n^{(\text{olasso})}]^T$, $\text{sign}(\mathbf{g}_n) = [\text{sign}(g_1) \ \text{sign}(g_2) \ \dots \ \text{sign}(g_n)]^T$, and $\mathbf{H}_n = \mathbf{W}_n^T \mathbf{W}_n$. Note that (17) is consistent to (10) for all terms with nonzero g_i . In the OFR procedure, any candidate terms \mathbf{w}_i producing zero $g_i^{(\text{olasso})}$ will not be selected since they will not contribute to any reduction in the LOOMSE.

The model residual is defined by

$$e^{(n)}(k, \lambda_n) = y(k) - (\mathbf{g}^{(\text{olasso})})^T \mathbf{w}(k) = y(k) - \left(\mathbf{y}^T \mathbf{W}_n - (\text{sign}(\mathbf{g}^{(\text{LS})}))^T \Lambda_n/2 \right) \mathbf{H}_n^{-1} \mathbf{w}(k), \quad (18)$$

where $\mathbf{w}(k)$ denotes the transpose of the k th row of \mathbf{W}_n . If the data sample indexed at k is removed from the estimation data set, the LOO parameter estimator obtained by using only the $(N-1)$ remaining data points is given by

$$\mathbf{g}_n^{(\text{olasso},-k)} = (\mathbf{H}_n^{(-k)})^{-1} \left((\mathbf{W}_n^{(-k)})^T \mathbf{y}^{(-k)} - \Lambda_n \text{sign}(\mathbf{g}^{(\text{LS},-k)})/2 \right) \quad (19)$$

where $\mathbf{H}_n^{(-k)} = (\mathbf{W}_n^{(-k)})^T \mathbf{W}_n^{(-k)}$, $\mathbf{W}_n^{(-k)}$ and $\mathbf{y}^{(-k)}$ are the resultant regression matrix and desired output vector, respectively, by removing $(\mathbf{x}(k), y(k))$, i.e., $(\mathbf{w}^T(k), y(k))$, from $\mathbf{W}(k)$ and $\mathbf{y}(k)$. Thus we have

$$\mathbf{H}_n^{(-k)} = \mathbf{H}_n - \mathbf{w}(k) \mathbf{w}^T(k), \quad (20)$$

$$(\mathbf{y}^{(-k)})^T \mathbf{W}_n^{(-k)} = \mathbf{y}^T \mathbf{W}_n - y(k) \mathbf{w}^T(k). \quad (21)$$

The LOO error evaluated at k is given by

$$\begin{aligned} e^{(n,-k)}(k, \lambda_n) &= y(k) - (\mathbf{g}^{(\text{lasso}, -k)})^T \mathbf{w}(k) \\ &= y(k) - \left((\mathbf{y}^{(-k)})^T \mathbf{W}_n^{(-k)} - \right. \\ &\quad \left. (\text{sign}(\mathbf{g}^{(\text{LS}, -k)}))^T \boldsymbol{\Lambda}_n / 2 \right) (\mathbf{H}_n^{(-k)})^{-1} \mathbf{w}(k). \end{aligned} \quad (22)$$

Applying the matrix inversion lemma to (20) yields

$$\begin{aligned} (\mathbf{H}_n^{(-k)})^{-1} &= (\mathbf{H}_n - \mathbf{w}(k) \mathbf{w}^T(k))^{-1} \\ &= \mathbf{H}_n^{-1} + \frac{\mathbf{H}_n^{-1} \mathbf{w}(k) \mathbf{w}^T(k) \mathbf{H}_n^{-1}}{1 - \mathbf{w}^T(k) \mathbf{H}_n^{-1} \mathbf{w}(k)} \end{aligned} \quad (23)$$

and

$$(\mathbf{H}_n^{(-k)})^{-1} \mathbf{w}(k) = \frac{\mathbf{H}_n^{-1} \mathbf{w}(k)}{1 - \mathbf{w}^T(k) \mathbf{H}_n^{-1} \mathbf{w}(k)}. \quad (24)$$

Substituting (21) and (24) into (22) yields

$$\begin{aligned} e^{(n,-k)}(k, \lambda_n) &= y(k) - \left(\mathbf{y}^T \mathbf{W}_n - y(k) \mathbf{w}^T(k) - \right. \\ &\quad \left. (\text{sign}(\mathbf{g}^{(\text{LS}, -k)}))^T \boldsymbol{\Lambda}_n / 2 \right) \frac{\mathbf{H}_n^{-1} \mathbf{w}(k)}{1 - \mathbf{w}^T(k) \mathbf{H}_n^{-1} \mathbf{w}(k)} \\ &= \frac{y(k) - \left(\mathbf{y}^T \mathbf{W}_n - (\text{sign}(\mathbf{g}^{(\text{LS}, -k)}))^T \boldsymbol{\Lambda}_n / 2 \right) \mathbf{H}_n^{-1} \mathbf{w}(k)}{1 - \mathbf{w}^T(k) \mathbf{H}_n^{-1} \mathbf{w}(k)}. \end{aligned} \quad (25)$$

Assuming that $\text{sign}(\mathbf{g}_n^{(\text{LS}, -k)}) = \text{sign}(\mathbf{g}_n^{(\text{LS})})$ holds for most data samples in D_N , and applying (18) to (25), we have

$$e^{(n,-k)}(k, \lambda_n) = \gamma_n(k) e^{(n)}(k, \lambda_n), \quad (26)$$

where $\gamma_n(k) = \frac{1}{1 - \sum_{i=1}^n (w_i(k))^2 / \mathbf{w}_i^T \mathbf{w}_i} > 0$, and $w_i(k)$ is the k th element of \mathbf{w}_i . The LOOMSE can then be calculated as

$$J(\lambda_n) = \frac{1}{N} \sum_{k=1}^N \gamma_n^2(k) (e^{(n)}(k, \lambda_n))^2. \quad (27)$$

Note that for $\text{sign}(\mathbf{g}_n^{(\text{LS}, -k)})$ and $\text{sign}(\mathbf{g}_n^{(\text{LS})})$ to be different, each element in $\mathbf{g}_n^{(\text{LS})}$ needs to be very close to zero, which is unlikely since only the model terms satisfying $|\mathbf{w}_n^T \mathbf{e}^{(n-1)}| > \varepsilon/2$ are considered. Hence we can treat $J(\lambda_n)$ given in (27) as the exact LOOMSE for any ε that is not too small.

We further represent (18) as

$$e^{(n)}(k, \lambda_n) = \eta(k) + \frac{\lambda_n}{2\mathbf{w}_n^T \mathbf{w}_n} w_n(k) \text{sign}(g_n^{(\text{LS})}), \quad (28)$$

where $\eta(k) = e^{(n-1)}(k) - g_n^{(\text{LS})} w_n(k)$ is the model residual obtained based on the least square estimate at the n th step stage. By setting $\frac{\partial J(\lambda_n)}{\partial \lambda_n} = 0$, we obtain λ_n in the form of the weighted least square estimate

$$\lambda_n = -2\text{sign}(g_n^{(\text{LS})}) \mathbf{w}_n^T \mathbf{w}_n \mathbf{w}_n^T \mathbf{\Gamma}^{(n)} \boldsymbol{\eta} / \mathbf{w}_n^T \mathbf{\Gamma}^{(n)} \mathbf{w}_n, \quad (29)$$

where $\mathbf{\Gamma}^{(n)} = \text{diag}\{\gamma_n^2(1), \gamma_n^2(2), \dots, \gamma_n^2(N)\}$ and $\boldsymbol{\eta} = [\eta(1) \ \eta(2) \ \dots \ \eta(N)]^T \in \mathbb{R}^N$. Finally we calculate

$$\lambda_n^{\text{opt}} = \max \left\{ \min \left\{ 2|\mathbf{w}_n^T \mathbf{e}^{(n-1)}|, \lambda_n \right\}, \varepsilon \right\}, \quad (30)$$

in order to satisfy the constraint that $\varepsilon \leq \lambda_n^{\text{opt}} \leq 2|\mathbf{w}_n^T \mathbf{e}^{(n-1)}|$. For λ_n^{opt} obtained using (30), we consider the following two cases:

- (1) If $\lambda_n^{\text{opt}} = 2|\mathbf{w}_n^T \mathbf{e}^{(n-1)}|$, then $g_n^{(\text{lasso})} = 0$, and this candidate regressor will not be selected.
- (2) If $\varepsilon \leq \lambda_n^{\text{opt}} < 2|\mathbf{w}_n^T \mathbf{e}^{(n-1)}|$, then calculate $J(\lambda_n^{\text{opt}})$ based on (27) as the LOOMSE for this candidate regressor.

3.3. Moving unselectable regressors to the inactive set

From Section 3.2 we noted that a candidate regressor satisfying $2|\mathbf{w}_n^T \mathbf{e}^{(n-1)}| < \varepsilon$ does not need to be considered at the n th stage of selection. To save computational cost, we define the inactive set \mathcal{S} as the index set of the unselectable regressors removed from the pool of candidates.

In the n th OFR stage, all the candidate regressors in the candidate pool are made orthogonal to the previously selected $(n-1)$ regressors, and the candidate with the smallest LOOMSE value is selected as the n th model term \mathbf{w}_n . Denote any other candidate regressor as $\mathbf{w}^{(-)}$.

Main Results: If $\|\mathbf{w}^{(-)}\| \cdot \|\mathbf{e}^{(n-1)}\| < \frac{\varepsilon}{2}$, then this candidate regressor will never be selected in further regression stages, and hence it can be moved to \mathcal{S} .

Proof: At the $(n+1)$ th OFR stage, consider making the regressor $\mathbf{w}^{(-)}$ orthogonal to \mathbf{w}_n , and define

$$\mathbf{w}^{(+)} = \mathbf{w}^{(-)} - \frac{\mathbf{w}_n^T \mathbf{w}^{(-)}}{\mathbf{w}_n^T \mathbf{w}_n} \mathbf{w}_n. \quad (31)$$

Clearly,

$$\begin{aligned} \|\mathbf{w}^{(+)}\|^2 &= \left(\mathbf{w}^{(-)} - \frac{\mathbf{w}_n^T \mathbf{w}^{(-)}}{\mathbf{w}_n^T \mathbf{w}_n} \mathbf{w}_n \right)^T \left(\mathbf{w}^{(-)} - \frac{\mathbf{w}_n^T \mathbf{w}^{(-)}}{\mathbf{w}_n^T \mathbf{w}_n} \mathbf{w}_n \right) \\ &= \|\mathbf{w}^{(-)}\|^2 - \frac{(\mathbf{w}_n^T \mathbf{w}^{(-)})^2}{\mathbf{w}_n^T \mathbf{w}_n} \leq \|\mathbf{w}^{(-)}\|^2. \end{aligned} \quad (32)$$

The model residual vector after the selection of \mathbf{w}_n is

$$\mathbf{e}^{(n)} = \mathbf{e}^{(n-1)} - g_n^{(\text{lasso})} \mathbf{w}_n, \quad (33)$$

where $g_n^{(\text{olasso})}$ can be written as

$$g_n^{(\text{olasso})} = \left(\mathbf{w}_n^T \mathbf{e}^{(n-1)} - \frac{\lambda_n}{2} \text{sign}(g_n^{(\text{LS})}) \right) / \mathbf{w}_n^T \mathbf{w}_n. \quad (34)$$

Thus we have

$$\begin{aligned} \|\mathbf{e}^{(n)}\|^2 &= \|\mathbf{e}^{(n-1)}\|^2 - 2g_n^{(\text{olasso})} \mathbf{w}_n^T \mathbf{e}^{(n-1)} \\ &\quad + (g_n^{(\text{olasso})})^2 \mathbf{w}_n^T \mathbf{w}_n, \end{aligned} \quad (35)$$

$$\begin{aligned} (g_n^{(\text{olasso})})^2 \mathbf{w}_n^T \mathbf{w}_n &= \left((\mathbf{w}_n^T \mathbf{e}^{(n-1)})^2 - \right. \\ &\quad \left. \lambda_n \text{sign}(g_n^{(\text{LS})}) \mathbf{w}_n^T \mathbf{e}^{(n-1)} + \frac{\lambda_n^2}{4} \right) / \mathbf{w}_n^T \mathbf{w}_n, \end{aligned} \quad (36)$$

and

$$\begin{aligned} 2g_n^{(\text{olasso})} \mathbf{w}_n^T \mathbf{e}^{(n-1)} &= \left(2(\mathbf{w}_n^T \mathbf{e}^{(n-1)})^2 - \right. \\ &\quad \left. \lambda_n \text{sign}(g_n^{(\text{LS})}) \mathbf{w}_n^T \mathbf{e}^{(n-1)} \right) / \mathbf{w}_n^T \mathbf{w}_n. \end{aligned} \quad (37)$$

Substituting (36) and (37) into (35) yields

$$\begin{aligned} \|\mathbf{e}^{(n)}\|^2 &= \|\mathbf{e}^{(n-1)}\|^2 - \left((\mathbf{w}_n^T \mathbf{e}^{(n-1)})^2 - \frac{\lambda_n^2}{4} \right) / \mathbf{w}_n^T \mathbf{w}_n \\ &< \|\mathbf{e}^{(n-1)}\|^2, \end{aligned} \quad (38)$$

due to the fact that $|\mathbf{w}_n^T \mathbf{e}^{(n-1)}| > \frac{\lambda_n}{2}$. From (32) and (38), it can be concluded that

$$\|\mathbf{w}^{(+)}\| \cdot \|\mathbf{e}^{(n)}\| < \|\mathbf{w}^{(-)}\| \cdot \|\mathbf{e}^{(n-1)}\| < \frac{\varepsilon}{2}. \quad (39)$$

Since $\|\mathbf{w}^{(+)}\| \cdot \|\mathbf{e}^{(n)}\|$ is the upper bound of $\left| (\mathbf{w}^{(+)})^T \mathbf{e}^{(n)} \right|$, this means that this regressor will not be selected at the $(n+1)$ th stage. By induction, it will never be selected in further regression stages, and hence it can be moved to \mathcal{S} .

4. The proposed l^1 -POFR algorithm

The proposed l^1 -POFR algorithm integrates (i) the model regressor selection based on minimizing the LOOMSE; (ii) regularization parameter optimization also based on minimizing the LOOMSE; and (iii) the mechanism of removing unproductive candidate regressors during the OFR procedure. Define

$$\Phi^{(n-1)} = [\mathbf{w}_1 \cdots \mathbf{w}_{n-1} \ \phi_n^{(n-1)} \cdots \phi_M^{(n-1)}] \in \mathbb{R}^{N \times M}, \quad (40)$$

with $\Phi^{(0)} = \Phi_M$. If some of the columns in $\Phi^{(n-1)}$ have been interchanged, this will still be referred as $\Phi^{(n-1)}$ for notational simplicity.

Table 1. The n th stage of the selection procedure.

For $\{n \leq j \leq M\} \cap \{j \notin \mathcal{S}\}$, denote the k th element of $\phi_j^{(n-1)}$ as $\phi_j^{(n-1)}(k)$ and compute $\alpha_j = (\phi_j^{(n-1)})^T \mathbf{e}^{(n-1)}$, and $\beta_j = \|\phi_j^{(n-1)}\| \cdot \|\mathbf{e}^{(n-1)}\|$.

Step 1): If $\beta_j < \varepsilon/2$, $\mathcal{S} = \mathcal{S} \cup j$; Else if $|\alpha_j| < \varepsilon/2$, set $J_n^{(j)}$ as a very large positive number so that it will not be selected in Step 4). Otherwise goto step 2).

Step 2): Calculate

$$\kappa_n^{(j)} = (\phi_j^{(n-1)})^T \phi_j^{(n-1)}, \quad (41)$$

$$g_n^{(\text{LS},j)} = \frac{\alpha_j}{\kappa_n^{(j)}}, \quad (42)$$

$$\mathbf{\Gamma}^{(n,j)} = \text{diag} \left\{ \frac{1}{\left(\zeta^{(n-1)}(1) - (\phi_j^{(n-1)}(1))^2 / \kappa_n^{(j)} \right)^2}, \right. \\ \left. \frac{1}{\left(\zeta^{(n-1)}(2) - (\phi_j^{(n-1)}(2))^2 / \kappa_n^{(j)} \right)^2}, \dots, \right. \\ \left. \frac{1}{\left(\zeta^{(n-1)}(N) - (\phi_j^{(n-1)}(N))^2 / \kappa_n^{(j)} \right)^2} \right\} \in \mathbb{R}^{N \times N}, \quad (43)$$

$$\boldsymbol{\eta}^{(j)} = \mathbf{e}^{(n-1)} - g_n^{(\text{LS},j)} \phi_j^{(n-1)}, \quad (44)$$

$$\lambda_n^{(\text{opt},j)} = \max \left\{ \min \left\{ 2|\alpha_j|, -2\text{sign}(g_n^{(\text{LS},j)}) \kappa_n^{(j)} \right. \right. \\ \left. \left. (\phi_j^{(n-1)})^T \mathbf{\Gamma}^{(n,j)} \boldsymbol{\eta}^{(j)} / (\phi_j^{(n-1)})^T \mathbf{\Gamma}^{(n,j)} \phi_j^{(n-1)} \right\}, \varepsilon \right\}. \quad (45)$$

Step 3): If $\lambda_n^{(\text{opt},j)} = 2|\alpha_j|$, set $J_n^{(j)}$ as a very large positive number so that it will not be selected in Step 4); Otherwise calculate

$$g_n^{(\text{olasso},j)} = \left(|g_n^{(\text{LS},j)}| - \frac{\lambda_n^{(\text{opt},j)}/2}{\kappa_n^{(j)}} \right)_+ \text{sign}(g_n^{(\text{LS},j)}), \quad (46)$$

$$\mathbf{e}^{(n,j)} = \mathbf{e}^{(n-1)} - g_n^{(\text{olasso},j)} \phi_j^{(n-1)}, \quad (47)$$

$$J_n^{(j)} = (\mathbf{e}^{(n,j)})^T \mathbf{\Gamma}^{(n,j)} \mathbf{e}^{(n,j)} / N. \quad (48)$$

Step 4): Find

$$J_n = J_n^{(j_n)} = \min \left\{ J_n^{(j)}, \{l \leq j \leq M\} \cap \{j \notin \mathcal{S}\} \right\}. \quad (49)$$

Then update $\mathbf{e}^{(n)}$ and $g_n^{(\text{olasso})}$ as $\mathbf{e}^{(n,j_n)}$ and $g_n^{(\text{olasso},j_n)}$, respectively. The j_n th and the n th columns of $\Phi^{(n-1)}$ are interchanged, while the j_n th column and the n th column of \mathbf{A}_M are interchanged up to the $(n-1)$ th row. This effectively selects the n th regressor in the subset model. The modified Gram-Schmidt orthogonalisation procedure (Chen et al., 1989) then calculates the n th row of the matrix \mathbf{A}_M and transfers $\Phi^{(n-1)}$ into $\Phi^{(n)}$ as follows

$$\left. \begin{aligned} \mathbf{w}_n &= \phi_n^{(n-1)}, \\ a_{n,j} &= \mathbf{w}_n^T \phi_j^{(n-1)} / \mathbf{w}_n^T \mathbf{w}_n, \{n+1 \leq j \leq M\} \cap \{j \notin \mathcal{S}\}, \\ \phi_j^{(n)} &= \phi_j^{(n-1)} - a_{n,j} \mathbf{w}_n, \{n+1 \leq j \leq M\} \cap \{j \notin \mathcal{S}\}. \end{aligned} \right\} \quad (50)$$

Then update $\zeta^{(n)}(k) = \zeta^{(n-1)}(k) - (w_n(k))^2 / \mathbf{w}_n^T \mathbf{w}_n$ for $1 \leq k \leq N$.

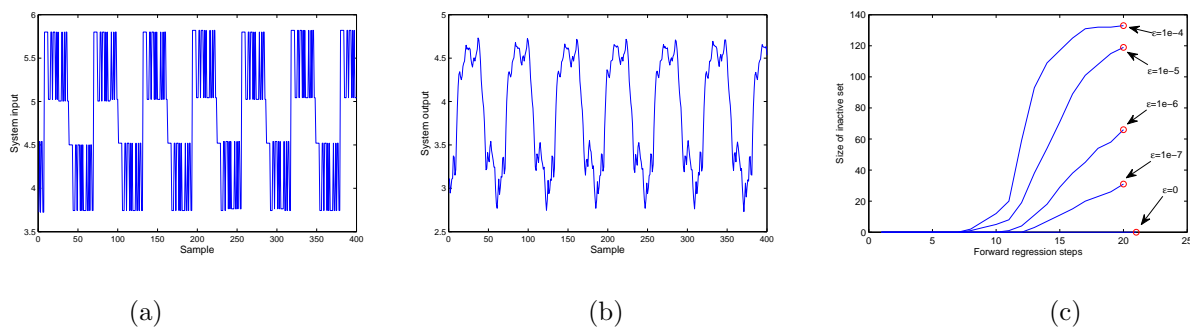


Figure 1. Engine Data: (a) the system input $u(k)$, (b) the system output $y(k)$, and (c) the evolution of the size of \mathcal{S} with respect to the chosen ε .

The initial conditions are as follows. Preset $\varepsilon > 0$ as a very small value. Set $\mathbf{e}^{(0)} = \mathbf{y}$, $\zeta^{(0)}(k) = 1$ for $1 \leq k \leq N$, and \mathcal{S} as the empty set \emptyset . The n th stage of the selection procedure is listed in Table 1. The OFR procedure is automatically terminated at the $(n_s + 1)$ th stage when the condition

$$J_{n_s+1} \geq J_{n_s} \quad (51)$$

is detected, yielding a subset model with n_s significant regressors. It is worth emphasizing that there always exists a model size n_s , and for $n \leq n_s$, the LOOMSE J_n decreases as n increases, while the condition (51) holds (Chen et al., 2004; Hong et al., 2003).

Note that the LOOMSE is used not only for deriving the closed form of the optimal regularization parameter estimate λ_n^{opt} but also for selecting the most significant model regressor. Specifically, a regressor is selected as the one that produces the smallest LOOMSE value as well as offering the reduction in the LOOMSE. After the n_s stage when there is no reduction in the LOOMSE criterion for a few consecutive OFR stages, the model construction procedure can be terminated. Thus, the l^1 -POFR algorithm automatically constructs a sparse n_s -term model, where typically $n_s \ll M$.

Also note that it is assumed that ε should not be too small such that the LOOMSE estimation formula can be considered to be accurate. This means that if ε is set too low, many insignificant candidate regressors will have inaccurate LOOMSE values for competition. However, we emphasize that these terms with inaccurate LOOMSE values will not be selected as the winner to enter the model. Hence in practice we only need to make sure that ε is not too large, which would introduce unnecessary bias to the model parameter estimates. Clearly, a relatively larger ε will save computational costs by 1) resulting in a sparser model, and 2) producing a larger sized inactive set during the OFR process.

Finally, regarding the computational complexity of the l^1 -POFR algorithm, if the unproductive regressors are not removed to the inactive set \mathcal{S} during the OFR procedure, it is well known that the computational cost is in the order of $O(N)$ for evaluating each candidate regressor (Chen et al., 2004). The total computational cost then needs to be scaled by the number of evaluations in forward regression, which is $M(M - n_s)/2$. By removing unproductive regressors to \mathcal{S} during the OFR procedure, the computational cost can obviously be reduced significantly. It is not possible to exactly assess the computational cost saving due to removing the unproductive regressors, as this is problem dependent.

5. Simulation Study

Example 1: This Engine Data set (Billings et al., 1989) contains the 410 data samples of the fuel rack position (the input $u(k)$) and the engine speed (the output $y(k)$), collected from a Leyland TL11

turbocharged, direct injection diesel engine which was operated at a low engine speed. The 410 input and output data points of the engine data set are plotted in Fig. 1 (a) and (b), respectively. The first 210 data samples were used in training and the last 200 data samples for model testing. The previous study has shown that the data set can be modeled adequately using the system input vector $\mathbf{x}(k) = [y(k-1) u(k-1) u(k-2)]^T$, and the best Gaussian RBF model was provided by the l^2 -norm local regularization assisted OLS (LROLS) algorithm based on the LOOMSE (LROLS-LOO) (Chen et al., 2004) which was quoted in Table 2 for comparison. The ε -SVM algorithm (Gun, 1998) and the LASSO were also experimented based on the Gaussian kernel with a common variance τ^2 . For the ε -SVM, the Matlab function *quadprog.m* was used with the algorithm option set as ‘interior-point-convex’. The tuning parameters in the ε -SVM algorithm, such as soft margin parameter C (Gun, 1998), were set empirically so that the best possible result was obtained after several trials. For the LASSO, the Matlab function *lasso.m* was used with 10-fold CV being used to select the associated regularization parameter. For both the ε -SVM and LASSO, we list the results obtained for a range of kernel width τ values in Table 2, for comparison.

Similar to the LROLS-LOO algorithm (Chen et al., 2004), we also used the Gaussian RBF kernel (3) for the proposed l^1 -POFR algorithm with an empirically set $\tau = 2.5$ and the RBF centers \mathbf{c}_i were formed using all the training data samples. With a preset value of ε , a sparse model of size n_s was automatically selected when the condition (51) was met. Fig. 1 (c) illustrates the evolution of the size of \mathcal{S} with respect to a range of the preset ε values. The test MSE values produced by the sparse models and the sizes of the models associated with the same range of ε values are recorded in Table 2, which show that the excellent model generalization capability of all the models generated by the proposed algorithm. Moreover, the l^1 -POFR algorithm produces the sparsest model.

Example 2: This regression benchmark data set, Boston Housing Data, is available at the UCI repository (Frank and Asuncion, 2010). The data set comprises 506 data points with 14 variables. The previous study (Chen et al., 2009) performed the task of predicting the median house value from the remaining 13 attributes using the ε -SVM (Gun, 1998), the LROLS-LOO (Chen et al., 2004) and the nonlinear OFR based on the LOOMSE (NonOFR-LOO) (Chen et al., 2009). The NonOFR-LOO algorithm (Chen et al., 2009) constructs a *nonlinear* RBF model in the OFR procedure,

Table 2. Comparison of the modeling performance for Engine Data. The computational cost saving is based on the same size of model without removing unproductive regressors in the l^1 -POFR.

Algorithm	MSE training set	MSE test set	Model size	Cost saving
LROLS-LOO (Chen et al., 2004)	0.000453	0.000490	22	NA
ε -SVM ($\tau = 3$)	0.000502	0.000482	208	NA
ε -SVM ($\tau = 2.5$)	0.000480	0.000475	208	NA
ε -SVM ($\tau = 2$)	0.000461	0.000486	208	NA
ε -SVM ($\tau = 1.5$)	0.000415	0.000579	208	NA
ε -SVM ($\tau = 1$)	0.000370	0.000794	208	NA
LASSO ($\tau = 1.5$)	0.000923	0.001010	70	NA
LASSO ($\tau = 1$)	0.000708	0.000748	44	NA
LASSO ($\tau = 0.5$)	0.000706	0.000842	54	NA
LASSO ($\tau = 0.2$)	0.000565	0.000800	81	NA
LASSO ($\tau = 0.1$)	0.000644	0.001907	76	NA
l^1 -POFR ($\varepsilon = 10^{-4}$)	0.000498	0.000502	20	27%
l^1 -POFR ($\varepsilon = 10^{-5}$)	0.000492	0.000480	20	18%
l^1 -POFR ($\varepsilon = 10^{-6}$)	0.000484	0.000485	20	8%
l^1 -POFR ($\varepsilon = 10^{-7}$)	0.000481	0.000476	20	3%
l^1 -POFR ($\varepsilon = 0$)	0.000452	0.000472	21	0%

where each stage of the OFR determines one RBF node's center vector and diagonal covariance matrix by minimizing the LOOMSE. In the experiment study presented in (Chen et al., 2009), 456 data points were randomly selected from the data set for training and the remaining 50 data points were used to form the test set. Average results were given over 100 realizations. For each realization, 13 input attributes were normalized so that each attribute had zero mean and standard deviation of one. We also experimented with the LASSO supplied by Matlab *lasso.m* with option set as 10-fold CV to select the associated regularization parameter. For the LASSO, a common kernel width τ was set for constructing the kernel model from the 456 candidate regressors of each realization, and a range of τ values were experimented.

For the l^1 -POFR, $\tau = 15$ was empirically set for constructing 456 candidate Gaussian RBF regressors of each realization. We experimented a range of the preset ε values for the l^1 -POFR algorithm, and the results obtained are as summarized in Table 3, in comparison with the results obtained by the ε -SVM and the LASSO, as well as the LROLS-LOO and NonOFR-LOO, which are quoted from the study (Chen et al., 2009).

6. Conclusions

We have developed an efficient data model algorithm, referred to as the l^1 -norm penalized orthogonal forward regression (l^1 -POFR), for linear-in-the-parameters nonlinear models based on a new l^1 -norm penalized cost function defined in the constructed orthogonal modeling space. The LOOMSE is used for simultaneous model term selection and regularization parameter estimation in a highly efficient OFR procedure. Additionally, we have proposed a lower bound of the regularisation parameters for robust LOOMSE estimation as well as detecting and removing insignificant regressors to an inactive set along the OFR process, further enhancing the efficiency of the OFR procedure. Numerical studies have been utilized to demonstrate the effectiveness of this new l^1 -POFR approach.

References

- Billings, S. A., Chen, S., and Backhouse, R. (1989). The identification of linear and nonlinear models of a turbocharged automotive diesel engine. *Mechanical Systems and Signal Processings*, 3(2):123–142.
- Chen, S. and Billings, S. A. (1989). Representation of nonlinear systems: The NARMAX model,. *International Journal of Control*, 49(3):1013–1032.

Table 3. Comparison of the modeling performance for Boston House Data. The results were averaged over 100 realizations and given as mean \pm standard deviation.

Algorithm	MSE training set	MSE test set	Model size
ε -SVM (Gun, 1998)	6.80 ± 0.44	23.18 ± 9.05	243 ± 5.3
LROLS-LOO (Chen et al., 2004)	12.97 ± 2.67	17.42 ± 4.67	58.6 ± 11.3
NonOFR-LOO (Chen et al., 2009)	10.10 ± 3.40	14.07 ± 3.62	34.6 ± 8.4
LASSO ($\tau = 2$)	8.52 ± 3.57	14.37 ± 8.15	76.8 ± 39.7
LASSO ($\tau = 3$)	8.55 ± 1.07	13.31 ± 6.65	68.6 ± 29.3
LASSO ($\tau = 5$)	10.45 ± 1.07	15.05 ± 8.37	85.9 ± 19.7
LASSO ($\tau = 10$)	16.42 ± 1.78	19.39 ± 8.31	29.9 ± 21.3
l^1 -POFR ($\varepsilon = 0.01$)	9.99 ± 1.37	14.47 ± 7.47	30.5 ± 5.3
l^1 -POFR ($\varepsilon = 0.001$)	9.24 ± 1.57	14.10 ± 7.02	34.9 ± 7.8
l^1 -POFR ($\varepsilon = 0.0001$)	9.07 ± 1.64	14.02 ± 6.85	36.6 ± 9.3
l^1 -POFR ($\varepsilon = 0.00001$)	9.08 ± 1.64	13.95 ± 6.76	36.5 ± 9.3

- Chen, S., Billings, S. A., and Luo, W. (1989). Orthogonal least squares methods and their applications to non-linear system identification. *International Journal of Control*, 50(5):1873–1896.
- Chen, S., Hong, X., and Harris, C. J. (2003). Sparse kernel regression modelling using combined locally regularised orthogonal least squares and D-optimality experimental design. *IEEE Trans. on Automatic Control*, 48(6):1029–1036.
- Chen, S., Hong, X., and Harris, C. J. (2009). Construction of tunable radial basis function networks using orthogonal forward selection. *IEEE Trans. on Systems, Man and Cybernetics, Part B: Cybernetics*, 39(2):457–466.
- Chen, S., Hong, X., Harris, C. J., and Sharkey, P. M. (2004). Sparse modelling using orthogonal forward regression with PRESS statistic and regularization. *IEEE Trans. on Systems, Man and Cybernetics, Part B: Cybernetics*, 34(2):898–911.
- Chen, S. S., Donoho, D. L., and Saunders, M. A. (1998). Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20(1):33–61.
- Efron, B., Johnstone, I., Hastie, T., and Tibshirani, R. (2004). Least angle regression. *Annals of Statistics*, 32:407–451.
- Frank, A. and Asuncion, A. (2010). UCI machine learning repository.
- Gun, S. R. (1998). Support vector machines for classification and regression. ISIS Res. Group, Dept. Electron. Comput. Sci., Univ. Southampton, Southampton, U.K.
- Hong, X., Sharkey, P. M., and Warwick, K. (2003). Automatic nonlinear predictive model construction using forward regression and the PRESS statistic. *IEE Proc.-Control Theory Applications*, 150(3):245–254.
- Ji, S., Xue, Y., and Carin, L. (2008) Bayesian compressive sensing. *IEEE Trans. Signal Processing*, 56(6):2346–2356.
- MacKay, D. J. C. (1991). *Bayesian Methods for Adaptive Models*. PhD thesis, California Institute of Technology, USA.
- Orr, M. J. L. (1995). Regularisation in the selection of radial basis function centers. *Neural Computation*, 7(3):954–975.
- Rao, R. B., Fung, G., and Rosales, R. (2008). On the dangers of cross-validation. An experimental evaluation. *Proc. 2008 SIAM Conf. Data Mining* (Atlanta, GA), Apr. 24–26, pp. 588–596.
- Sherman, J. and Morrison, W. J. Adjustment of an inverse matrix corresponding to a change in one element of a given matrix. *Annals of Mathematical Statistics*, 21(1):124–127.
- Stone, M. (1974). Cross validatory choice and assessment of statistical predictions. *Journal of the Royal Statistical Society, Series B*, 36(2):111–147.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of Royal Statistical Society, Series B*, 58(1):267–288.