

1-3-2011

Efficient techniques for simultaneous variable selection and sensor selection via convex selection inducing penalties

Ernest Fokoue

Follow this and additional works at: <http://scholarworks.rit.edu/article>

Recommended Citation

Fokoue, Ernest, "Efficient techniques for simultaneous variable selection and sensor selection via convex selection inducing penalties" (2011). Accessed from <http://scholarworks.rit.edu/article/137>

This Article is brought to you for free and open access by RIT Scholar Works. It has been accepted for inclusion in Articles by an authorized administrator of RIT Scholar Works. For more information, please contact ritscholarworks@rit.edu.

EFFICIENT TECHNIQUES FOR SIMULTANEOUS VARIABLE SELECTION AND SENSOR SELECTION VIA CONVEX SELECTION INDUCING PENALTIES

Ernest Fokoué

Center for Quality and Applied Statistics
Rochester Institute of Technology
Rochester, NY 14624, USA
eMail: epfeqa@rit.edu

January 3, 2011

Abstract

This paper extends results from the traditional D-optimality machinery to derive an efficient technique for simultaneous variable selection and sensor selection. An important advantage of the proposed technique is the convexity of the formulated optimization task along with a byproduct of straightforward sparsity. The theoretical foundation of the proposed method is explored at great length, and a variety of examples are provided to demonstrate the effectiveness of our technique. Comparisons with existing techniques are offered that provide evidence as to the superiority of our technique on a variety of indicators.

Keywords: *Variable Selection, Convex optimization, Bayesian analysis, Optimal Experimental Design, Sensor selection, Sparsity, D-Optimality.*

1 Introduction

Let $\mathbf{x}_j^\top \equiv (x_{j1}, x_{j2}, \dots, x_{jp})$ denote a p -dimensional vector of some observable characteristics of interest. Consider a p -dimensional vector $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_p)^\top$ of regression coefficients, then assume that a response (measurement) Y_j of interest at point \mathbf{x}_j can be written as

$$Y_j = \mathbf{x}_j^\top \boldsymbol{\beta} + \epsilon_j, \quad j = 1, \dots, n.$$

Throughout this paper, we shall assume that the ϵ_j 's are i.i.d $\mathbf{N}(0, \sigma^2)$. Note also that, for simplicity, we have restricted ourselves to a model that passes through the origin. Under this homoscedastic noise model, the maximum likelihood estimator $\hat{\boldsymbol{\beta}}_{\text{MLE}}$ of $\boldsymbol{\beta}$ is such that

$$\hat{\boldsymbol{\beta}}_{\text{MLE}} = \left(\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \right)^{-1} \sum_{j=1}^n y_j \mathbf{x}_j \quad \text{and} \quad \text{cov}(\hat{\boldsymbol{\beta}}_{\text{MLE}}) = \sigma^2 \left(\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \right)^{-1}$$

In traditional optimal experimental design, one has a set of n potential points of measurement or sensors, and the goal is to choose those k sensors or points of measurement that yield the "best" estimation of $\boldsymbol{\beta}$. For instance, with $\hat{\boldsymbol{\beta}}_{\text{MLE}}$ being an unbiased estimator, a reasonable criterion for measuring the goodness of $\hat{\boldsymbol{\beta}}_{\text{MLE}}$ will naturally be based on its covariance matrix. In fact, we will see later that all the three criteria used for measuring the optimality of the design will be based on

functions of the covariance matrix of $\hat{\beta}_{\text{MLE}}$. The problem in optimal experimental design is then two-fold: (i) Which $k \ll n$ sensors or points of measurement to choose out of the n possible ones; and (ii) How many times can each chosen sensor be used, while making the total number of uses at most equal to k . One of the most commonly used optimality criteria is the so-called D-optimality that seeks to choose those points that minimize the determinant of the covariance matrix of $\hat{\beta}_{\text{MLE}}$. In other words, if each π_j , $j = 1, 2, \dots, n$ represents the frequency of use of measurement point j , then a k -point D-optimal design is obtained as a solution to the relaxed sensor selection convex optimization problem

$$\begin{aligned} \text{Maximize} \quad & \log \det \left(\sum_{j=1}^n \pi_j \mathbf{x}_j \mathbf{x}_j^\top \right) \\ \text{Subject to} \quad & 0 \leq \pi_j \leq 1, \quad j = 1, \dots, n \quad \text{and} \quad \sum_{j=1}^n \pi_j = k \end{aligned} \quad (1)$$

The k -point D-optimal design is therefore the subset $\xi = \{i_1, i_2, \dots, i_k\} \subseteq \{1, 2, \dots, n\}$ that corresponds to set of sensors or measurements with the k largest values of π_j . (Joshi and Boyd 2009) proposes an approximate solution obtained by making the constraint $\pi_j \in (0, 1)$ implicit in the objective function so that the resulting convex optimization problem is

$$\begin{aligned} \text{Maximize} \quad & \log \det \left(\sum_{j=1}^n \pi_j \mathbf{x}_j \mathbf{x}_j^\top \right) + \kappa \left[\sum_{j=1}^n \log(\pi_j) + \sum_{j=1}^n \log(1 - \pi_j) \right] \\ \text{Subject to} \quad & \sum_{j=1}^n \pi_j = k. \end{aligned} \quad (2)$$

Let $w_k \in (0, 1)$ denote the importance (relevance) of variable x_k . Consider a diagonal matrix

$$\mathbf{W} = \begin{bmatrix} w_1 & & & \\ & w_2 & & \\ & & \ddots & \\ & & & w_p \end{bmatrix} = \text{diag}(w_1, w_2, \dots, w_p)$$

Instead of using the input vector (x_1, x_2, \dots, x_p) , consider using $\tilde{\mathbf{x}}_j = \mathbf{W} \mathbf{x}_j$, so that the response Y_j is now $Y_j = \mathbf{x}_j^\top \mathbf{W}^\top \beta + \epsilon_j$, $j = 1, \dots, n$. Now, let

$$\text{Pen}(\mathbf{W}) = \left[\sum_{k=1}^p \log(w_k) + \sum_{k=1}^p \log(1 - w_k) \right] \quad \text{and} \quad \mathbf{X}^\top \mathbf{X} = \left[\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \right]$$

Variable selection can be achieved by solving the problem formulated in equation (4).

$$\mathcal{E}(\mathbf{W}) = \log \det \left(\mathbf{W}^\top \mathbf{X}^\top \mathbf{X} \mathbf{W} \right) + \nu \text{Pen}(\mathbf{W}) \quad (3)$$

The corresponding optimization problem is

$$\left[\begin{array}{ll} \text{Maximize} & \mathcal{E}(\mathbf{W}) \\ \text{Subject to} & 0 \leq w_k \leq 1, \quad k = 1, \dots, p \\ & \sum_{k=1}^p w_k = \omega. \end{array} \right. \quad (4)$$

where ω is a fraction of p indicating the extend of parsimony desired. It is worth noting that that second part of the constraint need not be specified in practice during the optimization procedure. Typically,

one may want to use w_k 's that are indicator variables, so as to perform straightforward selection. However, it turns out that using real numbers between 0 and 1 has many advantages as discussed by (Joshi and Boyd 2009) and (Fokoue and Goel 2009). For instance, values of ν less than 0.5 force the resulting w_k 's to be close to 0 or 1, thereby providing a strong determination of the importance (relevance) of the corresponding variable. Another important advantage is both computational and numerical: to see, let's assume that $p = 2$ and $n = 2$. Then

$$g(\mathbf{W}) = \mathbf{W}^\top \mathbf{X}^\top \mathbf{X} \mathbf{W} = \begin{bmatrix} w_1^2 \sum_{i=1}^3 x_{1i} x_{i1} & w_1 w_2 \sum_{i=1}^3 x_{i1} x_{i2} \\ w_1 w_2 \sum_{i=1}^3 x_{i2} x_{i1} & w_2^2 \sum_{i=1}^3 x_{i2} x_{i2} \end{bmatrix}$$

Throughout the optimization procedure, it is crucial that this matrix remain well conditioned. Now, if the w_k 's are binary indicators, then the matrix cannot stay full rank. For instance, if $w_2 = 0$, then the matrix collapses to a scalar and the whole computation cannot continue. However, with $w_k \in (0, 1)$, the matrix remains well conditioned throughout. This advantage is crucial to obtaining the desired variable selection solution, but also makes the procedure doable throughout traditional methods like Newton's method.

Theorem 1 *If $A \in \mathbb{R}^{p \times p}$ and $B \in \mathbb{R}^{p \times p}$ are two symmetric matrices, then*

$$\frac{\partial \log \det A^\top B A}{\partial A} = \frac{\partial \log \det A^\top B A}{\partial A^\top B A} \frac{\partial A^\top B A}{\partial A} = (A^\top B A)^{-\top} \frac{\partial A^\top B A}{\partial A} = 2(A^\top B A)^{-\top} B A.$$

Therefore,

$$\frac{\partial g(\mathbf{W})}{\partial \mathbf{W}} = 2(\mathbf{W}^\top \mathbf{X}^\top \mathbf{X} \mathbf{W})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{W}$$

Let $\mathbf{D} = \mathbf{X}^\top \mathbf{X}$. Then, we need to derive

$$\begin{aligned} g''(\mathbf{W}) &= \frac{\partial(\mathbf{W}^\top \mathbf{D} \mathbf{W})^{-1} \mathbf{D} \mathbf{W}}{\partial \mathbf{W}} \\ &= \frac{\partial(\mathbf{W}^\top \mathbf{D} \mathbf{W})^{-1}}{\partial \mathbf{W}} \mathbf{D} \mathbf{W} + (\mathbf{W}^\top \mathbf{D} \mathbf{W})^{-1} \frac{\partial \mathbf{D} \mathbf{W}}{\partial \mathbf{W}} \\ &= \frac{\partial(\mathbf{W}^\top \mathbf{D} \mathbf{W})^{-1}}{\partial \mathbf{W}^\top \mathbf{D} \mathbf{W}} \frac{\partial \mathbf{W}^\top \mathbf{D} \mathbf{W}}{\partial \mathbf{W}} \mathbf{D} \mathbf{W} + (\mathbf{W}^\top \mathbf{D} \mathbf{W})^{-1} \mathbf{D} \\ &= -(\mathbf{W}^\top \mathbf{D} \mathbf{W})^{-\top} \otimes (\mathbf{W}^\top \mathbf{D} \mathbf{W})^{-1} 2 \mathbf{D} \mathbf{W} \mathbf{D} \mathbf{W} + (\mathbf{W}^\top \mathbf{D} \mathbf{W})^{-1} \mathbf{D} \end{aligned}$$

Perhaps the most importance advantage of our proposed scheme lies in the fact that we have a convex optimization problem, with the guarantee of a unique solution.

$$\hat{\mathbf{W}} = \arg \max_{\mathbf{W}} \mathcal{E}(\mathbf{W})$$

This is even a better advantage because once could even think of simultaneously achieving variable selection and the corresponding D-optimal design.

$$\mathcal{E}(\mathbf{W}, \mathbf{Z}) = \log \det \left(\mathbf{W}^\top \mathbf{X}^\top \mathbf{Z} \mathbf{X} \mathbf{W} \right) + \nu \text{Pen}(\mathbf{W}) + \kappa \text{Pen}(\mathbf{Z}) \quad (5)$$

where

$$\text{Pen}(\mathbf{Z}) = \left[\sum_{j=1}^n \log(z_j) + \sum_{j=1}^n \log(1 - z_j) \right] \quad (6)$$

with $\mathbf{Z} = \text{diag}(z_1, z_2, \dots, z_n)$ is the diagonal matrix whose entries act on the sensors. The convex optimization problem to solve in order to achieve that is

$$\left[\begin{array}{l} \text{Maximize} \quad \mathcal{E}(\mathbf{W}, \mathbf{Z}) \\ \text{Subject to} \quad 0 \leq w_k \leq 1, \quad k = 1, \dots, p \quad \text{and} \quad \sum_{k=1}^p w_k = \omega. \\ \quad \quad \quad 0 \leq z_j \leq 1, \quad j = 1, \dots, n \quad \text{and} \quad \sum_{j=1}^n z_j = \zeta. \end{array} \right. \quad (7)$$

The solution

$$(\hat{\mathbf{Z}}) = \arg \max_{\mathbf{W}, \mathbf{Z}} \mathcal{E}(\mathbf{W}, \mathbf{Z})$$

provides a total $n + p$ estimates, all in $(0, 1)$, with the \hat{w}_k 's allowing variable selection while the \hat{z}_j 's allow sensor selection. It is worth emphasizing once again that we have the great advantage of convex optimization.

Theorem 2 *If $A \in \mathbb{R}^{p \times p}$ and $B \in \mathbb{R}^{p \times p}$ are two symmetric matrices, then*

$$\frac{\partial \log \det A^\top B A}{\partial B} = \frac{\partial \log \det A^\top B A}{\partial A^\top B A} \frac{\partial A^\top B A}{\partial B} = (A^\top B A)^{-\top} \frac{\partial A^\top B A}{\partial B} = A(A^\top B A)^{-\top} A^\top.$$

From this proposition, the derivative of $h(\mathbf{Z}) = \log \det (\mathbf{W}^\top \mathbf{X}^\top \mathbf{Z} \mathbf{X} \mathbf{W})$ is given by

$$\frac{\partial h(\mathbf{Z})}{\partial \mathbf{Z}} = \mathbf{X} \mathbf{W} (\mathbf{W}^\top \mathbf{X}^\top \mathbf{Z} \mathbf{X} \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X}^\top.$$

Obviously, the case $\mathbf{W} = \mathbf{I}_p$ corresponds to the traditional D-optimality.

1.1 Traditional approximate D-optimality

Let $\mathbf{Z} = \text{diag}(z_1, z_2, \dots, z_n)$ be the diagonal matrix whose entries act on the sensors.

$$\mathcal{E}(\mathbf{Z}) = \log \det (\tau \mathbf{X}^\top \mathbf{Z} \mathbf{X} + \mathbf{Q}^{-1}) + \kappa \left[\sum_{j=1}^n \log(z_j) + \sum_{j=1}^n \log(1 - z_j) \right] \quad (8)$$

The convex optimization problem to solve in order to achieve that is

$$\left[\begin{array}{l} \text{Maximize} \quad \mathcal{E}(\mathbf{Z}) \\ \text{Subject to} \quad 0 \leq w_k \leq 1, \quad k = 1, \dots, p \quad \text{and} \quad \sum_{k=1}^p w_k = \omega. \\ \quad \quad \quad 0 \leq z_j \leq 1, \quad j = 1, \dots, n \quad \text{and} \quad \sum_{j=1}^n z_j = \zeta. \end{array} \right. \quad (9)$$

The solution

$$(\hat{\mathbf{W}}, \hat{\mathbf{Z}}) = \arg \max_{\mathbf{W}, \mathbf{Z}} \mathcal{E}(\mathbf{W}, \mathbf{Z})$$

provides a total $n + p$ estimates, all in $(0, 1)$, with the \hat{w}_k 's allowing variable selection while the \hat{z}_j 's allow sensor selection. It is worth emphasizing once again that we have the great advantage of convex optimization.

Theorem 3 If $A \in \mathbb{R}^{p \times p}$ and $B \in \mathbb{R}^{p \times p}$ are two symmetric matrices, then

$$\frac{\partial \log \det A^\top B A}{\partial B} = \frac{\partial \log \det A^\top B A}{\partial A^\top B A} \frac{\partial A^\top B A}{\partial B} = (A^\top B A)^{-\top} \frac{\partial A^\top B A}{\partial B} = A(A^\top B A)^{-\top} A^\top.$$

From this proposition, the derivative of $h(\mathbf{Z}) = \log \det (\mathbf{W}^\top \mathbf{X}^\top \mathbf{Z} \mathbf{X} \mathbf{W})$ is given by

$$\frac{\partial h(\mathbf{Z})}{\partial \mathbf{Z}} = \mathbf{X} \mathbf{W} (\mathbf{W}^\top \mathbf{X}^\top \mathbf{Z} \mathbf{X} \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X}^\top.$$

Obviously, the case $\mathbf{W} = \mathbf{I}_p$ corresponds to the traditional D-optimality.

Let $\mathbf{Z} = \text{diag}(z_1, z_2, \dots, z_n)$ denote the diagonal matrix of sensor pointers (weights).

$$f(\mathbf{Z}) = -\log(\det(\tau \mathbf{X}^\top \mathbf{Z} \mathbf{X} + \mathbf{Q}^{-1})) - \kappa \sum_{j=1}^n \{\log z_j + \log(1 - z_j)\}$$

$$\mathbf{W} = [\mathbf{X}^\top \mathbf{Z} \mathbf{X}]^{-1} \quad \text{and} \quad \mathbf{V} = \mathbf{X} \mathbf{W} \mathbf{X}^\top$$

Clearly,

$$\begin{aligned} \frac{\partial \log(\det(\tau \mathbf{X}^\top \mathbf{Z} \mathbf{X} + \mathbf{Q}^{-1}))}{\partial \mathbf{Z}} &= \frac{\partial \log(\det(\tau \mathbf{X}^\top \mathbf{Z} \mathbf{X} + \mathbf{Q}^{-1}))}{\partial \mathbf{X}^\top \mathbf{Z} \mathbf{X}} \frac{\partial \mathbf{X}^\top \mathbf{Z} \mathbf{X}}{\partial \mathbf{Z}} \\ &= -\text{vec}((\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-\top})^\top \mathbf{X}^\top \otimes \mathbf{X}^\top \\ &= \mathbf{X}(\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-1} \mathbf{X}^\top \end{aligned}$$

As a result, the Jacobian would then be

$$\nabla f = -\text{diag}(\mathbf{V}) - \kappa \left(\frac{1}{z_j} - \frac{1}{1 - z_j} \right)$$

Now,

$$\begin{aligned} \frac{\partial \mathbf{X}(\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-1} \mathbf{X}^\top}{\partial \mathbf{Z}} &= \frac{\partial \mathbf{X}(\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-1} \mathbf{X}^\top}{\partial (\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-1}} \frac{\partial (\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-1}}{\partial \mathbf{Z}} \\ &= \mathbf{X} \otimes \mathbf{X} \frac{\partial (\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-1}}{\partial \mathbf{X}^\top \mathbf{Z} \mathbf{X}} \frac{\partial \mathbf{X}^\top \mathbf{Z} \mathbf{X}}{\partial \mathbf{Z}} \\ &= \mathbf{X} \otimes \mathbf{X} (\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-\top} \otimes (\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-1} \mathbf{X}^\top \otimes \mathbf{X}^\top \\ &= [\mathbf{X}(\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-1} \mathbf{X}^\top] \otimes [\mathbf{X}(\mathbf{X}^\top \mathbf{Z} \mathbf{X})^{-1} \mathbf{X}^\top] \end{aligned}$$

Therefore, the Hessian in this case is given by

$$\mathbf{H} = \nabla \nabla f = \mathbf{V} \otimes \mathbf{V} + \kappa \text{diag} \left(\frac{1}{z_j^2} + \frac{1}{(1 - z_j)^2} \right)$$

2 Numerical demonstrations and simulations

Example 1: In order to gain insights into the similarities and the differences between D-optimal support points and relevant vectors, we first consider a simple univariate function

$$f(\mathbf{x}) = -\mathbf{x} + \sqrt{2} \sin \left(\pi^{3/2} \mathbf{x}^2 \right) \quad \text{with} \quad \mathbf{x} \in [-1, +1].$$

With this, our data consists of pairs $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, where the \mathbf{x}_i 's are equally spaced points in $[-1, +1]$. From a traditional D-optimal design standpoint, we need to specify a model in order to

form the data matrix. A natural candidate in this case is the polynomial regression model. A quick snoop at the scatterplot suggests that an 8th polynomial could capture the underlying function, i.e.,

$$Y_j = \beta_0 + \beta_1 \mathbf{x}_j + \beta_2 \mathbf{x}_j^2 + \cdots + \beta_8 \mathbf{x}_j^8 + \epsilon_j.$$

For the relevance vector machine, we used the gaussian radial basis function kernel, and found the bandwidth of $r = 0.5$ to be adequate for this data.

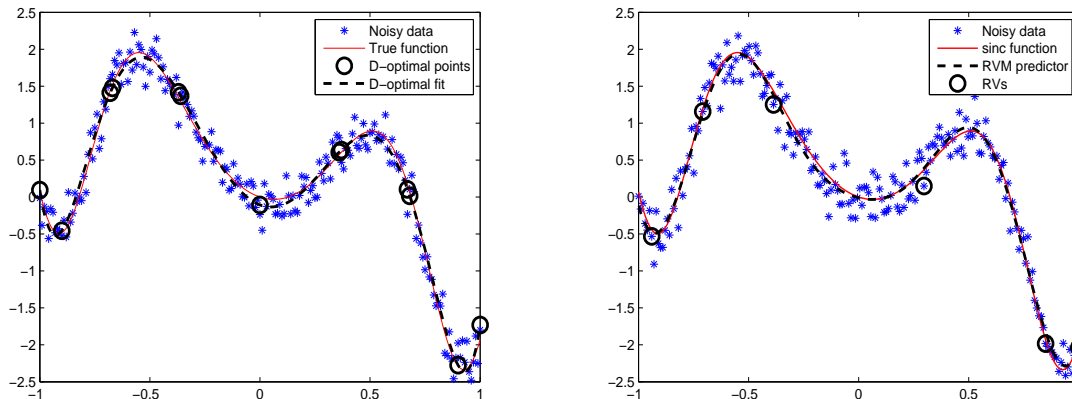


Figure 1: (left) D-optimal support points; (right) RVM relevant vectors.

For simplicity, the noise variance σ^2 is assumed known and fixed at 0.2^2 . As far as the similarities go, most of the points are identical for both methods. Regarding the differences, the relevance vector machine yields fewer points, for the obvious reason that it applies an extra constraint driven by the response values and therefore achieves more tuning. Besides, it is important to recall that the strongest motivation behind RVM is sparsity (fewer relevant vectors), while D-optimality sets out to find a k -point design. The number is fixed in one case, while the minimum number is sought in the other.

Example 2: As our second example, we take a look at the commonly used sinc function

$$f(\mathbf{x}) = \frac{\sin 10\mathbf{x}}{10\mathbf{x}} \quad \text{with } \mathbf{x} \in [-1, +1].$$

For this example, our noise variance is still $\sigma^2 = 0.2^2$, but our response variable is now expressed as a weighted sum of *Legendre* or *Chebyshev* orthogonal polynomials to which we add the homoscedastic gaussian noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$ as before. Figure (2) shows the results obtained from both the D-optimality criterion (left) and the Relevance vector machine approach (right). Again, while it is obvious that the two methods are looking for the points that most affect the variance of the estimates of the parameters, it seems clear that RVM retains fewer points than D-optimality. The reason is that the results presented here are obtained using the generic D-optimality criterion of equation (??). We solved this using *CVX*, a package for specifying and solving convex programs (Joshi and Boyd 2009). Once the D-optimality criterion is enriched with the selection inducing *Beta* as in equation (??), a more sparse solution should be expected. Also, the complete reformulation of equation should produce results that are fairly identical to the output from the Relevance Vector Machine of (Tipping 2001).

3 Conclusion, discussion and future work

We have shown in this paper that the statistical problem underlying the now very popular Relevance Vector machine can essentially be formulated as an adaptive D-optimal design problem. The formulation

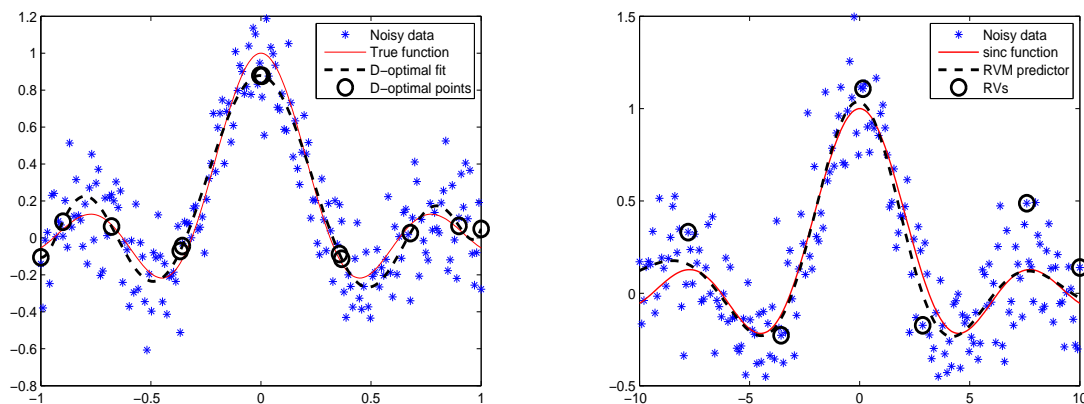


Figure 2: (left) D-optimal support points; (right) RVM relevant vectors.

derived in this paper provides a crucial advantage in that the problem is now a convex optimization task with the guarantee of a unique solution, as opposed to original RVM that is known not to yield a unique solution. Our immediate future work is to numerically implement the new formulation and also use our derived scheme on real life problems. Another aspect worth exploring is the reconstruction of the primal problem corresponding to the dual definition of the RVM. Much later, we hope to investigate the theoretical aspects of this connection a little further, and also consider exploring how this affects Relevance Vector Classification.

References

- Fokoue, E. and P. Goel (2009, December). An Optimal Experimental Design Perspective on the Relevance Vector Machine. Technical Report EPF-09-7-1, Rochester Institute of Technology, Rochester, New York, Center for Quality and Applied Statistics, 98 Lomb Memorial Drive, Rochester, NY 14623.
- Joshi, S. and S. Boyd (2009). Sensor selection via convex optimization. *IEEE Trans. Signal Processing* 57(2), 451–462.
- Tipping, M. E. (2001). Sparse bayesian learning and the relevance vector machine. *J. Mach. Learning Res.* 1, 211–244.