

Incremental Sparse Kernel Machine

Masa-aki Sato¹ and Shigeeyuki Oba²

¹ ATR, Human Information Science Laboratories, and CREST, JST,
masa-aki@atr.co.jp

² Nara Institute of Science and Technology, shige-o@is.aist-nara.ac.jp

Abstract. The Relevance Vector Machine (RVM) gives a probabilistic model for a sparse kernel representation. It achieves comparable performance to the Support Vector Machine (SVM) while using substantially fewer kernel bases. However, the computational complexity of the RVM in the training phase prohibits its application to large datasets. In order to overcome this difficulty, we propose an incremental Bayesian method for the RVM. The preliminary experiments showed the efficiency of our method for large datasets.

1 Introduction

The SVM [3] has been recognized as a state-of-the-art method for machine learning. It makes predictions based on a linear combination of kernel functions defined on a subset of the training data points called the support vectors. This sparse kernel representation avoids overfitting to the training data and increases the generalization ability. However, the main disadvantage of the SVM is its lack of probabilistic model. Therefore, it cannot estimate the uncertainty of a prediction, which gives important information in application to real-world problems. Recently, a probabilistic model for a kernel representation called the Relevance Vector Machine (RVM) [1] was proposed. The Bayes estimation of the RVM with an Automatic Relevance Determination (ARD) prior [1, 2] leads to a sparse kernel representation. The RVM achieves comparable performance to the SVM while using substantially fewer kernel bases [1]. However, the principal disadvantage of the RVM is its computational complexity in the training phase. It requires $O((\# \text{ of data})^3)$ computation time. This prohibits the application of the RVM to large datasets.

In order to overcome this difficulty, we propose an incremental Bayesian method for the RVM. The proposed method can learn a large amount of data incrementally and gives a sparse kernel representation. It requires $O((\# \text{ of bases})^2 * (\# \text{ of data}))$ computation time, which is much smaller than the computation time required by the RVM for a large amount of data.

The Gaussian Process (GP) [4] also provides a probabilistic model for kernel representation. An online learning method for the GP was proposed [5] recently. However, the GP has no inherent mechanism for basis selection. In order to get a sparse kernel representation, it is necessary to introduce heuristic mechanisms [5]. The advantage of our method is that the sparse kernel representation is an

inherent property of the Bayes estimation with the ARD prior. The preliminary experiments showed the efficiency of our approach for large datasets.

2 Relevance Vector Machine

In supervised learning, a set of input-output pairs $(\mathbf{X}, \mathbf{Y}) \equiv \{\mathbf{x}(n), \mathbf{y}(n) | n = 1, \dots, N\}$ is given to a learner, where \mathbf{x} and \mathbf{y} represent an L -dim. input vector and a D -dim. output vector, respectively. The task of the learner is to predict the output \mathbf{y} for a new input data \mathbf{x} based on the training data (\mathbf{X}, \mathbf{Y}) and some prior knowledge of the problem. The RVM [1] is defined as a probabilistic model for the kernel representation:

$$P(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{W}, \sigma) = \mathcal{N}(\mathbf{y}|\mathbf{f}(\mathbf{x}, \mathbf{X}, \mathbf{W}), \sigma^{-1}), \quad (1)$$

$$\mathbf{f}(\mathbf{x}, \mathbf{X}, \mathbf{W}) = \sum_{n=0}^N \mathbf{w}_n \phi_n(\mathbf{x}) = \sum_{n=0}^N \mathbf{w}_n K(\mathbf{x}, \mathbf{x}(n)), \quad (2)$$

where \mathbf{w}_n denotes a D -dim. weight vector, $\phi_n(\mathbf{x}) \equiv K(\mathbf{x}, \mathbf{x}(n))$ ($n = 1, \dots, N$) is a kernel basis function defined on an input point $\mathbf{x}(n)$, and $\phi_0(\mathbf{x}) \equiv 1$. $\mathcal{N}(\mathbf{y}|\boldsymbol{\mu}, \sigma^{-1})$ denotes a normal distribution over \mathbf{y} with mean $\boldsymbol{\mu}$ and variance σ^{-1} . Accordingly, the conditional likelihood for the training data can be written as $P(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \sigma) = \prod_{n=1}^N \mathcal{N}(\mathbf{y}(n)|\mathbf{f}(\mathbf{x}(n), \mathbf{X}, \mathbf{W}), \sigma^{-1})$. Applying the Bayesian method together with the ARD prior [1, 2], one can get a sparse kernel representation: $\mathbf{y} = \mathbf{f}(\mathbf{x}, \mathbf{X}_B, \mathbf{W}) = \sum_{j=0}^J \mathbf{w}_j \phi_j(\mathbf{x})$, where $\mathbf{X}_B \equiv \{\mathbf{x}_B(j) | j = 1, \dots, J\} \subset \mathbf{X}$ represents a reduced set of input data that defines the sparse kernel basis $\phi_j(\mathbf{x}) \equiv K(\mathbf{x}, \mathbf{x}_B(j))$ and is called the relevance vector [1].

3 Incremental Sparse Kernel Machine

3.1 Incremental Bayesian method

In the original RVM [1, 2], the posterior parameter distribution is calculated for the entire dataset (\mathbf{X}, \mathbf{Y}) . The incremental Bayesian method proposed here uses only a part of the dataset in one learning epoch. At the beginning of the learning, a dataset $(\mathbf{X}_D, \mathbf{Y}_D) \equiv \{\mathbf{x}_D(m), \mathbf{y}_D(m) | m = 1, \dots, M\}$ is selected from the dataset (\mathbf{X}, \mathbf{Y}) . The current basis set $\mathbf{X}_B \equiv \{\mathbf{x}_B(j) | j = 1, \dots, J\}$ is also selected from the input dataset \mathbf{X} . The conditional likelihood for the current dataset $(\mathbf{X}_D, \mathbf{Y}_D)$ under the current basis set \mathbf{X}_B is given by

$$P(\mathbf{Y}_D|\mathbf{X}_D, \mathbf{X}_B, \mathbf{W}, \sigma) = \prod_{m=1}^M \mathcal{N}(\mathbf{y}_D(m)|\mathbf{f}(\mathbf{x}_D(m), \mathbf{X}_B, \mathbf{W}), \sigma^{-1}), \quad (3)$$

where $\mathbf{f}(\mathbf{x}, \mathbf{X}_B, \mathbf{W})$ is defined in (2). We employ the hierarchical ARD prior [2] for the model parameters $\mathbf{W} \equiv \{\mathbf{w}_j | j = 0, \dots, J\}$:

$$P_0(\mathbf{W}|\sigma, \boldsymbol{\alpha}) = \prod_{j=0}^J \mathcal{N}(\mathbf{w}_j|\bar{\mathbf{w}}_{j0}, (\sigma\alpha_j)^{-1}), \quad P_0(\boldsymbol{\alpha}) = \prod_{j=0}^J \Gamma(\alpha_j|a_{j0}, \bar{\alpha}_{j0}^{-1}), \quad (4)$$

where the hyperparameters $\boldsymbol{\alpha} \equiv \{\alpha_j | j = 0, \dots, J\}$ are introduced for controlling the variance of each weight vector \mathbf{w}_j . The mean weight vector $\bar{\mathbf{w}}_{j0}$ of the ARD prior (4) is assumed to be zero, i.e., $\bar{\mathbf{w}}_{j0} \equiv \mathbf{0}$. This introduces a bias on \mathbf{w}_j toward a null vector. However, we have no idea about the proper value of the hyperparameter $\boldsymbol{\alpha}$. Therefore, the hyperparameter $\boldsymbol{\alpha}$ is integrated by introducing the hierarchical prior for $\boldsymbol{\alpha}$. The hyperprior for $\boldsymbol{\alpha}$ is assumed to be a product of the Gamma distribution defined by $\Gamma(\alpha|a, b) \equiv \alpha^{-1}(ab\alpha)^a \exp(-ab\alpha)/\Gamma(a)$, which is the conjugate prior for the inverse-variance parameter. $\Gamma(a)$ is the Gamma function defined by $\Gamma(a) \equiv \int_0^\infty dt t^{a-1} e^{-at}$. We also introduce a hierarchical prior for the inverse of the output variance parameter σ :

$$P_0(\sigma|\tau) = \Gamma(\sigma|\gamma_{\sigma 0}, \tau), \quad P_0(\tau) = \Gamma(\tau|\gamma_{\tau 0}, \bar{\tau}_0^{-1}). \quad (5)$$

3.2 Incremental Sparse Kernel Machine

The posterior distribution over the model parameters $\boldsymbol{\theta} = (\mathbf{W}, \sigma)$ and the hyperparameters $\boldsymbol{\xi} = (\boldsymbol{\alpha}, \tau)$ is given by

$$P(\boldsymbol{\theta}, \boldsymbol{\xi} | \mathbf{Y}_D, \mathbf{X}_D, \mathbf{X}_B) = \frac{P(\mathbf{Y}_D, \boldsymbol{\theta}, \boldsymbol{\xi} | \mathbf{X}_D, \mathbf{X}_B)}{P(\mathbf{Y}_D | \mathbf{X}_D, \mathbf{X}_B)}, \quad (6)$$

$$P(\mathbf{Y}_D, \boldsymbol{\theta}, \boldsymbol{\xi} | \mathbf{X}_D, \mathbf{X}_B) = P(\mathbf{Y}_D | \mathbf{X}_D, \mathbf{X}_B, \mathbf{W}, \sigma) P_0(\mathbf{W} | \sigma, \boldsymbol{\alpha}) P_0(\sigma | \tau) P_0(\boldsymbol{\alpha}) P_0(\tau),$$

$$P(\mathbf{Y}_D | \mathbf{X}_D, \mathbf{X}_B) = \int d\boldsymbol{\theta} d\boldsymbol{\xi} P(\mathbf{Y}_D, \boldsymbol{\theta}, \boldsymbol{\xi} | \mathbf{X}_D, \mathbf{X}_B). \quad (7)$$

Since the integration in the marginal likelihood (evidence) $P(\mathbf{Y}_D | \mathbf{X}_D, \mathbf{X}_B)$ is analytically intractable, we use the variational Bayes (VB) method developed in [6, 7, 2]. In order to approximate the posterior (6), the VB method introduces a trial (variational) posterior $Q(\boldsymbol{\theta}, \boldsymbol{\xi})$ and defines the free energy for Q :

$$\begin{aligned} F[Q] &= \int d\boldsymbol{\theta} d\boldsymbol{\xi} Q(\boldsymbol{\theta}, \boldsymbol{\xi}) \log(P(\mathbf{Y}_D, \boldsymbol{\theta}, \boldsymbol{\xi} | \mathbf{X}_D, \mathbf{X}_B) / Q(\boldsymbol{\theta}, \boldsymbol{\xi})) \\ &= \log P(\mathbf{Y}_D | \mathbf{X}_D, \mathbf{X}_B) - \int d\boldsymbol{\theta} d\boldsymbol{\xi} Q(\boldsymbol{\theta}, \boldsymbol{\xi}) \log \frac{Q(\boldsymbol{\theta}, \boldsymbol{\xi})}{P(\boldsymbol{\theta}, \boldsymbol{\xi} | \mathbf{Y}_D, \mathbf{X}_D, \mathbf{X}_B)} \end{aligned} \quad (8)$$

The second term in (8) is the Kullback-Leibler divergence between the trial posterior $Q(\boldsymbol{\theta}, \boldsymbol{\xi})$ and the true posterior $P(\boldsymbol{\theta}, \boldsymbol{\xi} | \mathbf{Y}_D, \mathbf{X}_D, \mathbf{X}_B)$ given by (6). Since the first term in (8) does not depend on $Q(\boldsymbol{\theta}, \boldsymbol{\xi})$, the free energy $F[Q]$ is maximized when $Q(\boldsymbol{\theta}, \boldsymbol{\xi})$ is equal to $P(\boldsymbol{\theta}, \boldsymbol{\xi} | \mathbf{Y}_D, \mathbf{X}_D, \mathbf{X}_B)$. In this case, the free energy $F[Q]$ becomes equal to the log-evidence $\log P(\mathbf{Y}_D | \mathbf{X}_D, \mathbf{X}_B)$. The VB method assumes a factorized form of $Q(\boldsymbol{\theta}, \boldsymbol{\xi})$:

$$Q(\boldsymbol{\theta}, \boldsymbol{\xi}) = Q_\theta(\boldsymbol{\theta}) Q_\xi(\boldsymbol{\xi}). \quad (9)$$

This factorization assumption is weaker than the factorization assumption used in [2], where each component of the parameters and the hyperparameters is assumed to be factorized. Under the assumption of (9), the free energy is alternately maximized with respect to $Q_\theta(\boldsymbol{\theta})$ and $Q_\xi(\boldsymbol{\xi})$. The maximized free energy gives a lower bound on the log-evidence.

The alternate free energy maximization leads to the functional form of the posterior given by

$$\begin{aligned}
Q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) &= Q_{\mathbf{W}}(\mathbf{W}|\sigma)Q_{\sigma}(\sigma), & Q_{\boldsymbol{\xi}}(\boldsymbol{\xi}) &= Q_{\boldsymbol{\alpha}}(\boldsymbol{\alpha})Q_{\tau}(\tau), \\
Q_{\mathbf{W}}(\mathbf{W}|\sigma) &= \prod_{j=1}^J \mathcal{N}(\mathbf{w}_j|\bar{\mathbf{w}}_j, \sigma^{-1}\Sigma^{-1}), & Q_{\boldsymbol{\alpha}}(\boldsymbol{\alpha}) &= \prod_{j=0}^J \Gamma(\alpha_j|a_j, \bar{\alpha}_j^{-1}), \\
Q_{\sigma}(\sigma) &= \Gamma(\sigma|\gamma_{\sigma}, \Delta), & Q_{\tau}(\tau) &= \Gamma(\tau|\gamma_{\tau}, \bar{\tau}^{-1}), \tag{10}
\end{aligned}$$

where the $((J+1) \times D)$ matrix $\bar{\mathbf{W}} \equiv (\bar{\mathbf{w}}_0, \bar{\mathbf{w}}_1, \dots, \bar{\mathbf{w}}_J)'$, the $((J+1) \times (J+1))$ inverse-covariance matrix Σ , $\gamma_{\sigma}, \Delta, \{a_j, \bar{\alpha}_j | j = 0, \dots, J\}, \gamma_{\tau}$, and $\bar{\tau}$ are determined in the following steps.

1) In the maximization step with respect to the posterior parameter distribution $Q_{\boldsymbol{\theta}}(\boldsymbol{\theta})$, the parameters in $Q_{\boldsymbol{\theta}}(\boldsymbol{\theta})$ are updated:

$$\begin{aligned}
\Sigma &= \Phi' \Phi + \bar{\mathbf{A}}, \quad \bar{\mathbf{W}} = \Sigma^{-1}(\Phi' \mathbf{Y}_D + \bar{\mathbf{A}} \bar{\mathbf{W}}_0), \quad \gamma_{\sigma} = \gamma_{\sigma 0} + DM/2, \tag{11} \\
\gamma_{\sigma} \Delta &= \frac{1}{2} \text{Tr}(\mathbf{Y}_D - \Phi \bar{\mathbf{W}})'(\mathbf{Y}_D - \Phi \bar{\mathbf{W}}) + \frac{1}{2} \text{Tr}(\bar{\mathbf{W}} - \bar{\mathbf{W}}_0)' \bar{\mathbf{A}}(\bar{\mathbf{W}} - \bar{\mathbf{W}}_0) + \gamma_{\sigma 0} \bar{\tau},
\end{aligned}$$

where the $(M \times (J+1))$ matrix Φ , the $((J+1) \times (J+1))$ matrix $\bar{\mathbf{A}}$, and the $(M \times D)$ matrix \mathbf{Y}_D are defined as $(\Phi)_{mj} = \phi_j(\mathbf{x}_D(m)) = K(\mathbf{x}_D(m), \mathbf{x}_B(j))$, $\bar{\mathbf{A}} = \text{diag}(\bar{\alpha}_0, \bar{\alpha}_1, \dots, \bar{\alpha}_J)$, and $\mathbf{Y}_D = (\mathbf{y}(1), \dots, \mathbf{y}(M))'$, respectively. The mean weight of the posterior is given by $\langle \mathbf{W} \rangle_{\boldsymbol{\theta}} \equiv \int d\boldsymbol{\theta} Q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) \mathbf{W} = \bar{\mathbf{W}}$, and the mean of the inverse-variance parameter σ is given by $\langle \sigma \rangle_{\boldsymbol{\theta}} = \Delta^{-1}$. The covariance matrix for the weight \mathbf{W} is given by $\sigma^{-1} \Sigma^{-1}$. γ_{σ} is a degree of freedom of the Gamma distribution $Q_{\sigma}(\sigma)$ and the variance of σ is proportional to γ_{σ}^{-1} .

2) In the maximization step with respect to the posterior hyperparameter distribution $Q_{\boldsymbol{\xi}}(\boldsymbol{\xi})$, the parameters in $Q_{\boldsymbol{\xi}}(\boldsymbol{\xi})$ are updated:

$$\begin{aligned}
a_j &= a_{j0} + \frac{1}{2}D, \quad a_j \bar{\alpha}_j^{-1} = a_{j0} \bar{\alpha}_{j0}^{-1} + \frac{1}{2} (\Delta^{-1} \|\bar{\mathbf{w}}_j - \bar{\mathbf{w}}_{j0}\|^2 + D(\Sigma^{-1})_{jj}), \tag{12} \\
\gamma_{\tau} &= \gamma_{\tau 0} + \gamma_{\sigma 0}, \quad \gamma_{\tau} \bar{\tau}^{-1} = \gamma_{\tau 0} \bar{\tau}_0^{-1} + \gamma_{\sigma 0} \Delta^{-1}.
\end{aligned}$$

The means of the hyperparameters are given by $\langle \alpha_j \rangle_{\boldsymbol{\xi}} \equiv \int d\boldsymbol{\xi} Q_{\boldsymbol{\xi}}(\boldsymbol{\xi}) \alpha_j = \bar{\alpha}_j$ and $\langle \tau \rangle_{\boldsymbol{\xi}} = \bar{\tau}$.

3) After the convergence, this step selects the weights that have norm $\|\bar{\mathbf{w}}_j\|^2$ larger than the threshold value W_{\min} . The corresponding basis points $\mathbf{x}_B(j)$ form an old basis set \mathbf{X}_{old} . In the next epoch, we select a new dataset $(\mathbf{X}_{\text{new}}, \mathbf{Y}_{\text{new}})$ from the entire dataset (\mathbf{X}, \mathbf{Y}) . The basis set and the dataset in the next epoch are given by $\mathbf{X}_B = \{\mathbf{X}_{\text{new}}, \mathbf{X}_{\text{old}}\}$ and $(\mathbf{X}_D, \mathbf{Y}_D) = (\{\mathbf{X}_{\text{new}}, \mathbf{X}_{\text{old}}\}, \{\mathbf{Y}_{\text{new}}, \mathbf{Y}_{\text{old}}\})$, respectively.

4) We have no specific information on the weights \mathbf{W}_{new} that correspond to the new basis points \mathbf{X}_{new} . Therefore, we use the ARD prior defined in (4) with a small a_{j0} and $\bar{\mathbf{w}}_{j0} = 0$. On the other hand, we do have information on the weights \mathbf{W}_{old} that correspond to the old basis points \mathbf{X}_{old} . Accordingly, we use the obtained value of $\bar{\mathbf{w}}_j, \bar{\alpha}_j$, and a_j as the prior parameter in the next

epoch, i.e., $\bar{\mathbf{w}}_{j0} = \bar{\mathbf{w}}_j, \bar{\alpha}_{j0} = \bar{\alpha}_j$, and $a_{j0} = a_j$. An alternative method is to use the obtained posterior $Q_{\mathbf{W}}(\mathbf{W}_{\text{old}}|\sigma)$ as the new prior for \mathbf{W}_{old} . In this case, the matrix $\bar{\mathbf{A}}$ in (11) becomes a block diagonal matrix and the block matrices are given by $\bar{\mathbf{A}}_{\text{old}} = \Sigma_{\text{old}}$ and $\bar{\mathbf{A}}_{\text{new}} = \text{diag}(\bar{\alpha}_{\text{new}})$ in an obvious notation. We also have information on the inverse-variance parameter σ , so we reset the prior parameter as $\bar{\tau}_0 = \bar{\tau}, \gamma_{\sigma 0} = \gamma_{\sigma}$, and $\gamma_{\tau 0} = \gamma_{\tau}$ or use the obtained posterior $Q_{\sigma}(\sigma)$ as the new prior for σ .

We repeat the above process until all data (\mathbf{X}, \mathbf{Y}) are processed. The predictive distribution after the learning can be calculated by using the posterior parameter distribution: $\hat{P}(\mathbf{y}|\mathbf{x}, \mathbf{X}_B) = \int d\mathbf{W} d\sigma P(\mathbf{y}|\mathbf{x}, \mathbf{X}_B, \mathbf{W}, \sigma) Q_{\mathbf{W}}(\mathbf{W}|\sigma) Q_{\sigma}(\sigma) = \mathcal{T}(\mathbf{y}|\bar{\mathbf{W}}' \cdot \phi(\mathbf{x}), \Delta(1 + \phi(\mathbf{x})' \Sigma^{-1} \phi(\mathbf{x})), 2\gamma_{\sigma})$, where $\mathcal{T}(\mathbf{y}|\boldsymbol{\mu}, C, \gamma)$ denotes the t-distribution with mean $\boldsymbol{\mu}$, variance $C/(1 - 2\gamma^{-1})$, and degree of freedom γ . $\phi(\mathbf{x})$ denotes the $(J + 1)$ -dim. vector defined by $\phi_j(\mathbf{x}) = K(\mathbf{x}, \mathbf{x}_B(j)) (j = 1, \dots, J)$ and $\phi_0(\mathbf{x}) \equiv 1$.

4 Experiments

The applicability of our incremental Sparse Kernel Machine (SKM) was investigated by using large datasets. The Gaussian kernel is used in all experiments. The problem is a prediction task for the chaotic Mackey-Glass (M-G) model, which is defined by the differential equation

$$ds(t)/dt = -bs(t) + as(t - \tau)/(1 + s(t - \tau)^{10}), \quad (13)$$

where $a = 0.2, b = 0.1$, and $\tau = 17$ [8]. The task was to predict $y = s(t + 85)$ from the delay coordinate $\mathbf{x} = (s(t), s(t - 6), s(t - 12), s(t - 18))$. The training data were prepared by adding the Gaussian noise with different levels to the generated time-series. The N/S (noise-to-signal) ratios of the added noise were 0.0, 0.11, and 0.22 in terms of standard deviation. Three datasets with 500, 1000, and 10000 data points were used to train the incremental SKM, the ν -SVM [9, 10] and the Kernel Principal Components Regression (KPCR) method [8]. The NRMSE (normalized root mean squared error) for the test dataset and the corresponding number of bases are listed in Table 1. The SKM showed better performance for the data with noise than did the ν -SVM or the KPCR, while the ν -SVM and the KPCR showed better performance for noiseless data. It should be noted that the SKM achieved good performance with a much fewer bases than the ν -SVM. The SKM selected moderate number of bases even when a large amount of data were given. Fig. 1 shows the basis points (relevance vector) obtained by the SKM for 10000 training data together with the M-G attractor.

5 Discussion

In this paper, we formulated the incremental Sparse Kernel Machine (SKM) and showed the efficiency of this method for large datasets. The Bayes estimation using the ARD prior automatically eliminated insignificant bases and gave a sparse

kernel representation for large datasets. We also examined another implementation of SKM, where the dataset in one epoch was the entire dataset and the basis set was given by $\mathbf{X}_B = \{\mathbf{X}_{\text{new}}, \mathbf{X}_{\text{old}}\}$. However, this implementation did not significantly improve performance, while it required $O((\# \text{ of bases}) * (\# \text{ of data})^2)$ computation time. In this paper, we only considered regression problems. However, our method could be easily extended to classification problems by using the same method developed in [2]. Application of our method to real-world problems remains a task for future research.

Table 1. NRMSE for SKM, ν -SVM and KPCR. The number of bases are presented in parentheses.

	n/s=0.0			n/s=0.11			n/s=0.22		
	500	1000	10000	500	1000	10000	500	1000	10000
SKM	0.130 (145)	0.109 (148)	0.088 (169)	0.238 (142)	0.201 (250)	0.110 (290)	0.333 (223)	0.322 (413)	0.146 (334)
ν -SVM	0.037 (462)	0.013 (704)	0.004 (3177)	0.275 (333)	0.197 (507)	0.148 (2201)	0.467 (334)	0.395 (590)	0.355 (2769)
KPCR	0.038	0.008	***	0.307	0.280	***	0.443	0.414	***

References

1. Tipping, M. E. (2000). The Relevance Vector Machine. *NIPS* **12**, pp. 652-658.
2. Bishop, C. M. & Tipping, M. E. (2000). Variational relevance vector machines. *Proc. of 16th Conf. UAI*, pp. 46-53.
3. V. N. Vapnik. (1995). *The Nature of Statistical Learning Theory*. Springer-Verlag.
4. Williams, C. K. I. & Rasmussen, C. E. (1996). Gaussian Process for Regression. *NIPS* **8**, pp. 514-520.
5. Csató, L. & Opper, M. (2001). Sparse Representation for Gaussian Process Models. *NIPS* **13**.
6. Attias, H. (1999). Inferring parameters and structure of latent variable models by variational Bayes. *Proc. of 15th Conf. UAI*, pp. 21-30.
7. Sato, M. (2001). On-line Model Selection Based on the Variational Bayes. *Neural Computation*, **13**, 1649-1681.
8. Rosipal, R. et al. (2001). Kernel PCA for Feature Extraction and De-Noising in Non-linear Regression. *Neural Computing & Applications*, 10(3).
9. Chang, C. and Lin, C., (2001), LIB-SVM: a library for SVM. (available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.)
10. Scholkopf, B. et al. (2000). New support vector algorithms. *Neural Computation*, **12**, 1207-1245.

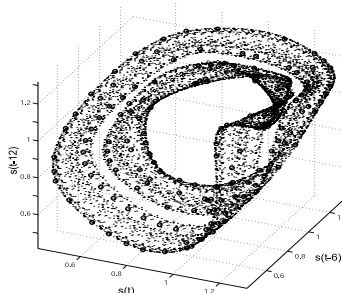


Fig. 1. Relevance vector

⁰ **Acknowledgement:** This research was supported in part by the Telecommunications Advancement Organization of Japan.