
A NOISE-ROBUST FAST SPARSE BAYESIAN LEARNING MODEL

Ingvild M. Helgøy *

Yushu Li †

ABSTRACT

This paper utilizes the hierarchical model structure from the Bayesian Lasso in the Sparse Bayesian Learning process to develop a new type of probabilistic supervised learning approach. This approach has several performance advantages, such as being fast, sparse and especially robust to the variance in random noise. The hierarchical model structure in this Bayesian framework is designed in such a way that the priors do not only penalize the unnecessary complexity of the model but also depend on the variance of the random noise in the data. The hyperparameters in the model are estimated by the Fast Marginal Likelihood Maximization algorithm and can achieve low computational cost and faster learning process. We compare our methodology with two other popular Sparse Bayesian Learning models: The Relevance Vector Machine and a sparse Bayesian model that has been used for signal reconstruction in compressive sensing. We show that our method will generally provide more sparse solutions and be more flexible and stable when data is polluted by high variance noise.

Keywords Sparse Bayesian Learning · Bayesian Lasso · hierarchical models · kernel basis function · type-II maximum likelihood

1 Introduction

1.1 Framework of Sparse Bayesian Learning

Supervised learning contains a set of training data $\{\mathbf{x}_i, y_i\}_{i=1}^N$, where for $i = 1, 2, \dots, N$, $\mathbf{x}_i \in \mathbb{R}^D$ is the i 'th observation of the D -dimensional input variable, \mathbf{x} , and $y_i \in \mathbb{R}$ is the corresponding scalar value of the output (target) variable, \mathbf{y} . Based on the training data, we are aiming at constructing a mathematical function $f(\mathbf{x})$ that can model the underlying relationship between the input covariates \mathbf{x} and the target variable \mathbf{y} . A common way to construct $f(\mathbf{x})$ is to approximate the function in the space that is linearly spanned by a set of M basis functions:

$$f(\mathbf{x}) \approx \hat{f}(\mathbf{x}) = \sum_{m=1}^M w_m \phi_m(\mathbf{x}). \quad (1)$$

Hence, the approximation $\hat{f}(\mathbf{x})$ is a weighted linear sum of M basis functions. A complete set of basis functions, $\{\phi_m(\mathbf{x})\}_{m=1}^M$, constructs the basis of a function space; some commonly used families of basis functions are polynomials, Fourier basis and different types of kernels functions. The approximation in Equation (1) can be in any nonlinear form, as although the formula is linear in the parameters, the basis functions can be nonlinear and M can be large or infinite.

*Department of Mathematics, University of Bergen, P.O. Box 7803, N-5020 Bergen, Norway, Ingvild.Helgoy@uib.no

†Department of Mathematics, University of Bergen, P.O. Box 7803, N-5020 Bergen, Norway, Yushu.Li@uib.no

Part of the learning process is to choose a suitable set of basis functions and the value of M , so that we can get the best approximation based on certain criteria such as Mean Squared Error (MSE).

In practical applications, the observed outputs y_i are always samples from the model with additive noise, and the Gaussian distribution is commonly utilized to model the noisy random error. Let $\Phi = [\phi_1, \dots, \phi_M]$ be an $N \times M$ design matrix whose column vectors are $\phi_j = [\phi_j(\mathbf{x}_1), \dots, \phi_j(\mathbf{x}_N)]^\top$; $j = 1, \dots, M$, $\mathbf{w} = (w_1, \dots, w_M)^\top$ be the weight vector which consists of M weight parameters, $\mathbf{y} = (y_1, \dots, y_N)^\top$ be the observed values of the target variable, $\epsilon = (\epsilon_1, \dots, \epsilon_N)^\top$ be the random error vector, and \mathbf{I}_N denotes the $N \times N$ identity matrix. Then, we have

$$\mathbf{y} = \Phi \mathbf{w} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_N), \quad (2)$$

where σ^2 is the variance of the error terms that are normally distributed. In the sparse learning framework, the model (2) is usually implicitly defined such that $M = N$ and $\phi_m(\mathbf{x}) = K(\mathbf{x}, \mathbf{x}_m)$; $m = 1, \dots, N$, where $K(\cdot, \cdot)$ is a positive definite kernel function centered at each of the training input vectors. Thus, each basis function $\phi_m(\mathbf{x})$ corresponds to one training input vector \mathbf{x}_m . After the type of kernel function is chosen, the learning task is then to estimate the weight parameters $\mathbf{w} = (w_1, \dots, w_N)^\top$ from the training data. Denoting the point estimation of the weight parameters as $\hat{\mathbf{w}} = (\hat{w}_1, \dots, \hat{w}_N)^\top$, for a new observed input \mathbf{x}^* , one reasonable point prediction for the target output is: $y^* = \sum_{i=1}^N \hat{w}_i \phi_i(\mathbf{x}^*)$. The key feature of the sparse learning is that a large part of the estimated weight parameters is set to zero during the learning process. Thus, the model achieves sparsity with only few \hat{w}_i being non-zero, and the corresponding basis functions can be used in prediction and approximation.

There already exist several studies that utilize different approaches to achieve this sparse estimation under the structure of the model (2) (Boser et al. 1992, Vapnik et al. 1997, Schölkopf et al. 1999, Tipping 2001, Ji et al. 2008, Babacan et al. 2010). In the field of kernel based machine learning, the Support Vector Machine (SVM) (Boser et al. 1992, Vapnik et al. 1997, Schölkopf et al. 1999) is one of the most popular methods. However, the kernels in SVM must satisfy Mercer's condition (Smola et al. 1998, Schölkopf 2001). Moreover, the SVM is purely deterministic as the SVM output is just a point estimate. Tipping (2001) further indicated several disadvantages of SVM and proposed a probabilistic sparse kernel learning approach in the Bayesian framework, that is called the Relevance Vector Machine (RVM). It can be shown that both SVM and RVM are related to Gaussian processes which are important Bayesian machine learning models (Seeger 2000, Rasmussen & Quinonero-Candela 2005, Williams & Rasmussen 2006). The RVM has been widely used in various applications (see, e.g., Agarwal & Triggs 2005, Ashburner 2007, Demir & Erturk 2007, Ghosh & Mujumdar 2008) and several extensions of the model can be found in the literature (see, e.g., Wipf & Rao 2004, Krishnapuram et al. 2005, Schmolck & Everson 2007, Ji et al. 2008). The hierarchical structure in RVM places an "Automatic Relevance Determination" (ARD) zero-mean Gaussian prior (MacKay 1992, Neal 2012) on the weight parameters $\mathbf{w} = (w_1, \dots, w_N)^\top$, where each weight parameter has its own variance in the Gaussian prior. The inverse of these variances are defined as a set of hyperparameters $\alpha = (\alpha_1, \dots, \alpha_N)^\top$. Those hyperparameters are viewed as precision variables and have their own hyperprior distribution, which is often a non-informative prior such as the Gamma distribution. Based on the ARD prior, if one basis function $\phi_m(\mathbf{x})$; $m \in 1, \dots, N$, can not directly contribute to the likelihood of the data, the joint likelihood of the data and \mathbf{w} is maximized by setting α_m to an infinitely large value. The corresponding posterior of w_m will then be infinitely peaked at zero so that $\phi_m(\mathbf{x})$ will be pruned. The related input vector \mathbf{x}_m will also be switched off and called "irrelevant vector". Further inferences and prediction can be carried out based on the posterior distributions of the rest non-zero weight parameters.

Other forms of priors on the weight parameter can be used to form an ARD prior. Recently, Babacan et al. (2010) utilized again a Gaussian prior on \mathbf{w} in a hierarchical manner where the hyperparameters $\gamma = (\gamma_1, \dots, \gamma_N)^\top$ are defined directly as the variances of the weight parameters, and an exponential hyperprior is set to those hyperparameters. In this way, if the hyperparameter γ_m is estimated to be zero, the corresponding weight parameter w_m is also set to be zero and the related basis function $\phi_m(\mathbf{x})$ is pruned. As explained by Babacan et al. (2010), the basic prior for \mathbf{w} after integrating all the hyperparameters in γ is a Laplace distribution. Furthermore, Babacan et al. (2010) showed that the

resulting model is at least as sparse as the RVM proposed by Tipping (2001), where the basic prior in the RVM is a Student's t -distribution after integrating all the hyperparameters in α . The Laplace prior is also log-concave, which can lead to unimodal posterior distribution and eliminate local minima (Seeger & Nickisch 2008, Wipf et al. 2007). The learning approach proposed by Babacan et al. (2010) can be used in the reconstruction of signals in compressive sensing, and we refer to Candes et al. (2004) and Donoho et al. (2006) for more detailed algorithms.

No matter which ARD prior is utilized, the essential step in the learning process is to estimate the hyperparameters in the ARD prior, and both Tipping (2001) and Babacan et al. (2010) utilize a type-II maximum likelihood method (Bishop 2006, Williams & Rasmussen 2006) in estimation. Furthermore, in the process of optimizing the marginal likelihood function to estimate the hyperparameters, Faul & Tipping (2002) and Tipping et al. (2003) proposed a highly accelerated algorithm so that the whole computation process performs much faster. With this algorithm, the methods proposed by Tipping (2001) and Babacan et al. (2010) can be utilized for large dataset.

1.2 Our Bayesian Lasso Based Sparse Learning

In the pure classic linear regression setting, the Lasso of Tibshirani (Tibshirani 1996) managed to achieve sparsity by using a L_1 regularization in the model (2). The estimates of the weight parameters are the solutions of the penalized regression, and sparsity in the Lasso model is obtained since some of the estimated weight parameters in \mathbf{w} will be zero (Tibshirani 1996, Efron et al. 2004). The Lasso solution can also be interpreted as a Bayesian posterior mode estimate when the prior for the weight parameters, $p(\mathbf{w})$, is a Laplace distribution (Tibshirani 1996, Friedman et al. 2001). However, the posterior mode is not the natural choice to obtain point estimates in the Bayesian framework, since a fully Bayesian analysis would instead suggest using the mean or median of the posterior as point estimates. A fully Bayesian model of the Lasso was later introduced by Park & Casella (2008), where they use a conditional prior on \mathbf{w} of the form $p(\mathbf{w}|\sigma^2)$. The approximation of the posterior distribution for the weight parameters can, thereafter, be obtained by using the Gibbs sampler. However, almost none of the estimated weight parameters from the Gibbs sampling will be set exactly to zero. Hence, the Bayesian Lasso does not perform variable selection, thus, it is not a sparse model. Nevertheless, the Bayesian Lasso by Park & Casella (2008) has several attractive properties when it comes to parameter and hyperparameter estimation in the Bayesian framework. This includes the joint posterior distribution for \mathbf{w} , and σ^2 generally has one mode instead of multiple posterior modes. Recently, Balakrishnan & Madigan (2009) argued for the advantages of Bayesian Lasso and tried to combine the conditional prior in Bayesian Lasso model with the Sparse Bayesian learning algorithms to achieve feature selection. However, the prior distribution for \mathbf{w} by Balakrishnan & Madigan (2009) is not the same conditional prior $p(\mathbf{w}|\sigma^2)$ from the Bayesian Lasso by Park & Casella (2008). Instead, their prior for \mathbf{w} is the Laplace distribution, same as the prior by Babacan et al. (2010).

In our paper, we combine the hierarchical structure in the Bayesian Lasso proposed by Park & Casella (2008) with the estimation and inference procedure proposed by Tipping (2001), Faul & Tipping (2002) and Tipping et al. (2003) to achieve a fast sparse probabilistic learning process. As the prior distribution of the weight parameters from Park & Casella (2008) will be conditional on the variance of the random noise, our method should theoretically be more flexible to the noisy dataset. More specifically, in our method, the variances of the weights are associated with independent hyperparameters, and we prove that these hyperparameters will be set to zero when the estimated values based on the type-II maximum likelihood is lower than a threshold. This threshold is related to the variance of the random error, which measures the extent of the noise. Consequently, our method is robust to the noise in the dataset. Furthermore, we conduct a comprehensive simulation study to compare our method with the two fast sparse Bayesian learning methods derived from Tipping et al. (2003) and Babacan et al. (2010). The results indicate that our method generally performs better than the other two methods when the dataset is noisy.

The remainder of the paper is divided into the following sections: Section 2 introduces the Bayesian framework for the sparse supervised learning, and special attention is given to the two sparse Bayesian methods given by Tipping (2001)

and Babacan et al. (2010). The hierarchical structure in the Bayesian Lasso as described by Park & Casella (2008) will be presented in this section as well. Section 3 contains a detailed description of our method, including the fast optimization algorithm when using the type-II maximum likelihood method to estimate the hyperparameters. Section 4 presents the simulation study, and the conclusion of the study is presented in Section 5.

2 The Bayesian Framework

When using a Bayesian approach for the regression problem in Equation (2), the unknown parameters in the vector \mathbf{w} and σ^2 are treated as stochastic variables that have their own prior probability distribution. The prior distribution expresses our prior belief about the values that those parameters might take. In sparse Bayesian learning, the weight parameters in \mathbf{w} are often given a prior distribution $p(\mathbf{w}|\gamma)$, where γ is a vector of the hyperparameters which are assigned a hyperprior distribution. Further, the likelihood function of the complete data set \mathbf{y} in model (2) can be written as

$$p(\mathbf{y}|\mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2),$$

where Φ is the kernel matrix, with the matrix elements $\Phi_{ij} = \phi_j(\mathbf{x}_i)$; $i = 1, \dots, N$, $j = 1, \dots, M$, for the input vector, \mathbf{x}_i . From Bayes' rule, the posterior distribution for all the unknown parameters given the observed data \mathbf{y} is given by

$$p(\mathbf{w}, \gamma, \sigma^2|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{w}, \gamma, \sigma^2)p(\mathbf{w}, \gamma, \sigma^2)}{p(\mathbf{y})}. \quad (3)$$

Given a new test input \mathbf{x}^* , predictions for the output y^* can be achieved by using the predictive distribution:

$$p(y^*|\mathbf{y}) = \int p(y^*|\mathbf{w}, \gamma, \sigma^2)p(\mathbf{w}, \gamma, \sigma^2|\mathbf{y}) d\mathbf{w} d\gamma d\sigma^2. \quad (4)$$

However, the predictive distribution (4) must often be approximated as the normalizing constant, $p(\mathbf{y})$, in Equation (3) is almost never possible to calculate directly. The RVM in Tipping (2001) and the Fast Laplace method in Babacan et al. (2010) try to approximate the predictive distribution (4) analytically, while the Bayesian Lasso in Park & Casella (2008) uses a numerical approach. We now give a brief description of the learning and inference process in the RVM described by Tipping (2001), the Fast Laplace method described by Babacan et al. (2010), as well as the Bayesian Lasso from Park & Casella (2008) in Sections 2.1-2.3, before we go to our method in Section 3.

2.1 The Relevance Vector Machine

In the RVM (Tipping 2001), where $M = N$ for the zero constant regression, a Gaussian prior is placed on the weight parameters:

$$p(\mathbf{w}|\alpha) = \prod_{i=1}^N \mathcal{N}(w_i|0, \alpha_i^{-1}),$$

where $\alpha = (\alpha_1, \dots, \alpha_N)^\top$ denote the inverse of the variances for the weights and can be viewed as precision hyperparameters. These hyperparameters in α are again defined as stochastic variables, and a Gamma distribution is used as the hyperprior. The Gamma distribution is also set as the prior for the inverse of the random error variance σ^2 . Based on those prior assumptions, Tipping (2001) shows that the underlying marginal prior, $p(\mathbf{w})$, is a Student's t-distribution which will enforce sparsity. The joint posterior distribution of the parameters, similar to Equation (3), can not be computed analytically from this form due to the normalizing constant. Therefore, Tipping (2001) utilize the following decomposition:

$$p(\mathbf{w}, \alpha, \sigma^2|\mathbf{y}) = p(\mathbf{w}|\mathbf{y}, \alpha, \sigma^2)p(\alpha, \sigma^2|\mathbf{y}). \quad (5)$$

The first term on the right-hand side of Equation (5) is the posterior distribution over the weights. By using Bayes' rule, this term can be computed analytically and Tipping (2001) shows that it is a Gaussian distribution with the following

mean vector and covariance matrix:

$$\boldsymbol{\mu} = \sigma^{-2} \boldsymbol{\Sigma} \boldsymbol{\Phi}^\top \mathbf{y}, \quad (6)$$

$$\boldsymbol{\Sigma} = (\sigma^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \boldsymbol{\Lambda})^{-1}, \quad (7)$$

where $\boldsymbol{\Lambda} = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_N)$. The second term on the right-hand side of Equation (5), however, can not be computed analytically. Tipping (2001) proposed a type-II maximum likelihood procedure by getting the maximization of the second term to obtain point estimates for $\boldsymbol{\alpha}$ and σ^2 . Denoting the estimates as $\boldsymbol{\alpha}_{MP}$ and σ_{MP}^2 , the estimates for the mean vector and covariance matrix for the Gaussian distribution in the first term are thereafter achieved by replacing $\boldsymbol{\alpha}$ and σ^2 in Equations (6) and (7) by $\boldsymbol{\alpha}_{MP}$ and σ_{MP}^2 . Further probabilistic inference for the weight parameter \mathbf{w} based on the posterior $p(\mathbf{w}|\mathbf{y}, \boldsymbol{\alpha}_{MP}, \sigma_{MP}^2)$ is also possible. Moreover, when using the type-II maximum likelihood procedure, some of the estimates for α will be set to an infinitely large value and the corresponding weights' posteriors will be infinitely peaked at 0. Thus, the related basis functions will be pruned from the model. As when $M = N$, each basis function corresponds to one input vector; the remaining input vectors related to the remainder non-zero weights are called "relevance vectors". Furthermore, given a new test input vector \mathbf{x}^* , the predictive distribution (4) can then be approximated by

$$p(y^*|\mathbf{y}, \boldsymbol{\alpha}_{MP}, \sigma_{MP}^2) = \int p(y^*|\mathbf{w}, \sigma_{MP}^2) p(\mathbf{w}|\mathbf{y}, \boldsymbol{\alpha}_{MP}, \sigma_{MP}^2) d\mathbf{w}.$$

Thus, one obvious advantage of RVM, compared to SVM, is that the RVM can give out probabilistic predictions and, thereby, capture the uncertainty of the prediction. SVM, on the other hand, can only give out deterministic predictions. This advantage of probabilistic predictions will be retained in other sparse Bayesian learning methods, including the method proposed in this paper.

2.2 The Fast Laplace Model

The Fast Laplace method introduced by Babacan et al. (2010) uses a Gaussian prior for the weights:

$$p(\mathbf{w}|\boldsymbol{\gamma}) = \prod_{i=0}^N \mathcal{N}(w_i|0, \gamma_i),$$

where $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_N)^\top$ is a vector of independent hyperparameters that are the individual variances for the weights. The hyperprior distribution for $\boldsymbol{\gamma}$ is the exponential distribution. The Fast Laplace model uses a similar approach as Faul & Tipping (2002) and Tipping et al. (2003) when it comes to inference and parameter estimation, where a fast type-II maximum likelihood method is utilized to estimate the hyperparameters in $\boldsymbol{\gamma}$. For the Fast Laplace method, some of the estimates γ_i are set to 0 and the corresponding weights are pruned from the model. Again, only input vectors related to the remainder non-zero weights are kept; they are the relevance vectors. Furthermore, after integrating out the hyperparameters γ_i , the basic prior for \mathbf{w} is a Laplace distribution:

$$p(\mathbf{w}) = \prod_{i=1}^N \frac{\lambda}{2} e^{-\lambda|w_i|}. \quad (8)$$

Compared to the Student's t-distribution, which is the basic prior in the RVM, using the Laplace prior results in a more sparse weight vector, \mathbf{w} . Babacan et al. (2010) give detailed mathematical proofs and provide comprehensive experiments to show that their method is at least as sparse as the RVM method.

2.3 The Bayesian Lasso

The Lasso of Tibshirani (1996) estimates the weights in the regression problem (2) with M basis functions by solving the optimization problem:

$$\min_{\mathbf{w}} (\mathbf{y} - \boldsymbol{\Phi} \mathbf{w})^\top (\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}) + \lambda \sum_{i=1}^M |w_i|$$

with $\lambda \geq 0$. This form indicates that the Lasso solution can be interpreted as a Bayesian posterior mode estimate when the prior for the weights w_i is the Laplace prior (8) (Tibshirani 1996, Friedman et al. 2001). However, a fully Bayesian approach would integrate over the posterior distribution or use the mean to obtain point estimates, instead of finding a posterior mode estimate. The Bayesian Lasso described by Park & Casella (2008) suggests a fully Bayesian analysis which estimates the posterior distribution numerically. Moreover, the Bayesian lasso uses the following prior for the weights:

$$p(\mathbf{w}|\gamma, \sigma^2) = \prod_{i=1}^M \mathcal{N}(w_i|0, \gamma_i \sigma^2), \quad (9)$$

where $\gamma_1, \dots, \gamma_M$ are individual variance parameters associated independently with every weight. Park & Casella (2008) use an exponential hyperprior for the hyperparameters γ_i and an inverse gamma hyperprior for σ^2 in Equation (9). After integrating out the hyperparameters γ_i , the following conditional Laplace prior for the weights is obtained:

$$p(\mathbf{w}|\sigma^2) = \prod_{i=1}^M \frac{\sqrt{\lambda}}{2\sqrt{\sigma^2}} e^{-\sqrt{\lambda}|w_i|/\sqrt{\sigma^2}}. \quad (10)$$

This follows from the representation of the Laplace distribution as a scaled mixture of Gaussians with an exponential mixing density (Andrews & Mallows 1974). Notice that the prior given in Equation (10) differs slightly from the prior given in Equation (8) since the prior in Equation (10) is conditioned on σ^2 . Conditioning on σ^2 is important because it results in a unimodal full posterior distribution (Park & Casella 2008). Based on the conditional posterior distribution for \mathbf{w} , σ^2 and γ , Park & Casella (2008) use the Gibbs sampler to obtain numerical estimates of the parameters \mathbf{w} and σ^2 . Since this is a fully Bayesian approach, it will, however, produce shrinkage of the weights and not provide a sparse solution.

3 The Bayesian Lasso-based Sparse Learning Model

We now propose a sparse Bayesian learning process to estimate the model in Equation (2). We will utilize the hierarchical structure of the Bayesian Lasso and introduce our methodology as the BLS method.

3.1 Hierarchical representation of the BLS model

Let Φ be a $N \times N$ matrix with matrix elements $\Phi_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$, $i = 1, \dots, N$; $j = 1, \dots, N$ and $\mathbf{w} = (w_1, \dots, w_N)^\top$ be the weight parameters where each weight corresponds to one input vector. The hierarchical representation of the full model in BLS is from the Bayesian Lasso described by Park & Casella (2008):

$$\begin{aligned} p(\mathbf{y}|\mathbf{w}, \sigma^2) &= \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2), \\ p(\mathbf{w}|\gamma, \sigma^2) &= \prod_{i=1}^N \mathcal{N}(w_i|0, \gamma_i \sigma^2), \end{aligned} \quad (11)$$

$$\begin{aligned} p(\gamma|\lambda) &= \prod_{i=1}^N \frac{\lambda_i}{2} \exp\left(-\frac{\lambda_i \gamma_i}{2}\right) \\ p(\lambda) &= \frac{\delta^v}{\Gamma(v)} (\lambda)^{v-1} e^{-\delta \lambda}, \quad (v > 0, \delta > 0) \end{aligned} \quad (12)$$

$$p(\sigma^2) = \frac{b^a}{\Gamma(a)} (\sigma^2)^{a-1} e^{-b\sigma^2}, \quad (a > 0, b > 0) \quad (13)$$

After integrating out the hyperparameters $\gamma_1, \dots, \gamma_N$ in Equation (11), the conditional prior distribution for the weights is given by

$$p(\mathbf{w}|\sigma^2) = \prod_{i=1}^N \frac{\sqrt{\lambda}}{2\sqrt{\sigma^2}} e^{-\sqrt{\lambda}|w_i|/\sqrt{\sigma^2}}. \quad (14)$$

Thus, the prior in Equation (14) is a Laplace prior conditioning on σ^2 and is the same as the conditional prior in the Bayesian Lasso given by Equation (10). For the parameter λ in Equation (14), Park & Casella (2008) offer two different methods to obtain the estimate, where one of them uses the hyperprior given by Equation (12). For the parameter σ^2 , Park & Casella (2008) use the scale invariant prior $p(\sigma^2) = 1/\sigma^2$; this can be obtained from the hyperprior given by Equation (13) by setting $a = 0$ and $b = 0$, which are the parameters in the Gamma distribution.

Having defined the prior, the posterior of all the parameters, given the data, is $p(\mathbf{w}, \gamma, \sigma^2, \lambda | \mathbf{y})$. Like the RVM (Tipping 2001), this posterior can not be found directly from Bayes' rule. We therefore use the same decomposition as proposed by Tipping (2001):

$$p(\mathbf{w}, \gamma, \sigma^2, \lambda | \mathbf{y}) = p(\mathbf{w} | \mathbf{y}, \gamma, \sigma^2, \lambda) p(\gamma, \sigma^2, \lambda | \mathbf{y}). \quad (15)$$

Similar to Tipping (2001), the distribution $p(\mathbf{w} | \mathbf{y}, \gamma, \sigma^2, \lambda)$ in Equation (15) can be calculated analytically and is again a Gaussian distribution with the following mean vector and covariance matrix:

$$\begin{aligned} \boldsymbol{\mu} &= \sigma^{-2} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{y} \\ \boldsymbol{\Sigma} &= [\sigma^{-2} \boldsymbol{\Phi}^T \boldsymbol{\Phi} + \Lambda^{-1}]^{-1} \end{aligned}$$

where $\Lambda = \text{diag}(\gamma_i \sigma^2)$. Finally, after we get point estimates for γ , σ^2 and λ , further probabilistic inference of \mathbf{w} is possible based on $p(\mathbf{w} | \mathbf{y}, \gamma, \sigma^2, \lambda)$.

To estimate γ , we can search for the local maximization with respect to the individual hyperparameters γ_i in the second term of Equation (15), by using the type-II maximum likelihood procedure described by Tipping et al. (2003). Furthermore, by using $p(\gamma, \sigma^2, \lambda | \mathbf{y}) = p(\mathbf{y}, \gamma, \sigma^2, \lambda) / p(\mathbf{y}) \propto p(\mathbf{y}, \gamma, \sigma^2, \lambda)$, we can instead maximize the following joint distribution $p(\mathbf{y}, \gamma, \sigma^2, \lambda)$ to get the type-II maximum likelihood estimation of γ . This joint distribution can be obtained by integrating out \mathbf{w} as follows:

$$\begin{aligned} p(\mathbf{y}, \gamma, \sigma^2, \lambda) &= \int p(\mathbf{y} | \mathbf{w}, \sigma^2) p(\mathbf{w} | \gamma) p(\gamma | \lambda) p(\lambda) p(\sigma^2) d\mathbf{w} \\ &= \left(\frac{1}{2\pi} \right)^{N/2} |\mathbf{C}|^{-1/2} e^{-\frac{1}{2} \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y}} p(\gamma | \lambda) p(\lambda) p(\sigma^2), \end{aligned}$$

where $\mathbf{C} = (\sigma^2 \mathbf{I}_N + \boldsymbol{\Phi} \Lambda \boldsymbol{\Phi}^T)$. Equivalent to maximizing the joint distribution is the maximization of its logarithm:

$$\begin{aligned} \mathcal{L} = & - \frac{1}{2} \log |\mathbf{C}| - \frac{1}{2} \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} + N \log \frac{\lambda}{2} - \frac{\lambda}{2} \sum_i \gamma_i \\ & + v \log \delta - \log \Gamma(v) + (v - 1) \log \lambda - \delta \lambda \\ & + a \log b - \log \Gamma(a) + (a - 1) \log \sigma^2 - b \sigma^2. \end{aligned} \quad (16)$$

In the following section, we prove that this gives a sparse model since some of the γ_i from the type-II maximum likelihood estimate will be set to zero. Thereafter, the corresponding weights and input vectors are pruned. Thus, the idea of the sparse setting is similar to the Fast Laplace method described by Babacan et al. (2010). The main difference is that, by using the conditional prior in Equation (14) for the weights in the BLS model, the criteria for letting $\gamma_i = 0$ will now also depend on σ^2 . As σ^2 is a measurement of the extent of the noise in our dataset, we expect that our method will be more flexible to the data noise.

3.2 Fast optimization algorithm

One disadvantage of the original RVM method described by Tipping (2001) is that it is computationally slow in the maximization of the type-II likelihood. The RVM begins with all the M basis functions included in the model and updates the hyperparameters iteratively. During the updates, some of the basis functions are pruned. However, the first few iterations require $O(M^3)$ computations. Faul & Tipping (2002) overcome this problem by introducing a Fast Marginal Likelihood Maximization algorithm for Sparse Bayesian Models. Instead of updating the whole

hyperparameter vector γ , only a single parameter γ_i is updated at each iteration. This fast maximization process is utilized in many sparse learning studies including the work by Babacan et al. (2010). In this paper, we also utilize this algorithm for the maximization of the log-likelihood function given in Equation (16).

In order to obtain the fast marginal likelihood solutions, the covariance matrix in the log-likelihood in Equation (16) must be decomposed as follows:

$$\begin{aligned}\mathbf{C} &= \sigma^2 \mathbf{I} + \sum_{m \neq i} \sigma^2 \gamma_m \phi_m \phi_m^T + \sigma^2 \gamma_i \phi_i \phi_i^T \\ &= \mathbf{C}_{-i} + \sigma^2 \gamma_i \phi_i \phi_i^T,\end{aligned}\quad (17)$$

where \mathbf{C}_{-i} denotes \mathbf{C} without the inclusion of basis function i . Using the Woodbury identity on the expression for the covariance matrix given in Equation (17), the inverse of the covariance matrix can be written as follows:

$$\mathbf{C}^{-1} = \mathbf{C}_{-i}^{-1} - \frac{\mathbf{C}_{-i}^{-1} \phi_i \phi_i^T \mathbf{C}_{-i}^{-1}}{\sigma^{-2} \gamma_i^{-1} + \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i}.$$

Further, by using the determinant identity, we can also obtain the following decomposition of the determinant:

$$|\mathbf{C}| = |\mathbf{C}_{-i}| |1 + \sigma^2 \gamma_i \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i|.$$

These last two expressions can then be inserted in Equation (16), and by considering only those terms that involve γ , we obtain

$$\begin{aligned}\mathcal{L}(\gamma) &= -\frac{1}{2} \left[\log |\mathbf{C}_{-i}| + \mathbf{y}^T \mathbf{C}_{-i}^{-1} \mathbf{y} + \lambda \sum_{j \neq i} \gamma_j \right] \\ &\quad + \frac{1}{2} \left[\log \frac{\sigma^2}{\sigma^2 + \gamma_i s_i} + \frac{q_i^2 \sigma^2 \gamma_i}{1 + \sigma^2 \gamma_i s_i} - \lambda \gamma_i \right] \\ &= \mathcal{L}(\gamma_{-i}) + l(\gamma_i),\end{aligned}$$

where

$$s_i = \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i, \quad \text{and} \quad q_i = \phi_i^T \mathbf{C}_{-i}^{-1} \mathbf{y}. \quad (18)$$

The log-likelihood function $\mathcal{L}(\gamma)$ has now been decomposed into $\mathcal{L}(\gamma_{-i})$, which is the marginal likelihood with ϕ_i excluded, and $l(\gamma_i)$ which contains the terms that involve γ_i . From this decomposition, we are now appropriately positioned to find the derivative of $\mathcal{L}(\gamma)$ with respect to γ_i , where the other parameters are considered as fixed:

$$\begin{aligned}\frac{d\mathcal{L}(\gamma)}{d\gamma_i} = \frac{dl(\gamma_i)}{d\gamma_i} &= \frac{1}{2} \left[-\frac{s_i}{\sigma^{-2} + \gamma_i s_i} + \frac{q_i^2 \sigma^{-2}}{(\sigma^{-2} + \gamma_i s_i)^2} - \lambda \right] \\ &= -\frac{1}{2} \left[\frac{\gamma_i^2 (s_i^2 \lambda) + \gamma_i (s_i^2 + 2\lambda \sigma^{-2} s_i) + \sigma^{-2} (\lambda \sigma^{-2} + s_i - q_i^2)}{(\sigma^{-2} + \gamma_i s_i)^2} \right].\end{aligned}$$

The numerator has a quadratic form while the denominator is always positive so that $dl(\gamma_i)/d\gamma_i = 0$ is satisfied at

$$\begin{aligned}\gamma_i &= \frac{-s_i(s_i + 2\lambda \sigma^{-2}) \pm s_i \sqrt{(s_i + 2\lambda \sigma^{-2})^2 - 4\lambda \sigma^{-2}(s_i - q_i^2 + \lambda \sigma^{-2})}}{2\lambda s_i^2} \\ &= \frac{-s_i(s_i + 2\lambda \sigma^{-2}) \pm s_i \sqrt{\Delta}}{2\lambda s_i^2},\end{aligned}\quad (19)$$

where $\Delta = (s_i + 2\lambda \sigma^{-2})^2 - 4\lambda \sigma^{-2}(s_i - q_i^2 + \lambda \sigma^{-2}) > 0$. The solution given by Equation (19) has a similar form as the expression obtained by Babacan et al. (2010). However, the terms s_i and q_i are different since we are using the covariance matrix given by Equation (17). Additionally, we also obtain the inclusion of the parameter σ^2 explicitly. When analyzing the solution given by Equation (19), we see that if $q_i^2 - s_i < \lambda \sigma^{-2}$, then $\Delta^2 < s_i + 2\lambda \sigma^{-2}$, and both solutions of Equation (19) are negative. Furthermore, since $dl(\gamma_i)/d\gamma_i$ evaluated at $\gamma_i = 0$ is negative, the maximum occurs at $\gamma_i = 0$. In the other situation, when $q_i^2 - s_i > \lambda \sigma^{-2}$, there are two real solutions of Equation (19), one

negative and one positive. Since $dl(\gamma_i)/d\gamma_i$ is positive when evaluated at $\gamma_i = 0$ and negative at $\gamma_i = \infty$, the positive solution from Equation (19) maximizes $l(\gamma_i)$ as well as $\mathcal{L}(\gamma)$. The maximum of $\mathcal{L}(\gamma)$, when holding the remaining components fixed, is therefore obtained at:

$$\gamma_i = \begin{cases} \frac{-s_i(s_i + 2\lambda\sigma^{-2}) + s_i\sqrt{(s_i + 2\lambda\sigma^{-2})^2 - 4\lambda\sigma^{-2}(s_i - q_i^2 + \lambda\sigma^{-2})}}{2\lambda s_i^2} & \text{if } q_i^2 - s_i > \lambda\sigma^{-2} \\ 0 & \text{otherwise.} \end{cases} \quad (20)$$

To estimate λ in Equation (20), we can take the derivative of Equation (16) with respect to λ and set it to zero:

$$\lambda = \frac{2(N + v - 1)}{\sum_i \gamma_i + 2\delta}. \quad (21)$$

In order to use the above estimate, we have to either estimate or set the values for the parameters v and δ . Park & Casella (2008)'s first method to estimate λ is to find an empirical Bayes estimate of λ by using the Gibbs sampler. In the second method, they use a hyperprior for λ and obtain the estimate given in Equation (21) by setting values for the parameters v and δ based on the value of λ from the first method (Park & Casella 2008, Ch. 5.2).

However, the maximum likelihood estimates of v and δ can actually be calculated from the following expressions (Choi & Wette 1969):

$$\delta = \frac{r}{\bar{\lambda}}, \quad (22)$$

$$\log v = \log \bar{\lambda} - \overline{\log \lambda} + \psi(v). \quad (23)$$

It is anyhow not possible to find an analytical solution of v from Equation (23). Hence, v must be estimated numerically. From the expression in Equation (23), we also see that we need more than one observation of λ in order to estimate v . If we derive the estimates from the log-likelihood given in Equation (16), we would obtain similar expressions. However, the term $\log \bar{\lambda} - \overline{\log \lambda}$ in Equation (23) would disappear with only one λ value. What we have done during implementation is to simulate a small sample of λ by using the Gibbs sampler as described by Park & Casella (2008) and Casella (2001) to get the initial values of v and δ . After a couple of iterations we obtain a new sample of estimated λ s from Equation (21). These new λ s are then used in the Equations (22) and (23) to update v and δ , which again are used in Equation (21) to update λ . Finally, in the simulation section, we also tried to use the Gibbs estimate of λ directly instead of using Equation (21). However, we obtain the best results by using the above procedure and get a stable value of λ . Thus, in the simulation study, we only present the result where λ is estimated based on the procedure developed from Equation (21) to (23).

In the optimization algorithm, we also have to update the expressions for s_i and q_i in Equation (18). Instead of computing and updating s_i and q_i directly, one can first find the values of

$$S_i = \phi_i^\top \mathbf{C}^{-1} \phi_i, \quad Q_i = \phi_i^\top \mathbf{C}^{-1} \mathbf{y}, \quad (24)$$

and from Equation (24), we can obtain:

$$s_i = \frac{S_i}{1 - \gamma_i \sigma^2 S_i}, \quad q_i = \frac{Q_i}{1 - \gamma_i \sigma^2 S_i}.$$

Notice that when γ_i is set to zero, we get $s_i = S_i$ and $q_i = Q_i$. Further, the Woodbury identity will be used on Equation (24) such that S_i and Q_i can be calculated from:

$$S_i = \sigma^{-2} \phi_i^\top \phi_i - \sigma^{-2} \phi_i^\top \phi \Sigma \phi^\top \phi_i \sigma^{-2}, \quad (25)$$

$$Q_i = \sigma^{-2} \phi_i^\top \mathbf{y} - \sigma^{-2} \phi_i^\top \phi \Sigma \phi^\top \mathbf{y} \sigma^{-2}, \quad (26)$$

where Σ and ϕ contain only those basis functions that are currently included in the model. This computation is therefore much faster compared to if we had started out with all the N basis functions. Based on these results, we obtain Algorithm 1.

In the first step of Algorithm 1, we need to fix σ^2 to a sensible value so that the Fast algorithm can be utilized. As Faul & Tipping (2002) mentioned, it is only with fixed σ^2 that the quantities in Equations (25) and (26) can be calculated via an update formulae. The alternative is to find an estimate from the log-likelihood given in Equation (16) as we did for λ . However, Babacan et al. (2010) argued that the estimation of σ^2 based on iterations will be unreliable and the result will be significantly affected. Thus, they fix σ^2 as the initialized value in the whole algorithm, same as Faul & Tipping (2002) and the above mentioned Algorithm 1.

Algorithm 1 The Noise-Robust Fast Sparse Bayesian Learning Model

```

1: Fix  $\sigma^2$  to some sensible value (e.g.  $0.01 \|\mathbf{y}\|$ )
2: Initialize all  $\gamma_i = 0$ ,  $\lambda = 0$ 
3: while convergence criteria are not met, do
4:   Choose a  $\gamma_i$ 
5:   if  $q_i^2 - s_i > \lambda \sigma^{-2}$  and  $\gamma_i = 0$  then
6:     Add  $\gamma_i$  to the model
7:   else if  $q_i^2 - s_i > \lambda \sigma^{-2}$  and  $\gamma_i > 0$ , then
8:     Re-estimate  $\gamma_i$ 
9:   else if  $q_i^2 - s_i < \lambda \sigma^{-2}$ , then
10:    Prune  $i$  from the model (set  $\gamma_i = 0$ )
11:   end if
12:   Update  $\Sigma$  and  $\mu$ 
13:   Update  $s_i$  and  $q_i$ 
14:   Update  $\lambda$  using Equation (21)
15:   Update  $\delta$  and  $v$  using Equations (22) and (23)
16: end while

```

From Equation (20) and Algorithm 1, we see that the criteria for setting $\gamma_i = 0$ depends on λ and the variance term σ^2 . However, in order to see the true relation, we need to rewrite q_i and s_i since both also include σ^2 . The expressions for s_i and q_i are given by Equation (18) where we find \mathbf{C}_{-i} from Equation (17). By using the vector notation, we can also write \mathbf{C}_{-i} as:

$$\begin{aligned}
\mathbf{C}_{-i} &= (\sigma^2 \mathbf{I} + \Phi_{-i} \Lambda_{-i} \Phi_{-i}^\top) \\
&= \sigma^2 (\mathbf{I} + \Phi_{-i} \tilde{\Lambda}_{-i} \Phi_{-i}^\top) \\
&= \sigma^2 \tilde{\mathbf{C}}_{-i},
\end{aligned}$$

where Φ_{-i} is the $N \times N - 1$ design matrix where basis function i is removed, $\tilde{\Lambda}_{-i}$ is the diagonal matrix Λ where the single element γ_i and the component σ^2 on the diagonal are removed, and $\tilde{\mathbf{C}}_{-i}$ denotes \mathbf{C}_{-i} where the component σ^2 is excluded. We can now decompose s_i and q_i to obtain

$$s_i = \sigma^{-2} \tilde{s}_i \quad q_i = \sigma^{-2} \tilde{q}_i$$

where $\tilde{s}_i = \phi_i^\top \tilde{\mathbf{C}}_{-i}^{-1} \phi_i$ and $\tilde{q}_i = \phi_i^\top \tilde{\mathbf{C}}_{-i} \mathbf{y}$. The inequality from Equation (20), which decides when $\gamma_i = 0$, can now be decomposed as follows:

$$\begin{aligned}
q_i^2 - s_i &\leq \lambda \sigma^{-2}, \\
(\sigma^{-2} \tilde{q}_i)^2 - \sigma^{-2} \tilde{s}_i &\leq \lambda \sigma^{-2}, \\
\sigma^{-2} \tilde{q}_i^2 - \tilde{s}_i &\leq \lambda.
\end{aligned} \tag{27}$$

Hence, we see that as $\sigma^2 \rightarrow \infty$, γ_i will be set to zero since $\lambda > 0$, and the inequality from Equation (27) must hold. Thus, we can see that in BLS, the information of σ^2 is utilized to adjust the number of zero hyperparameters when we

estimate γ , otherwise, the noisy information might often be confused with the real signal information such that only a small proportion of the γ_i s is set to zero. In the RVM and the Fast Laplace method, the variance information of the random error can not be extracted as explicitly as in the BLS method, which often leads to more dense models.

3.3 Prediction

After the convergence of the learning Algorithm 1, we end up with L ($L < N$) non-zero γ_i 's and each of them correspond to a "relevance basis function" and a related "relevance input vector" from the training data. Let γ_{MP} denote the vector that contains those L non-zero γ_i 's. For a new input data \mathbf{x}^* , we can now make predictions based on the posterior of the weights conditioning on γ_{MP} and $\hat{\sigma}^2$. The predictive distribution for the output y^* can be approximated by

$$p(y^*|\mathbf{y}, \gamma_{MP}, \hat{\sigma}^2) = \int p(y^*|\mathbf{w}, \gamma_{MP}, \hat{\sigma}^2)p(\mathbf{w}|\mathbf{y}, \gamma_{MP}, \hat{\sigma}^2) d\mathbf{w}.$$

This distribution is analytically tractable and is also Gaussian with the following predictive mean and predictive variance:

$$\begin{aligned} y^* &= \boldsymbol{\mu}_{MP}^\top \boldsymbol{\phi}(\mathbf{x}^*), \\ \sigma^{2*} &= \hat{\sigma}^2 + \boldsymbol{\phi}(\mathbf{x}^*)^\top \boldsymbol{\Sigma}_{MP} \boldsymbol{\phi}(\mathbf{x}^*), \end{aligned}$$

where $\boldsymbol{\mu}_{MP}$ and $\boldsymbol{\Sigma}_{MP}$ are calculated by

$$\boldsymbol{\mu}_{MP} = \hat{\sigma}^{-2} \boldsymbol{\Sigma}_{MP} \boldsymbol{\Phi}_{MP}^\top \mathbf{y}, \quad (28)$$

$$\boldsymbol{\Sigma}_{MP} = [\hat{\sigma}^{-2} \boldsymbol{\Phi}_{MP}^\top \boldsymbol{\Phi}_{MP} + \Lambda_{MP}^{-1}]^{-1}. \quad (29)$$

Here, $\boldsymbol{\phi}(\mathbf{x}^*) = [\phi_1(\mathbf{x}^*), \dots, \phi_L(\mathbf{x}^*)]^\top$; $\phi_j(\mathbf{x}^*) = K(\mathbf{x}^*, \mathbf{x}_j)$, $j = 1, \dots, L$, where \mathbf{x}_j is the j 'th relevance input vector among the total L relevance input vectors. Furthermore, in Equations (28) and (29), $\boldsymbol{\Phi}_{MP} = [\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_L]$ is the $N \times L$ design matrix whose column vectors are $\boldsymbol{\phi}_j = [\phi_j(\mathbf{x}_1), \dots, \phi_j(\mathbf{x}_N)]^\top$; $j = 1, \dots, L$. Moreover, the estimated covariance matrix $\Lambda_{MP} = \text{diag}(\hat{\sigma}^2 \gamma_{MP})$ is an $L \times L$ matrix. Thus, $\boldsymbol{\mu}_{MP}$ and $\boldsymbol{\Sigma}_{MP}$, which are separately the estimated posterior mean vector and covariance matrix over the weight, contain only L non-zero elements that correspond to those non-zero elements in γ_{MP} . In practice, the predictive mean can be used as a point prediction, and the predictive variance can be used to construct the prediction interval.

4 The Simulation Results

In this section, we compare the proposed BLS method with the RVM described by Tipping et al. (2003) and the Fast Laplace method by Babacan et al. (2010) by using several examples with a synthetic dataset, where the input variables vary from one dimension to multi-dimensions. Since both Tipping et al. (2003) and Babacan et al. (2010) utilize the fast marginal likelihood maximization to estimate the hyperparameters, we call the method described by Tipping et al. (2003) as FRVM and the method depicted by Babacan et al. (2010) as FLAP.

4.1 One dimensional input – two cases

We first consider the simplest situation where the dimension of the input variable is $D = 1$. Two cases of different forms of $f(\cdot)$ are utilized. Case 1 uses the Sinc function with $f(x) = \sin(x)/x$, a benchmark function that is frequently utilized to evaluate how the kernel-based learning methods perform (Vapnik et al. 1997, Tipping 2001, Schmolck & Everson 2007). As the choice of hyperparameters in the kernel function also affects the result, we first utilize a hyperparameter-free univariate "linear spline" kernel as the basis function. The univariate linear spline kernel has the

following representation:

$$K(x_m, x_n) = 1 + x_m x_n + x_m x_n \min(x_m, x_n) - \frac{x_m + x_n}{2} \min(x_m, x_n)^2 + \frac{x_m + x_n}{3} \min(x_m, x_n)^3.$$

As the Sinc function generates a very smooth signal, the linear spline kernel can approximate smooth functions without a serious overfitting problem (Schmolck & Everson 2007). To compare the performance of BLS, FRVM and FLAP when data are exposed to different extents of noise, we set ϵ as the zero-mean Gaussian noise with standard deviation at different levels. Figure 1 illustrates the approximation results from the three methods with $\sigma = 0.1, 0.5$ and 1 and $N = 200$. The black dots indicate the training data from the model $\mathbf{y} = f(x) + \epsilon$ where the x values lie within the interval $[-10, 10]$. We generate the same training dataset for all the three methods for the same value of σ . The green line corresponds to the Sinc function $f(x)$, which we call the signal function. The blue line is the approximation $\mathbf{y}^* = \Phi_{MP} \boldsymbol{\mu}_{MP}$, where $\boldsymbol{\mu}_{MP}$ from Equation (28) is the mean of the posterior for \mathbf{w} , which only contains L non-zero weight estimations. The location of the non-zero weighted input vectors are the red circles. The values in $\boldsymbol{\mu}_{MP}$ varies when using different methods and different σ , which result in different shapes of the blue line in Figure 1.

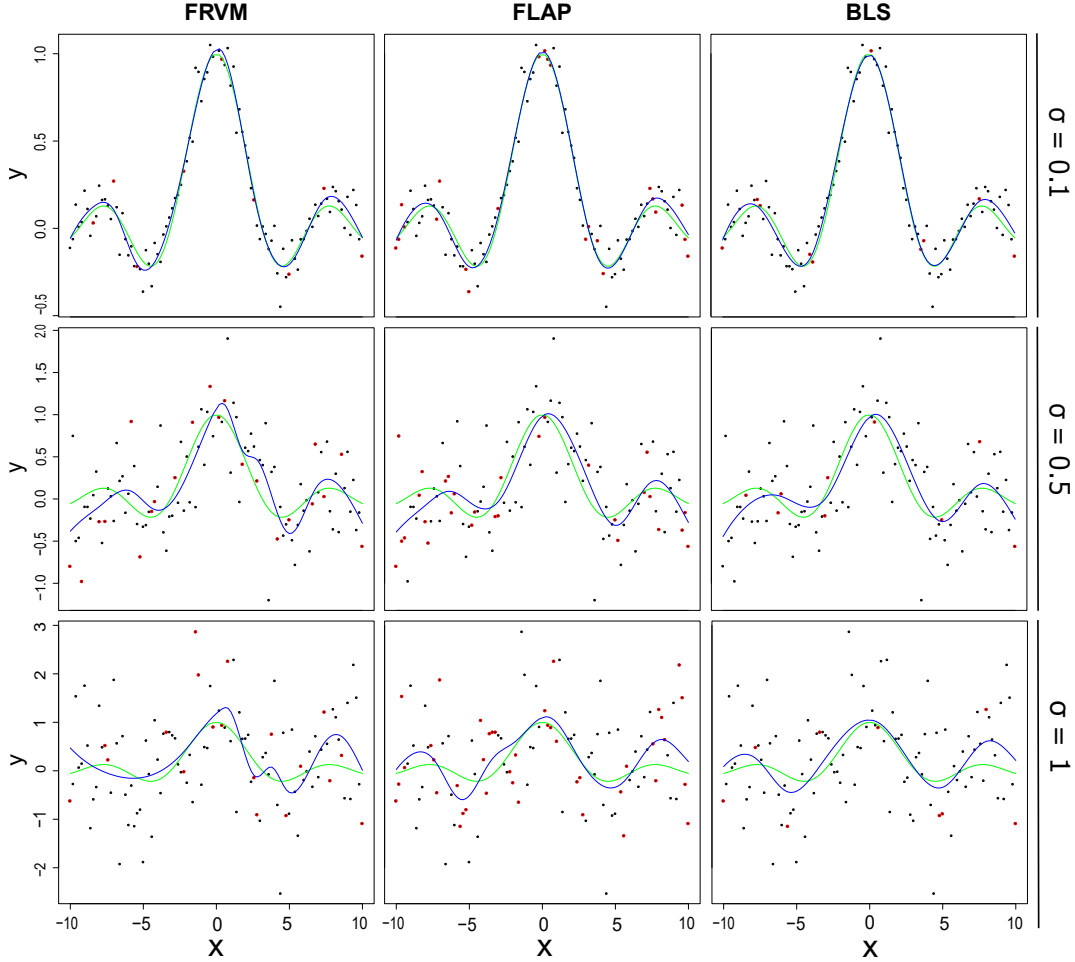


Figure 1: The Sinc function (green line) and its reconstruction (blue line) from the training data that are generated for different values of σ . The red dots are the relevance vectors and the black dots are the remaining training data. The same data generation is used for all methods to better compare them.

Figure 1 shows that when $\sigma = 0.1$, the approximations from all three methods almost overlap with the original signal. The locations of the relevance vectors are mainly toward the end and the turning points of the signals. This means that only the most informative input vectors are utilized in approximation. When σ increases to 1, the BLS method can still capture the general form of the original signal, while the RVM produces a rougher approximation. The FLAP method provides a good fit but uses significantly more relevance vectors for the approximation compared to the BLS method.

We next run 100 data generations for each learning method with $\sigma = 0.01, 0.1, 0.3, 0.5$ and 1 and utilize several criteria to compare the approximation results for the three methods. For each generation, we get the value of the mean squared error $\|y_* - f(x)\|^2/N$ and the number of relevant vectors L . Based on those 100 generations, we use MSE and NOV to denote the average of the mean squared error values and the number of relevance vectors, respectively, while using $\text{sd}(\text{MSE})$ and $\text{sd}(\text{NOV})$ to denote the sample standard deviation. We summarize the simulation result for Case 1 in Table 1. Notice that we have rounded off to three decimal places, which means that the zeros represent values lower than this threshold.

Table 1: 1D function Case 1

σ	NOV			sd(NOV)			MSE			sd(MSE)		
	FRVM	FLAP	BLS	FRVM	FLAP	BLS	FRVM	FLAP	BLS	FRVM	FLAP	BLS
0.01	7.63	21.03	9.95	1.228	0.771	0.77	0	0	0	0	0	0
0.1	18.79	21.22	9.71	6.753	1.936	1.038	0.001	0.001	0.001	0	0	0
0.3	22.76	26.77	9.44	3.245	5.035	1.157	0.013	0.009	0.008	0.005	0.004	0.003
0.5	22.23	28.76	9.06	3.084	6.808	1.144	0.034	0.026	0.021	0.012	0.012	0.009
1	22.16	28.51	8.5	4.616	9.093	1.642	0.143	0.082	0.086	0.053	0.039	0.037

From the first row in Table 1, we see that when σ is small, all three methods perform quite well with very low MSE. When σ increases, the BLS obtains the lowest values in almost all the different criteria. The average number of relevance vectors is also quite stable for the BLS method for all values of σ . We see that it is equal to 9.95 when $\sigma = 0.01$ and decreases slightly to 8.5 when $\sigma = 1$. The trend for the other two methods is however the opposite, where the number of relevance vectors increases when the value of σ increases. This is due to the noisy information that is confused with the real signal in the approximations. For the BLS method, the number of relevance vectors is adjusted accordingly by including σ explicitly in the learning algorithm. From Equation (27), we see that more relevance vectors are set to zero as σ^2 increases. Furthermore, Figure 1 also illustrates that even when $\sigma=1$, the location of those relevance vectors of the BLS are centered around the most informative places: the turning points of the curve and the end. The other two methods have more relevance vectors which are spread over the whole curve.

We also estimate λ for FLAP and BLS based on 100 data generations, and get the following table where “Mean λ ” denotes the average value of the estimate and “sd λ ” is the sample standard deviation.

Table 2: Estimates of λ for 1D function Case 1

σ	Mean λ		sd λ	
	FLAP	BLS	FLAP	BLS
0.01	567.08	19.174	31.8	0.014
0.1	620.183	19.042	138.954	0.058
0.3	568.127	19.079	242.224	0.155
0.5	615.83	19.376	477.004	0.285
1	17647.283	19.541	147887.951	0.354

Thus, compared with FLAP, BLS gives a stable estimate for λ , and it seems to not depend on the value of σ . Additionally, this may even be explained by the fact that the BLS method can extract the variance information explicitly in the

learning process and the estimation result will be flexible to the extent of the variance, while the other two methods can not manage this extraction. We can say that generally, the BLS method is sparser and more stable, especially for higher values of σ .

The dataset generated by the Sinc function has a smooth shape. Next we investigate how the three methods perform on a dataset with genuine high frequency and several bumps. Case 2 in this section uses the Bumps data (Donoho & Johnstone 1994) where the signal consists of zero values, except for a few non-zero bumps, and the x values again lie within the interval $[-10, 10]$. The values of the output y are the values of the Bumps data in (Donoho & Johnstone 1994) along with the values of the random noise. In Case 2, the linear spline kernel is not suitable as it will seriously smoothen all the bumps. Hence, we choose the Gaussian kernel $K(x_m, x_n) = \exp(-r^{-2}||x_m - x_n||^2)$ where the width hyperparameter r is set to 0.1. The results for all methods are shown in Figure 2 for the variances $\sigma = 0.1, 0.5$ and 1. The same data generation is used for all methods to better compare them.

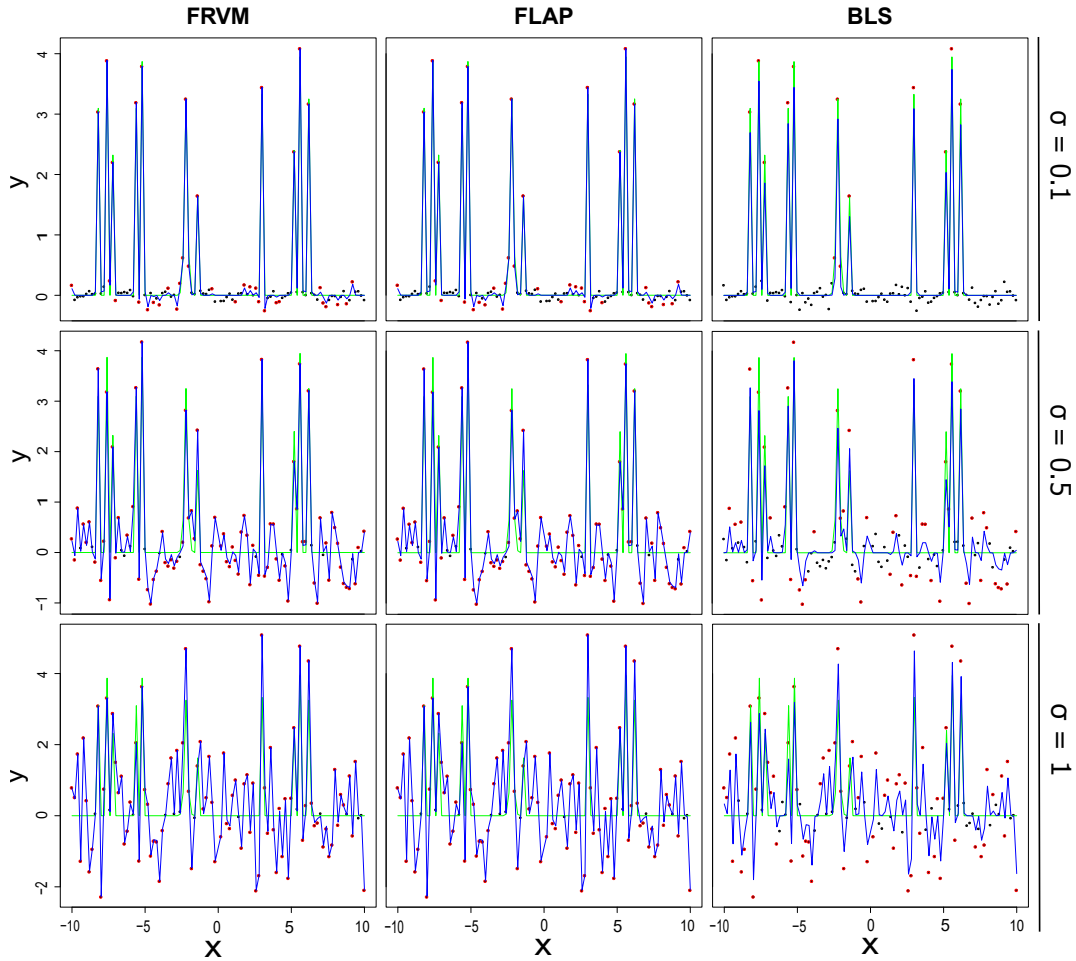


Figure 2: The Bump data (green line) and its reconstruction (blue line) plotted for different values of σ . The red dots are the relevance vectors and the black dots are the remaining training data. The same data generation is used for all methods to better compare them.

Figure 2 shows that when σ is 0.1, the locations of the relevance vectors for BLS are mostly the locations of the bump points which include the most important “spike” information, while FRVM and FLAP use more relevance vectors that do not lie at the bumps. However, from Figure 2, we also see that the BLS tends to smoothen out the signal in this situation since the prediction does not reach the top of the bumps. When σ increases, the output values become

increasingly “polluted” by the noise and the predictions of all models show small spikes at the places where the original noise free data is zero. Figure 2 shows that all the three methods tend to use more relevance vectors as the noise increases. In other words, when σ increases, the noisy data points mimic the pattern of bumps and all the three sparse learning methods “misunderstand” the noisy values as the bump signal and use more relevance vectors. However, BLS still uses lesser relevance vectors when compared to the other two methods. This is better seen in Table 3 which shows the NOV based on the average of 100 generations. In addition, the table shows the $\text{sd}(\text{NOV})$, MSE and $\text{sd}(\text{MSE})$.

Table 3: 1D function Case 2

σ	NOV			sd(NOV)			MSE			sd(MSE)		
	FRVM	FLAP	BLS	FRVM	FLAP	BLS	FRVM	FLAP	BLS	FRVM	FLAP	BLS
0.01	15.92	15.82	13	0.307	0.386	0	0.001	0.001	0.016	0	0	0
0.1	41.6	41.86	13.05	4.66	4.797	0.297	0.005	0.005	0.018	0.001	0.001	0.002
0.3	76.11	76.68	31.12	3.527	4.447	4.105	0.078	0.076	0.034	0.012	0.013	0.008
0.5	84.53	83.7	51.1	3.41	3.658	4.554	0.236	0.228	0.101	0.036	0.035	0.026
1	90.19	89.17	68.36	3.014	2.671	4.672	0.946	0.964	0.538	0.144	0.138	0.109

Table 3 shows that for the smaller values of the standard deviation, $\sigma = 0.01$ and 0.1 , all three methods give quite precise approximations, however, BLS performs slightly worse than the other two methods with slightly higher MSE. When σ is 0.3 , 0.5 and 1 , we see that the BLS performs better than FRVM and FLAP with a lower MSE and fewer relevance vectors.

The results from Case 1 and Case 2 for one dimensional input data show that, generally, the BLS produces a more stable approximation for higher values of σ . The BLS method is not only sparser but also more robust to the noise and can provide better resolutions when the underlying signal is “polluted” by the random noise.

4.2 Multidimensional input

We now apply other synthetic examples to compare all three methods with multidimensional input where $D = 2$ and $D = 5$. For $D = 2$, we again use two cases for comparison. Case 1 has the following Data Generating Process (DGP):

$$y_i = f(\mathbf{x}_i) + \epsilon = \frac{\sin(x_{i1})}{x_{i1}} + \frac{\sin(x_{i2})}{x_{i2}} + \epsilon, \quad (30)$$

where the first dimension’s input is x_{i1} while the second dimension’s input is x_{i2} ; $i = 1, \dots, N$, and they are both uniformly spaced in $[-5, 5]$ with an interval length of 0.3 . We first use Figure 3 to illustrate the shape of the training data generated from Equation (30) with $\sigma = 0, 0.1, 0.5$ and 1 , where $\sigma = 0$ corresponds to the data without noise and shows the true shape of the function $f(\mathbf{x})$.

Figure 3 shows that when $\sigma = 0.1$, we can still see the original shape of the function, while when $\sigma = 1$, the data become very noisy and lose the shape of the original image. The following multivariate Gaussian kernel function with $D = 2$ is utilized with width hyperparameters $r_1 = r_2 = 0.1$:

$$K(\mathbf{x}_m, \mathbf{x}_n) = \prod_{d=1}^D \exp(-r_d^{-2} \|x_{md} - x_{nd}\|^2).$$

Figure 4 shows the approximation of the function $f(\mathbf{x})$ by using the three methods for $\sigma = 0.1, 0.5$ and 1 , with the red points being the estimates obtained from the relevance input vectors. Figure 4 shows again that when $\sigma = 0.1$, all the three methods give out quite good approximations and can capture the shape of the original function well. However, when σ is 1 , FRVM gives a rougher approximation, while FLAP and BLS give closer approximations. The number of

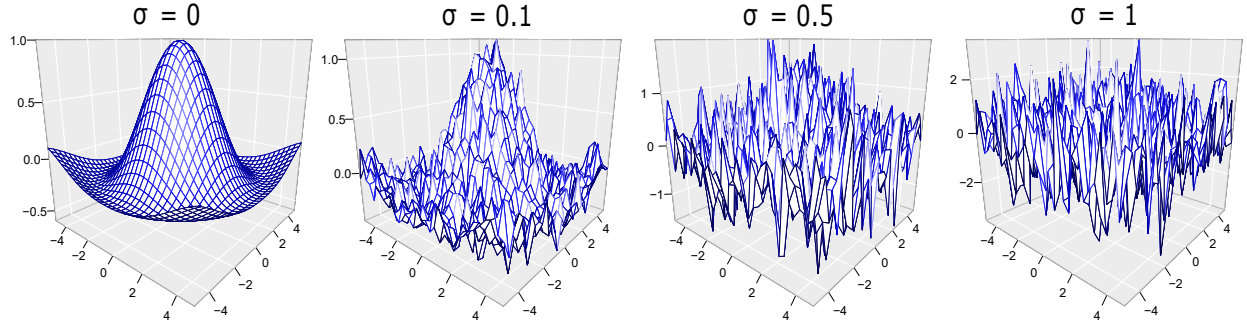


Figure 3: Training data generated from the two dimensional Sinc function with different values of σ . The figure to the left shows the true shape of the function, and the remainder plots shows the function with additive noise.

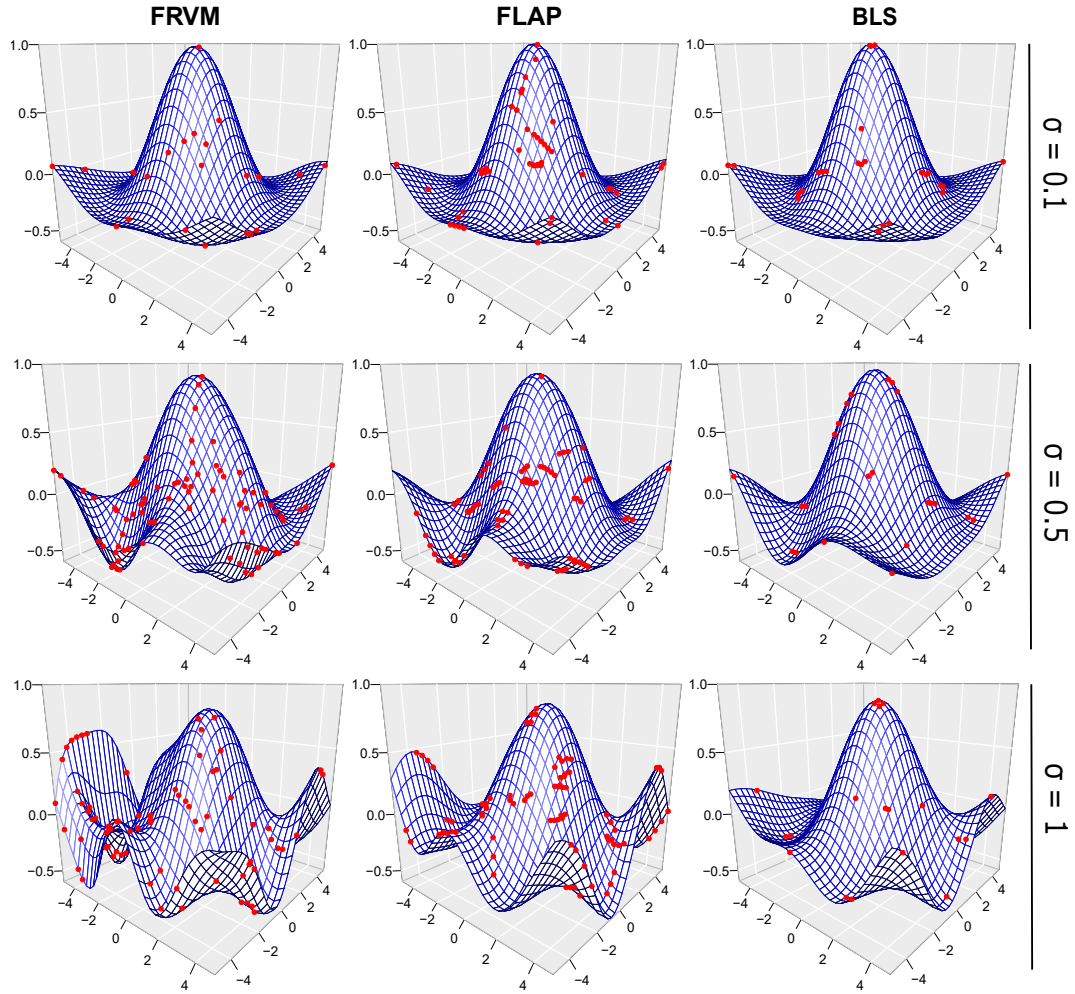


Figure 4: The reconstructions of the two dimensional Sinc function from the training data plotted in Figure 3. The red points are the estimated outputs from the relevance input vectors.

relevance vectors are fewest for the BLS method for all values of σ , and the location of the relevance vectors are at the most informative points such as the top and curving places of the surface.

To compare the results based on different criteria, we carry out 100 generations and show the key NOV, MSE, sd(NOV), and sd(MSE) for the three methods in Table 4. The random error's standard deviation is set as $\sigma = 0.1, 0.5, 1, 2$ and 3 . The motivation for setting a higher σ than for $D = 1$ is that one source of the random error comes from the measurement error in input variable and, generally, variance to the random error is higher for higher dimensional input cases.

Table 4: 2D function Case 1

σ	NOV			sd(NOV)			MSE			sd(MSE)		
	FRVM	FLAP	BLS	FRVM	FLAP	BLS	FRVM	FLAP	BLS	FRVM	FLAP	BLS
0.1	21.8	48.27	29.76	2.412	9.053	5.578	0	0	0	0	0	0
0.5	88.63	84.88	24.7	15.982	19.604	4.239	0.01	0.007	0.004	0.002	0.002	0.001
1	91.47	97.77	21.51	14.811	22.972	3.307	0.04	0.027	0.016	0.007	0.008	0.005
2	88.68	103.26	20.71	15.094	28.765	3.033	0.163	0.109	0.058	0.04	0.032	0.02
3	87.79	99.17	19.97	17.045	32.162	3.07	0.36	0.246	0.128	0.07	0.079	0.045

The second case for $D = 2$ is a combination of the two examples in the one-dimensional situation, and the data is generated from:

$$y_i = x_{i1} + \frac{\sin(x_{i2})}{x_{i2}} + \epsilon, \quad (31)$$

where x_{i1} is the same Bumps data with most zero values and few non-zero bumps, while x_{i2} , $i = 1, 2, \dots, 20$ are uniformly spaced in $[-10, 10]$ with an interval length of 1. The results for $\sigma = 0.1, 0.5, 1, 2$ and 3 , using the data generation given by Equation (31), are shown in Table 5.

Table 5: 2D function Case 2

σ	NOV			sd(NOV)			MSE			sd(MSE)		
	FRVM	FLAP	BLS	FRVM	FLAP	BLS	FRVM	FLAP	BLS	FRVM	FLAP	BLS
0.1	19.04	29.34	20.03	2.666	4.854	2.794	0.173	0.174	0.174	0.004	0.004	0.004
0.5	27.91	34.01	25.03	7.679	6.206	4.416	0.19	0.188	0.186	0.019	0.022	0.021
1	47.82	41.52	27.96	13.813	9.803	5.335	0.259	0.224	0.221	0.047	0.043	0.042
2	53.61	52.5	29.3	10.858	14.616	5.212	0.522	0.409	0.345	0.11	0.11	0.109
3	52.74	52.59	29.94	11.747	13.88	6.94	0.932	0.649	0.557	0.215	0.189	0.179

The BLS method has the lowest value for almost all four criteria in both the 2D test cases, shown in Table 4 and 5. Thus, these cases indicate that the BLS method performs better than the FRVM and FLAP methods. When σ increases, BLS outperforms the other two methods with a significantly smaller number of relevance vectors and a lower MSE. Furthermore, Table 4 shows the same trend as Table 1, where the number of relevance vectors for BLS decrease when σ increases. In the first 2D test case the NOV approaches a stable value around 20, which is the number of relevance vectors for FRVM when $\sigma = 0.1$. The number of relevance vectors for BLS in Table 5 increases since the one-dimensional value of the input variable comes from the Bumps data. However, the BLS method still uses the fewest number of relevance vectors and has the least problems of overfitting.

Finally, we show an example where the dimension of the input vectors is $D = 5$. The DGP is from Friedman et al. (1991), with $N = 400$, where the inputs are five-dimensional data where each dimension corresponds to an independent variable uniformly distributed within the interval $[0, 1]$:

$$y_i = 10 \sin(\pi x_{i1} x_{i2}) + 20(x_{i3} - 0.5)^2 + 10x_{i4} + 5x_{i5} + \epsilon.$$

Here, if we use the Gaussian Kernel, then we need to choose five width parameters, one for each dimension. Thus, we use the multivariate linear spline kernel with $D = 5$:

$$K(\mathbf{x}_m, \mathbf{x}_n) = \prod_{d=1}^D \left(1 + x_{md}x_{nd} + x_{md}x_{nd} \min(x_{md}, x_{nd}) - \frac{x_{md} + x_{nd}}{2} \min(x_{md}, x_{nd})^2 + \frac{x_{md} + x_{nd}}{3} \min(x_{md}, x_{nd})^3 \right).$$

The result is presented in Table 6. Again, we see that the BLS gives the lowest MSE when σ is 1, 2, and 3. Although it is not much lower than FLAP, it is significantly better than FRVM. Moreover, it also uses the fewest relevance vectors and has the lowest standard error on the number of relevance vectors among all three methods.

Table 6: 5D function with Linear spline kernel

σ	NOV			sd(NOV)			MSE			sd(MSE)		
	FRVM	FLAP	BLS	FRVM	FLAP	BLS	FRVM	FLAP	BLS	FRVM	FLAP	BLS
0.1	29.74	40.09	20.64	7.076	4.262	1.872	0.21	0.986	0.901	0.097	0.22	0.152
0.5	31.83	41.36	20.79	8.321	4.416	2.095	0.311	1.03	0.888	0.101	0.209	0.183
1	37.1	41.39	20.89	8.113	4.948	1.938	0.583	1.2	1.073	0.105	0.25	0.186
2	50.77	42.25	20.53	6.241	5.237	1.872	2.079	1.767	1.807	0.385	0.31	0.319
3	51.71	43.82	19.72	5.286	5.813	1.804	4.867	2.879	2.949	0.868	0.655	0.674

5 Conclusion

In this paper, we propose a new sparse Bayesian learning method and compare this method with the wellknown RVM (Tipping 2001) as well as the method proposed by Babacan et al. (2010), where the latter is used within the field of compressive sensing. The proposed method combines the hierarchical Bayesian framework as described by Park & Casella (2008) with the estimation and inference process proposed by Tipping (2001). The prior distribution of the weight parameters in our paper is conditional on the variance of the random error, which leads to the posterior distribution of the weights parameter being adjusted by the extent of noise. By this conditional prior distribution, our method becomes more sparse and robust to the dataset that is polluted by the high variance noise. We demonstrate analytically how the sparsity and robustness can be achieved when using the type-II maximum likelihood method to estimate the hyperparameters. A fast optimization algorithm is utilized in the maximization of the type-II likelihood so that the learning process is effective, and the whole algorithm is illustrated in detail. We carry out a comprehensive study in stimulation where the dataset is suffered from various extents of noise. Our simulation results show that the method proposed in this paper is both sparse and stable to the high variance noise. In implementations, our method is especially suitable to carry out signal reconstructions, when dataset is large and noisy.

References

- Agarwal, A. & Triggs, B. (2005), ‘Recovering 3d human pose from monocular images’, *IEEE transactions on pattern analysis and machine intelligence* **28**(1), 44–58.
- Andrews, D. F. & Mallows, C. L. (1974), ‘Scale mixtures of normal distributions’, *Journal of the Royal Statistical Society: Series B (Methodological)* **36**(1), 99–102.
- Ashburner, J. (2007), ‘A fast diffeomorphic image registration algorithm’, *Neuroimage* **38**(1), 95–113.
- Babacan, S. D., Molina, R. & Katsaggelos, A. K. (2010), ‘Bayesian compressive sensing using laplace priors’, *IEEE Transactions on Image Processing* **19**(1), 53–63.

- Balakrishnan, S. & Madigan, D. (2009), ‘Priors on the variance in sparse bayesian learning: the demi-bayesian lasso’, *Frontiers of Statistical Decision Making and Bayesian Analysis: In Honor of James O. Berger* pp. 346–359.
- Bishop, C. M. (2006), *Pattern recognition and machine learning*, springer.
- Boser, B. E., Guyon, I. M. & Vapnik, V. N. (1992), A training algorithm for optimal margin classifiers, in ‘Proceedings of the fifth annual workshop on Computational learning theory’, ACM, pp. 144–152.
- Candes, E., Romberg, J. & Tao, T. (2004), ‘Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information’, *arXiv preprint math/0409186*.
- Casella, G. (2001), ‘Empirical bayes gibbs sampling’, *Biostatistics* **2**(4), 485–500.
- Choi, S. C. & Wette, R. (1969), ‘Maximum likelihood estimation of the parameters of the gamma distribution and their bias’, *Technometrics* **11**(4), 683–690.
- Demir, B. & Erturk, S. (2007), ‘Hyperspectral image classification using relevance vector machines’, *IEEE Geoscience and Remote Sensing Letters* **4**(4), 586–590.
- Donoho, D. L. & Johnstone, J. M. (1994), ‘Ideal spatial adaptation by wavelet shrinkage’, *biometrika* **81**(3), 425–455.
- Donoho, D. L. et al. (2006), ‘Compressed sensing’, *IEEE Transactions on information theory* **52**(4), 1289–1306.
- Efron, B., Hastie, T., Johnstone, I., Tibshirani, R. et al. (2004), ‘Least angle regression’, *The Annals of statistics* **32**(2), 407–499.
- Faul, A. C. & Tipping, M. E. (2002), Analysis of sparse bayesian learning, in ‘Advances in neural information processing systems’, pp. 383–389.
- Friedman, J. H. et al. (1991), ‘Multivariate adaptive regression splines’, *The annals of statistics* **19**(1), 1–67.
- Friedman, J., Hastie, T. & Tibshirani, R. (2001), *The elements of statistical learning*, Vol. 1, Springer series in statistics New York.
- Ghosh, S. & Mujumdar, P. P. (2008), ‘Statistical downscaling of gcm simulations to streamflow using relevance vector machine’, *Advances in water resources* **31**(1), 132–146.
- Ji, S., Xue, Y., Carin, L. et al. (2008), ‘Bayesian compressive sensing’, *IEEE Transactions on signal processing* **56**(6), 2346.
- Krishnapuram, B., Carin, L., Figueiredo, M. A. & Hartemink, A. J. (2005), ‘Sparse multinomial logistic regression: Fast algorithms and generalization bounds’, *IEEE transactions on pattern analysis and machine intelligence* **27**(6), 957–968.
- MacKay, D. J. C. (1992), ‘The evidence framework applied to classification networks’, *Neural Computation* **4**, 720–736.
- Neal, R. M. (2012), *Bayesian learning for neural networks*, Vol. 118, Springer Science & Business Media.
- Park, T. & Casella, G. (2008), ‘The bayesian lasso’, *Journal of the American Statistical Association* **103**(482), 681–686.
- Rasmussen, C. E. & Quinonero-Candela, J. (2005), Healing the relevance vector machine through augmentation, in ‘Proceedings of the 22nd international conference on Machine learning’, ACM, pp. 689–696.
- Schmolck, A. & Everson, R. (2007), ‘Smooth relevance vector machine: a smoothness prior extension of the rvm’, *Machine Learning* **68**(2), 107–135.
- Schölkopf, B. (2001), The kernel trick for distances, in ‘Advances in neural information processing systems’, pp. 301–307.
- Schölkopf, B., Burges, C. J., Smola, A. J. et al. (1999), *Advances in kernel methods: support vector learning*, MIT press.

- Seeger, M. (2000), ‘Relationships between gaussian processes, support vector machines and smoothing splines’, *Machine Learning* .
- Seeger, M. W. & Nickisch, H. (2008), Compressed sensing and bayesian experimental design, in ‘Proceedings of the 25th international conference on Machine learning’, ACM, pp. 912–919.
- Smola, A. J., Schölkopf, B. & Müller, K.-R. (1998), ‘The connection between regularization operators and support vector kernels’, *Neural networks* **11**(4), 637–649.
- Tibshirani, R. (1996), ‘Regression shrinkage and selection via the lasso’, *Journal of the Royal Statistical Society: Series B (Methodological)* **58**(1), 267–288.
- Tipping, M. E. (2001), ‘Sparse bayesian learning and the relevance vector machine’, *Journal of machine learning research* **1**(Jun), 211–244.
- Tipping, M. E., Faul, A. C. et al. (2003), Fast marginal likelihood maximisation for sparse bayesian models., in ‘AISTATS’.
- Vapnik, V., Golowich, S. E. & Smola, A. J. (1997), Support vector method for function approximation, regression estimation and signal processing, in ‘Advances in neural information processing systems’, pp. 281–287.
- Williams, C. K. & Rasmussen, C. E. (2006), *Gaussian processes for machine learning*, Vol. 2, MIT press Cambridge, MA.
- Wipf, D. P. & Rao, B. D. (2004), ‘Sparse bayesian learning for basis selection’, *IEEE Transactions on Signal processing* **52**(8), 2153–2164.
- Wipf, D., Palmer, J., Rao, B. & Kreutz-Delgado, K. (2007), Performance analysis of latent variable models with sparse priors, in ‘Proceedings of ICASSP 2007’.