

A Bayesian Regression Model and Applications

Yijun Yu

Department of Mathematics, Tuskegee University, Tuskegee, AL, USA

Email: yyu@tuskegee.edu

How to cite this paper: Yu, Y.J. (2020) A Bayesian Regression Model and Applications. *Journal of Applied Mathematics and Physics*, 8, 1877-1887.

<https://doi.org/10.4236/jamp.2020.89141>

Received: July 17, 2020

Accepted: September 19, 2020

Published: September 22, 2020

Copyright © 2020 by author(s) and Scientific Research Publishing Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

A sparse vector regression model is developed. The model is established by employing Bayesian formulation and trained by using a set of data $D = \{\mathbf{x}_n, \mathbf{t}_n\}_{n=1}^N$. The parameters needed to be determined in the algorithm are reduced by a special prior hyperparameter setting, and therefore the algorithm is simpler than similar type of Bayesian vector regression models. The examples of applications to the function approximation and inverse scattering problem are presented.

Keywords

Bayesian, Regression, Applications

1. Introduction

There has been a lot of interest in studying the Bayesian vector regression and its application on various classification and regression problems [1] [2] [3] [4]. The Bayesian approach considers probability distributions with the observed data; prior distributions are converted to posterior distribution through the use of Bayes' theorem. Let \mathbf{x} be an input vector and \mathbf{t} be a vector of target parameters. In a regression formulation our goal is to define a model $\mathbf{y}(\mathbf{x}; \mathbf{w})$ that yields an approximation to the true target \mathbf{t} , with the model defined by the parameters \mathbf{w} . The model is typically designed using a set of "training" data $D = \{\mathbf{x}_n, \mathbf{t}_n\}_{n=1}^N$. Although we initially consider a finite set D , the goal is for the subsequent model $\mathbf{y}(\mathbf{x}; \mathbf{w})$ to be applicable to arbitrary $(\mathbf{x}, \mathbf{t}) \notin D$, over the anticipated range of \mathbf{t} . When developing a regression model one must address the bias-variance tradeoff. A bias is introduced by restricting the form that $\mathbf{y}(\mathbf{x}; \mathbf{w})$ may take, while the variance represents the error between the model $\mathbf{y}(\mathbf{x}; \mathbf{w})$ and true target parameters \mathbf{t} . Models with minimal bias typically have significant flexibility, and therefore the model parameters may vary significantly as a function of the spe-

cific training set \mathcal{D} employed. To obtain good model generalization, which may be connected to the variation in the model parameters as a function of \mathcal{D} , one must introduce a bias. The utilization of a small number of non-zero parameters \mathbf{w} often yields a good balance between bias and variance; such models are termed “sparse”. This has led to development of the relevance vector machine [5].

The rest of this paper is organized as follows. The theory of the vector-regression formulation is presented in Section 2, with application example provided in Section 3. The work is summarized in Section 4.

2. Sparse Bayesian Vector Regression

2.1. Model Specification

Assume we have available a set of training data $\mathcal{D} = \{\mathbf{x}_n, \mathbf{t}_n\}_{n=1}^N$, where $\mathbf{x}_n = [x_n^{(1)} \ x_n^{(2)} \ \cdots \ x_n^{(L)}]^\top$ and $\mathbf{t}_n = [t_n^{(1)} \ t_n^{(2)} \ \cdots \ t_n^{(M)}]^\top$. Our objective is to develop a function $\mathbf{y}(\mathbf{x}; \mathbf{w})$ that is dependent on the parameters \mathbf{w} . After $\mathbf{y}(\mathbf{x}; \mathbf{w})$ is so designed, it may be used to map an arbitrary \mathbf{x} to an approximation of the target parameters \mathbf{t} .

The specific vector-regression function

$\mathbf{y}(\mathbf{x}; \mathbf{w}) = [y^{(1)}(\mathbf{x}; \mathbf{w}) \ y^{(2)}(\mathbf{x}; \mathbf{w}) \ \cdots \ y^{(M)}(\mathbf{x}; \mathbf{w})]^\top$ employed here is defined as

$$\mathbf{y}(\mathbf{x}; \mathbf{w}) = \sum_{i=1}^N w_i \mathbf{t}_i K(\mathbf{x}, \mathbf{x}_i) + \mathbf{w}_0 \quad (1)$$

where $\mathbf{w}_0 = [w_0^{(1)} \ w_0^{(2)} \ \cdots \ w_0^{(M)}]^\top$, and $K(\mathbf{x}, \mathbf{x}_i)$ is a kernel function that is designed such that $K(\mathbf{x}, \mathbf{x}_i)$ is large if $\mathbf{x}_i \approx \mathbf{x}$ and otherwise $K(\mathbf{x}, \mathbf{x}_i)$ is small. Hence in (1) only those $\mathbf{x}_i \approx \mathbf{x}$ are important in defining $\mathbf{y}(\mathbf{x}; \mathbf{w})$.

Let

$$\mathbf{w} = [w_1 \ w_2 \ \cdots \ w_N \ w_0^{(1)} \ w_0^{(2)} \ \cdots \ w_0^{(M)}]^\top,$$

$$\boldsymbol{\psi}_i(\mathbf{x}) = [\phi_i^{(1)} \ \phi_i^{(2)} \ \cdots \ \phi_i^{(M)}]^\top, \quad i = 1, 2, \dots, N$$

with

$$\phi_i^{(k)} = t_i^{(k)} K(\mathbf{x}, \mathbf{x}_i), \quad i = 1, 2, \dots, N; \quad k = 1, 2, \dots, M \quad (2)$$

and $M \times (N + M)$ matrix

$$\boldsymbol{\Psi}(\mathbf{x}) = [\boldsymbol{\psi}_1(\mathbf{x}) \ \boldsymbol{\psi}_2(\mathbf{x}) \ \cdots \ \boldsymbol{\psi}_N(\mathbf{x}) \ \mathbf{I}_M], \quad (3)$$

where \mathbf{I}_M is $M \times M$ identity matrix, then (1) can be expressed in matrix form

$$\mathbf{y}(\mathbf{x}; \mathbf{w}) = \boldsymbol{\Psi}(\mathbf{x}) \mathbf{w} \quad (4)$$

Assume that target is from the model with additive noise

$$\mathbf{t} = \mathbf{y}(\mathbf{x}; \mathbf{w}) + \boldsymbol{\varepsilon} = \boldsymbol{\Psi}(\mathbf{x}) \mathbf{w} + \boldsymbol{\varepsilon}, \quad (5)$$

where model error $\boldsymbol{\varepsilon} = [\varepsilon^{(1)} \ \varepsilon^{(2)} \ \cdots \ \varepsilon^{(M)}]^\top$ and $\varepsilon^{(k)}, k = 1, 2, \dots, M$ are inde-

pendent samples from a zero-mean Gaussian process with variance α_0^{-1}

$$p(\varepsilon^{(k)}) = N(\varepsilon^{(k)} | 0, \alpha_0^{-1}), \quad k = 1, 2, \dots, M \tag{6}$$

We therefore have

$$\begin{aligned} p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \alpha_0) &= \left(\frac{2\pi}{\alpha_0}\right)^{\frac{M}{2}} \exp\left(-\frac{\alpha_0}{2} \|\mathbf{t} - \Psi(\mathbf{x})\mathbf{w}\|_2^2\right) \\ &= N(\mathbf{t} | \Psi(\mathbf{x})\mathbf{w}, \alpha_0^{-1}\mathbf{I}_M) \end{aligned} \tag{7}$$

We wish to constrain the weights \mathbf{w} such that a simple model is favored, this accomplished by invoking a prior distribution on \mathbf{w} that favors most of the weights being zero. In this context, only the most relevant members of the training set $\mathbf{D} = \{\mathbf{x}_n, \mathbf{t}_n\}_{n=1}^N$, those with nonzero weights w_n , are ultimately used in the final regression model. This simplicity allows improved regression performance for $(\mathbf{x}, \mathbf{t}) \notin \mathbf{D}$ [5] [6].

We employ a zero-mean Gaussian prior distribution for \mathbf{w}

$$p(\mathbf{w} | \alpha_0, \alpha) = N(\mathbf{w} | \mathbf{0}_{N+M}, \alpha_0^{-1}\alpha^{-1}\mathbf{I}_{N+M}), \tag{8}$$

where $\mathbf{0}_{N+M}$ is a $(N+M)$ -dimensional zero vector, \mathbf{I}_{N+M} is a $(N+M) \times (N+M)$ identity matrix, and suitable priors over hyperparameters α_0 and α are Gamma distributions [7]

$$p(\alpha_0 | a, b) = \text{Gamma}(\alpha_0 | a, b) \tag{9}$$

$$p(\alpha | c, d) = \text{Gamma}(\alpha | c, d) \tag{10}$$

where $\text{Gamma}(\alpha_0 | a, b) = \Gamma(a)^{-1} b^a \alpha_0^{a-1} e^{-b\alpha_0}$ with $\Gamma(a) = \int_0^\infty t^{a-1} e^{-t} dt$.

The hierarchical prior over \mathbf{w} favors a sparse model and the prior over α_0 will be used to favor small model error on the training data \mathbf{D} .

2.2. Inference

For training data $\mathbf{D} = \{\mathbf{x}_n, \mathbf{t}_n\}_{n=1}^N$ we introduce LN -dimensional vector

$$\mathbf{X} = [\mathbf{x}_1^\top \ \mathbf{x}_2^\top \ \dots \ \mathbf{x}_N^\top]^\top$$

and MN -dimensional vector

$$\mathbf{T} = [\mathbf{t}_1^\top \ \mathbf{t}_2^\top \ \dots \ \mathbf{t}_N^\top]^\top$$

and let $(MN) \times (M+N)$ matrix

$$\Phi = [\Phi_1^\top \ \Phi_2^\top \ \dots \ \Phi_N^\top]^\top \quad \text{with} \quad \Phi_i = \Psi(\mathbf{x}_i), \quad i = 1, 2, \dots, N,$$

then by (7), we have

$$\begin{aligned} p(\mathbf{T} | \mathbf{w}, \alpha_0, \mathbf{X}) &= \left(\frac{2\pi}{\alpha_0}\right)^{\frac{MN}{2}} \exp\left(-\frac{\alpha_0}{2} \|\mathbf{T} - \Phi\mathbf{w}\|_2^2\right) \\ &= N(\mathbf{T} | \Phi\mathbf{w}, \alpha_0^{-1}\mathbf{I}_{MN}) \end{aligned} \tag{11}$$

Noting that $p(\mathbf{T} | \alpha_0, \alpha, \mathbf{X}) = \int p(\mathbf{T} | \mathbf{w}, \alpha_0, \mathbf{X}) p(\mathbf{w} | \alpha_0, \alpha) d\mathbf{w}$ is a convolu-

tion of Gaussians, the posterior distribution over the weights \mathbf{w} can be derived as

$$p(\mathbf{w} | \alpha_0, \alpha, \mathbf{X}, \mathbf{T}) = \frac{p(\mathbf{T} | \mathbf{w}, \alpha_0, \mathbf{X}) p(\mathbf{w} | \alpha_0, \alpha)}{p(\mathbf{T} | \alpha_0, \alpha, \mathbf{X})} = N(\mathbf{w} | \boldsymbol{\mu}, \alpha_0^{-1} \boldsymbol{\Sigma}) \quad (12)$$

where

$$\boldsymbol{\Sigma} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \alpha \mathbf{I}_{M+N})^{-1} = \left(\sum_{i=1}^N \boldsymbol{\Phi}_i^T \boldsymbol{\Phi}_i + \alpha \mathbf{I}_{M+N} \right)^{-1} \quad (13)$$

$$\boldsymbol{\mu} = \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{T} = \boldsymbol{\Sigma} \sum_{i=1}^N (\boldsymbol{\Phi}_i^T t_i) \quad (14)$$

2.3. Hyperparameter Optimization

We determine α in (13) by maximizing $p(\alpha | \mathbf{T}, \mathbf{X}) \propto p(\mathbf{T} | \alpha, \mathbf{X}) p(\alpha)$ with respect to α . It is equivalent to maximize the \ln of this quantity. In addition, we can choose to maximize with respect to $\ln \alpha$ as we can assume hyperpriors over a logarithmic scale.

Since

$$\begin{aligned} & \ln p(\mathbf{T} | \alpha, \mathbf{X}) \\ &= \ln \int p(\mathbf{T} | \mathbf{w}, \alpha_0, \mathbf{X}) p(\mathbf{w} | \alpha_0, \alpha) p(\alpha_0 | a, b) d\mathbf{w} d\alpha_0 \\ &= -\frac{1}{2} \left[\ln |\mathbf{B}| + (MN + 2a) \ln (\mathbf{T}^T \mathbf{B}^{-1} \mathbf{T} + 2b) \right] + const \end{aligned}$$

where $\mathbf{B} = \mathbf{I}_{MN} + \alpha^{-1} \boldsymbol{\Phi} \boldsymbol{\Phi}^T$, and $p(\ln \alpha) = \alpha p(\alpha)$, we obtain objective function

$$L(\alpha) = -\frac{1}{2} \left[\ln |\mathbf{B}| + (MN + 2a) \ln (\mathbf{T}^T \mathbf{B}^{-1} \mathbf{T} + 2b) \right] + c \ln \alpha - d\alpha \quad (15)$$

By the determinant identity [8], we have

$$\begin{aligned} |\mathbf{B}| &= \left| \mathbf{I}_{MN} + \alpha^{-1} \boldsymbol{\Phi} \boldsymbol{\Phi}^T \right| \\ &= \alpha^{-(M+N)} \left| \alpha \mathbf{I}_{M+N} + \boldsymbol{\Phi}^T \boldsymbol{\Phi} \right| \\ &= \alpha^{-(M+N)} \left| \boldsymbol{\Sigma}^{-1} \right|, \end{aligned}$$

and so

$$\ln |\mathbf{B}| = -(M + N) \ln \alpha + \ln \left| \boldsymbol{\Sigma}^{-1} \right| \quad (16)$$

Using the Woodbury formula, we obtain

$$\begin{aligned} \mathbf{B}^{-1} &= \left(\mathbf{I}_{MN} + \alpha^{-1} \boldsymbol{\Phi} \boldsymbol{\Phi}^T \right)^{-1} \\ &= \mathbf{I}_{MN} - \boldsymbol{\Phi} \left(\alpha \mathbf{I}_{M+N} + \boldsymbol{\Phi}^T \boldsymbol{\Phi} \right)^{-1} \boldsymbol{\Phi}^T \\ &= \mathbf{I}_{MN} - \boldsymbol{\Phi} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T, \end{aligned}$$

thus

$$\begin{aligned} \mathbf{T}^T \mathbf{B}^{-1} \mathbf{T} &= \mathbf{T}^T \left(\mathbf{T} - \boldsymbol{\Phi} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{T} \right) \\ &= \mathbf{T}^T \left(\mathbf{T} - \boldsymbol{\Phi} \boldsymbol{\mu} \right) \end{aligned} \quad (17)$$

$$= \|\mathbf{T}\|^2 - \mathbf{T}^T \boldsymbol{\Phi} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{T} \quad (18)$$

Then by (16) and Jacobi's formula, we have

$$\begin{aligned} \frac{d \ln |\mathbf{B}|}{d \ln \alpha} &= -(M+N) + \frac{1}{|\boldsymbol{\Sigma}^{-1}|} \frac{d|\boldsymbol{\Sigma}^{-1}|}{d \ln \alpha} \\ &= -(M+N) + \text{tr} \left(\boldsymbol{\Sigma} \frac{d\boldsymbol{\Sigma}^{-1}}{d \ln \alpha} \right) \\ &= -(M+N) + \alpha \sum_{j=1}^{M+N} \boldsymbol{\Sigma}_{jj} \end{aligned} \tag{19}$$

where $\boldsymbol{\Sigma}_{jj}$ is the j -th diagonal element of matrix $\boldsymbol{\Sigma}$.

By (18)

$$\begin{aligned} \frac{d\mathbf{T}^T \mathbf{B}^{-1} \mathbf{T}}{d \ln \alpha} &= -\frac{d\mathbf{T}^T \boldsymbol{\Phi} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{T}}{d \ln \alpha} \\ &= -\mathbf{T}^T \boldsymbol{\Phi} \frac{d\boldsymbol{\Sigma}}{d \ln \alpha} \boldsymbol{\Phi}^T \mathbf{T} \\ &= -\mathbf{T}^T \boldsymbol{\Phi} \boldsymbol{\Sigma} \frac{d\boldsymbol{\Sigma}^{-1}}{d \ln \alpha} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{T} \\ &= \alpha \|\boldsymbol{\mu}\|^2 \end{aligned} \tag{20}$$

Using (17), (19) and (20), we have

$$\begin{aligned} \frac{dL(\alpha)}{d\alpha} &= \frac{1}{2} \left(M+N - \alpha \sum_{j=1}^{M+N} \boldsymbol{\Sigma}_{jj} \right) - \frac{(MN+2a)}{2(\mathbf{T}^T \mathbf{B}^{-1} \mathbf{T} + 2b)} \frac{d\mathbf{T}^T \mathbf{B}^{-1} \mathbf{T}}{d \ln \alpha} + c - d\alpha \\ &= \frac{1}{2} \left(M+N - \alpha \sum_{j=1}^{M+N} \boldsymbol{\Sigma}_{jj} \right) - \frac{(MN+2a)\|\boldsymbol{\mu}\|^2 \alpha}{2[\mathbf{T}^T (\mathbf{T} - \boldsymbol{\Phi} \boldsymbol{\mu}) + 2b]} + c - d\alpha \end{aligned} \tag{21}$$

Setting (21) to zero, followed by algebra operations, yield

$$\alpha = \frac{M+N+2c}{\sum_{j=1}^{M+N} \boldsymbol{\Sigma}_{jj} + 2d + (MN+2a)\|\boldsymbol{\mu}\|^2 / [\mathbf{T}^T (\mathbf{T} - \boldsymbol{\Phi} \boldsymbol{\mu}) + 2b]} \tag{22}$$

The algorithm consists of (13), (14) and (22) with iteration for $\alpha, \boldsymbol{\Sigma}$ and $\boldsymbol{\mu}$.

2.4. Making Predictions

Assume α_{MP} and $\alpha_{0,MP}$ are maximizing values obtained by maximizing $p(\alpha | \mathbf{T}, \mathbf{X})$ (Sec. 2.3) and $p(\alpha_0 | \mathbf{T}, \mathbf{X})$, respectively. Assume

$$p(\alpha_0, \alpha | \mathbf{X}, \mathbf{T}) \approx \delta(\alpha_0 - \alpha_{0,MP}) \delta(\alpha - \alpha_{MP})$$

then

$$\begin{aligned} p(\mathbf{t} | \mathbf{x}, \mathbf{X}, \mathbf{T}) &= \int p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \alpha_0, \alpha) p(\mathbf{w}, \alpha_0, \alpha | \mathbf{X}, \mathbf{T}) d\mathbf{w} d\alpha_0 d\alpha \\ &= \int p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \alpha_0) p(\mathbf{w} | \alpha_0, \alpha, \mathbf{X}, \mathbf{T}) p(\alpha_0, \alpha | \mathbf{X}, \mathbf{T}) d\mathbf{w} d\alpha_0 d\alpha \\ &\approx \int p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \alpha_0) p(\mathbf{w} | \alpha_0, \alpha, \mathbf{X}, \mathbf{T}) \delta(\alpha_0 - \alpha_{0,MP}) \delta(\alpha - \alpha_{MP}) d\mathbf{w} d\alpha_0 d\alpha \\ &= \int p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \alpha_{0,MP}) p(\mathbf{w} | \alpha_{0,MP}, \alpha_{MP}, \mathbf{X}, \mathbf{T}) d\mathbf{w} \\ &= N \left(\mathbf{t} | \mathbf{y}(\mathbf{x}; \boldsymbol{\mu}), (\alpha_{0,MP})^{-1} \boldsymbol{\Omega} \right) \end{aligned} \tag{23}$$

with

$$\mathbf{y}(\mathbf{x}; \boldsymbol{\mu}) = \boldsymbol{\Psi}(\mathbf{x})\boldsymbol{\mu} \quad (24)$$

$$\boldsymbol{\Omega} = \mathbf{I}_M + \boldsymbol{\Psi}(\mathbf{x})\boldsymbol{\Sigma}\boldsymbol{\Psi}(\mathbf{x})^T \quad (25)$$

3. Applications

In examples we employ a radial-basis-function kernel

$K(x, x_i) = \exp(-\|x - x_i\|^2 / r^2)$, and just parameters a, b, c and d by training and testing on given training data, finally we take $a = b = c = d = 0.05$ for all examples in this section. In all figures the horizontal axis is the index of samples and the vertical axis is output.

3.1. Regression: Function Approximation

The model can be used to establish the relation between independent variables and dependent variables of a function.

Example 1 2-dimensional vector function with two variables

$$t_1 = \operatorname{sinc}\left(\frac{x_1 + x_2}{4}\right)$$

$$t_2 = -0.5\operatorname{sinc}\left(\frac{x_1 + x_2}{4}\right)\sin\left(\frac{x_1 x_2}{20}\right) - 0.4$$

in domain $\{(x_1, x_2) \mid -10 \leq x_1 \leq 10, 0 \leq x_2 \leq 20\}$, where $\operatorname{sinc}(x) = \sin(x)/x$.

Figure 1 and **Figure 2** illustrate the results. **Figure 1** is learning from 100 noise-free training samples. **Figure 2** is based on 100 noisy training samples. The noise is generated from zero-mean Gaussian with 5% of average training data $\|t\|$ as standard deviation. Both test on 100 examples that are not in training data.

Example 2 3-dimensional vector function with 200 variables

$$(x_1, x_2, \dots, x_{200}) \rightarrow (t_1, t_2, t_3).$$

$$t_1 = \sum_{k=1}^{200} \sin\left((x_k)^{5/7}\right) + \frac{x_{50}}{100}$$

$$t_2 = \frac{x_{200}}{800}t_1 + \frac{x_{50}}{200} + \cos\left(\frac{x_{100}}{5}\right) - 10$$

$$t_3 = \operatorname{atan}\left(\frac{t_1 + t_2}{6}\right) + \frac{t_2 - t_1}{2} - 10$$

We choose samples at point $\mathbf{x}^n = (x_1^n, x_2^n, \dots, x_{200}^n)$ with $x_k^n = k + (n-1)\pi/4$. 100 samples at points \mathbf{x}^n with $n = 1, 3, 5, \dots, 199$ used as training data, and 100 samples at points \mathbf{x}^n with $n = 2, 4, 6, \dots, 200$ used as testing data.

Figure 3 is from noise-free training samples. **Figure 4** is based on noisy training samples. The noise is generated from zero-mean Gaussian with 5% of average training data $\|t\|$ as standard deviation.

3.2. Regression: Inverse Scattering

The model can be used to characterize the connection between measured vector

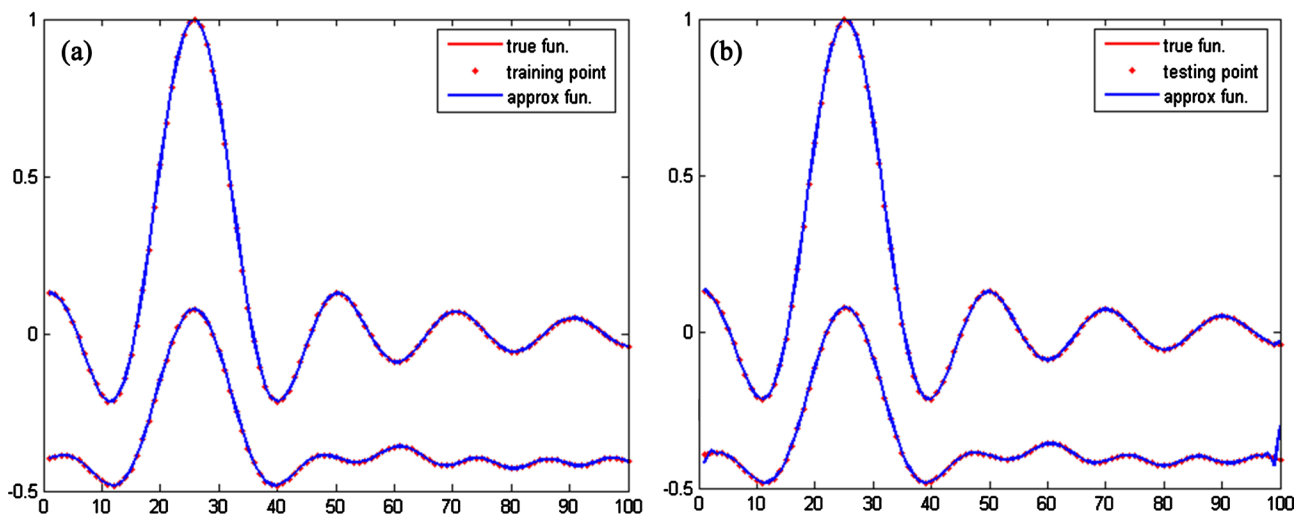


Figure 1. Results for 2-dim vector function with noise-free data: (a) predict on training points; (b) predict on testing points.

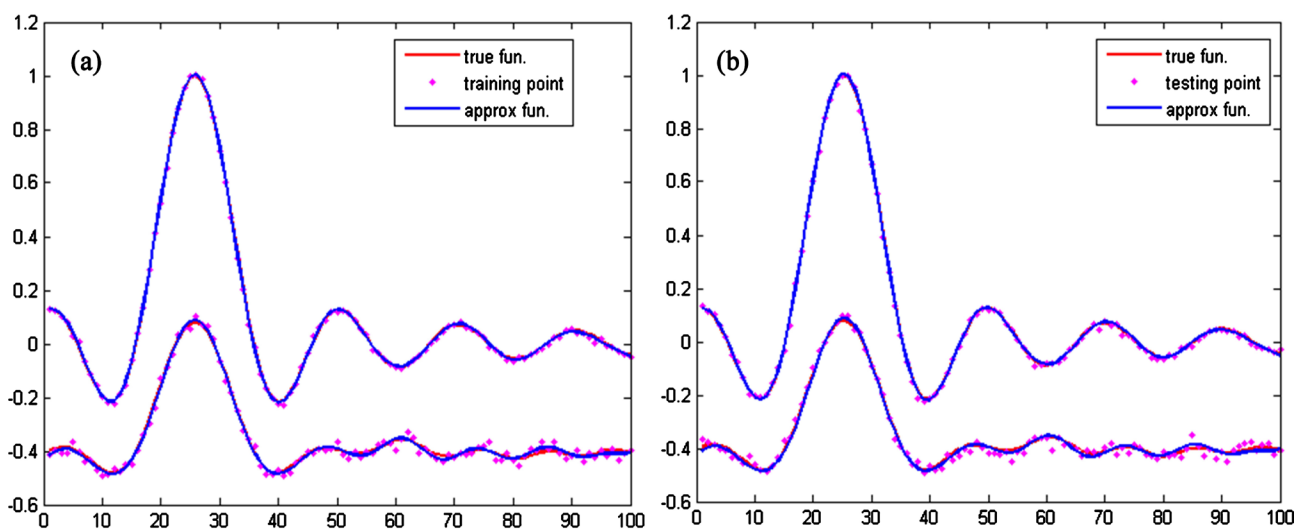


Figure 2. Results for 2-dim vector function with noisy data: (a) predict on training points; (b) predict on testing points.

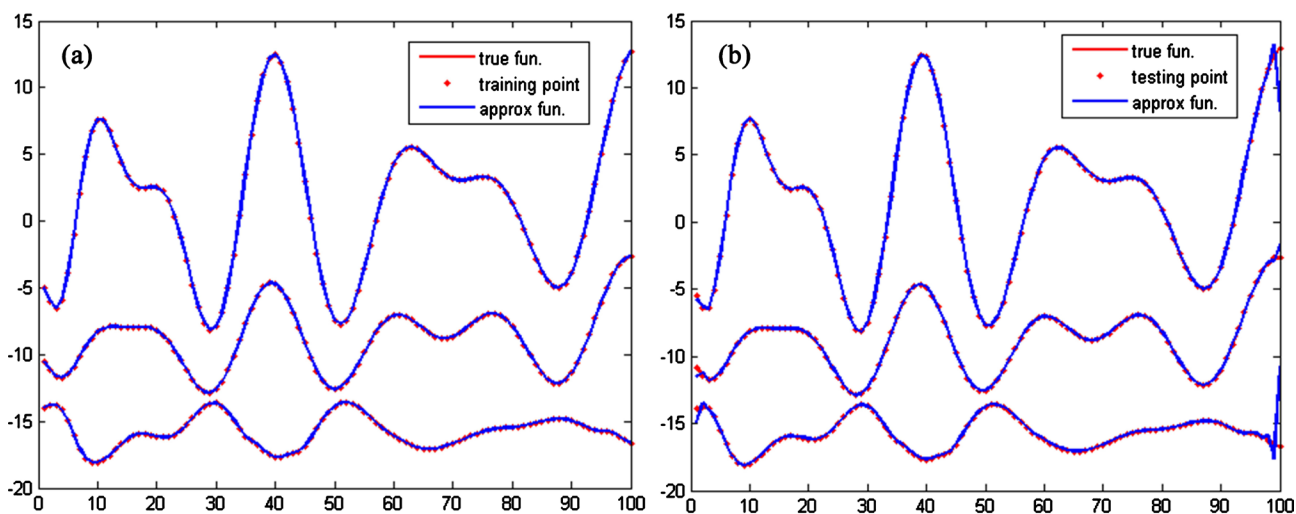


Figure 3. Results for 3-dim vector function with noise-free data: (a) predict on training points; (b) predict on testing points.

scattered-field data \mathbf{x} and the underlying target responsible for these fields, characterized by the parameter vector \mathbf{t} . The scattering data \mathbf{x} may be measured at multiple positions. In the examples the measure data is simulated by forward model.

We consider a homogeneous lossless dielectric target buried in a lossy dielectric half space. The objective is to invert for the parameters of the target. In the examples, the parameter vector \mathbf{t} is composed of three real numbers: the depth of target, the size of target, and the dielectric constant of target. For each target there are 100 simulated measure data. Training data $\mathbf{D} = \{\mathbf{x}_n, \mathbf{t}_n\}_{n=1}^N$ is composed of $N = 180$ examples and testing data is composed of 125 examples that are not in \mathbf{D} .

Example 1 We consider cube target in this example. **Figure 5** and **figure 6** illustrate the results. **Figure 5** is from noise-free data. **Figure 6** is based on noisy data. The noise is generated from zero-mean Gaussian with 10% of average training data $\|\mathbf{x}\|$ as standard deviation. The “size” is the width of cube.

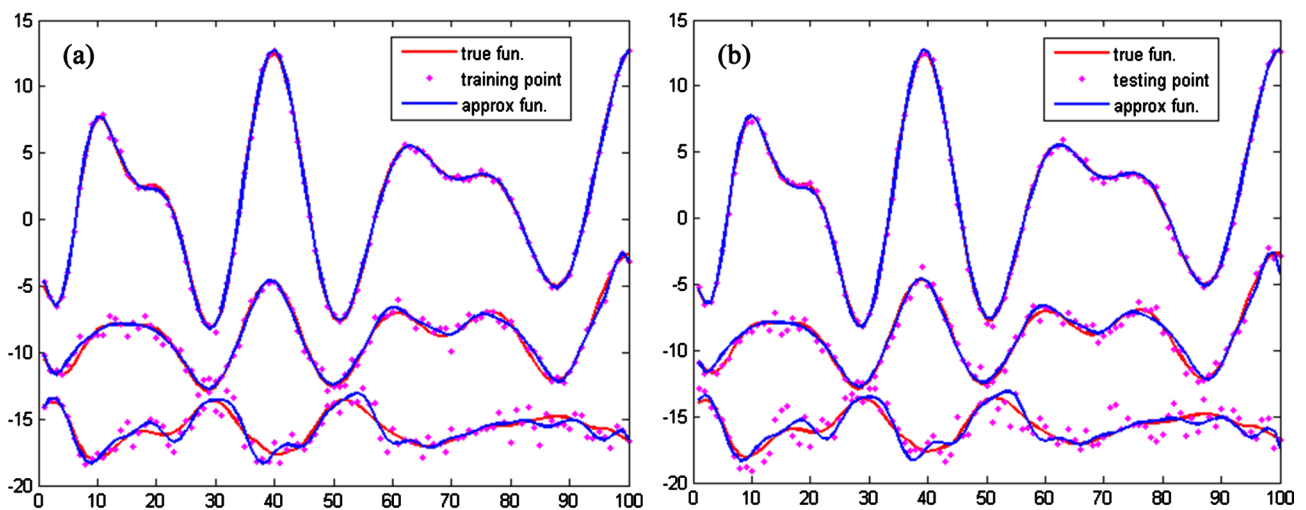


Figure 4. Results for 3-dim vector function with noisy data: (a) predict on training points; (b) predict on testing points.

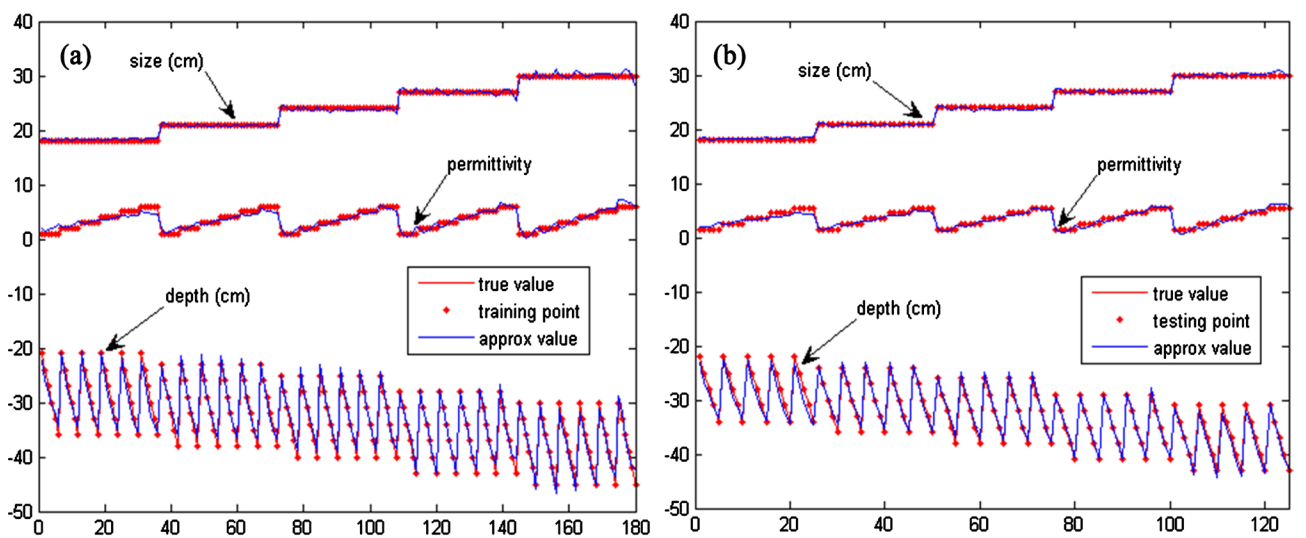


Figure 5. Results for cube target with noise-free data: (a) predict on training points; (b) predict on testing points.

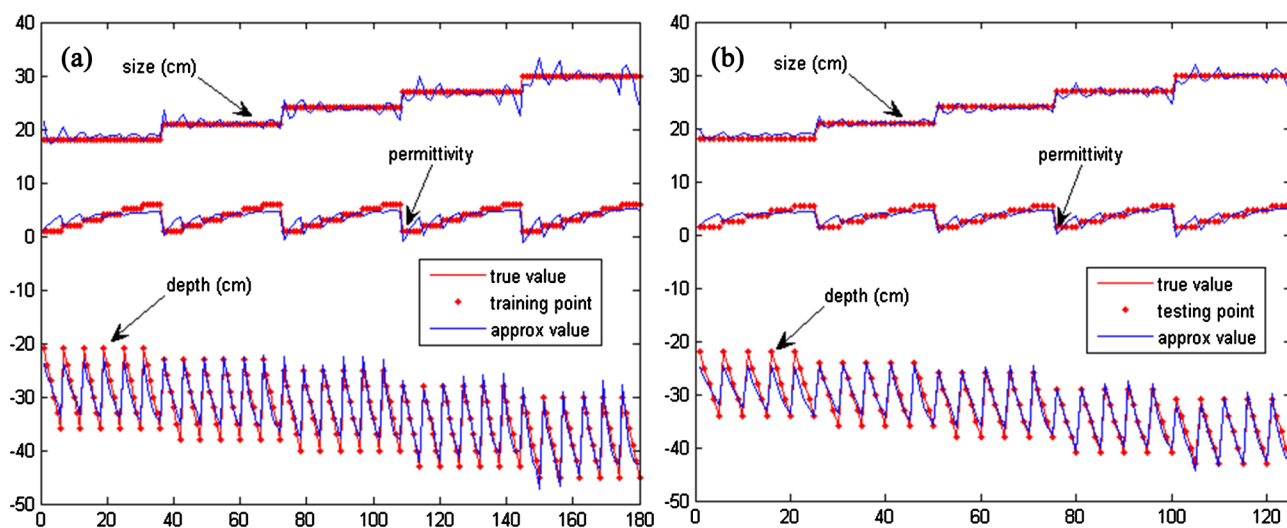


Figure 6. Results for cube target with noisy data: (a) predict on training points; (b) predict on testing points.

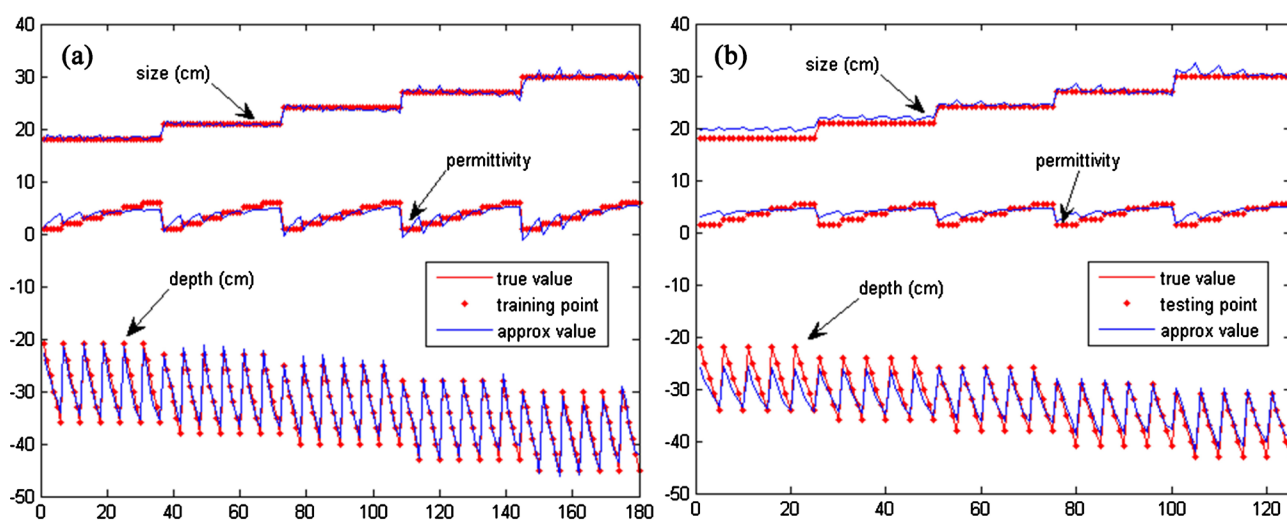


Figure 7. Results for sphere target with noise-free data: (a) predict on training points; (b) predict on testing points.

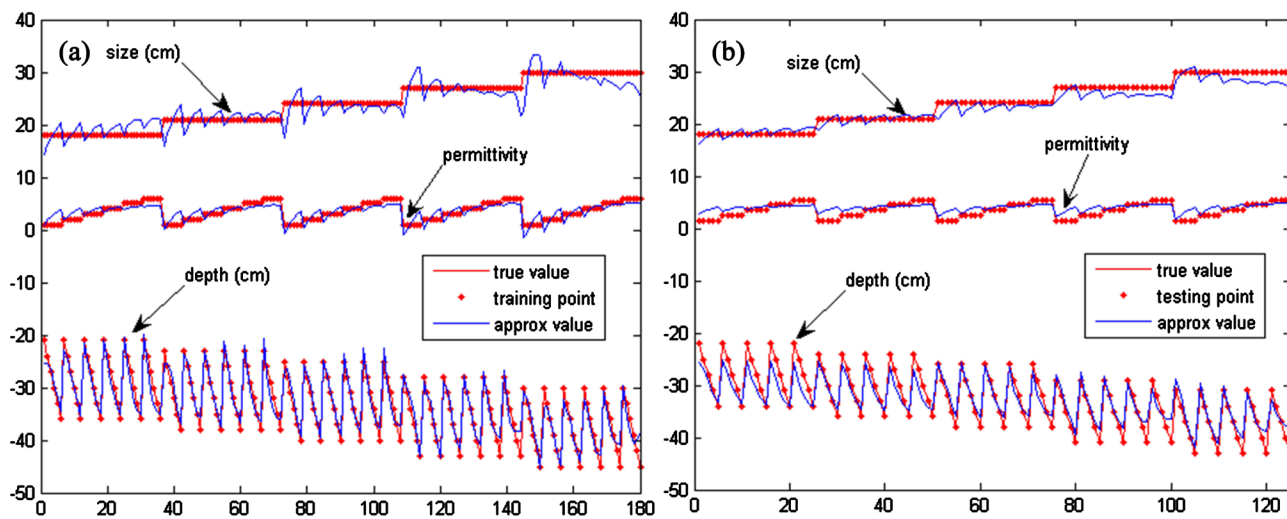


Figure 8. Results for sphere target with noisy data: (a) predict on training points; (b) predict on testing points.

Example 2 We consider sphere target in this example. **Figure 7** and **figure 8** illustrate the results. **Figure 7** is from noise-free data. **Figure 8** is based on noisy data. The noise is generated from zero-mean Gaussian with 10% of average training data $\|\mathbf{x}\|$ as standard deviation. The “size” is the diameter of sphere.

We applied the model to two completely different types of problems, the model works well for both application. The results display this regression model can apply to various types of regression problems.

4. Conclusion

A Bayesian vector-regression algorithm has been developed. The model employs a statistical prior that favors a sparse model, for which most of its weights are zero [5]. This model improves the algorithm in [9], and reduces the number of hyperparameters, which need to be calculated in the algorithm, from two to one. The model is not established for one specific problem, and so can be applied to different regression problems. We have discussed the theoretical development of the model and have presented several example results for two different applications. One is for function approximation, and the other is for inverse scattering of dielectric targets buried in a lossy half space. It has been demonstrated that the algorithm works well for different applications.

Conflicts of Interest

The author declares no conflicts of interest regarding the publication of this paper.

References

- [1] Law, T. and Shawe-Taylor, J. (2017) Practical Bayesian Support Vector Regression for Financial Time Series Prediction and Market Condition Change Detection. *Quantitative Finance*, **17**, 1403-1416. <https://doi.org/10.1080/14697688.2016.1267868>
- [2] Yu, J. (2012) A Bayesian Inference Based Two-Stage Support Vector Regression Framework for Soft Sensor Development in Batch Bioprocesses. *Computers & Chemical Engineering*, **41**, 134-144. <https://doi.org/10.1016/j.compchemeng.2012.03.004>
- [3] Jacobs, J.P. (2012) Bayesian Support Vector Regression with Automatic Relevance Determination Kernel for Modeling of Antenna Input Characteristics. *IEEE Transactions on Antennas and Propagation*, **60**, 2114-2118. <https://doi.org/10.1109/TAP.2012.2186252>
- [4] Hans, C. (2009) Bayesian Lasso Regression. *Biometrika*, **96**, 835-845. <https://doi.org/10.1093/biomet/asp047>
- [5] Tipping, M.E. (2001) Sparse Bayesian Learning and the Relevance Vector Machine. *Journal of Machine Learning Research*, **1**, 211-244.
- [6] Scholkopf, B. and Smola, A.J. (2001) Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. MIT Press, Cambridge.
- [7] Berger, J.O. (1985) Statistical Decision Theory and Bayesian Analysis. 2nd Edition, Springer, Berlin. <https://doi.org/10.1007/978-1-4757-4286-2>

- [8] Mardia, K.V., Kent, J.T. and Bibby, J.B. (1979) *Multivariate Analysis*. Academic Press, New York.
- [9] Yu, Y., Krishnapuram, B. and Carin, L. (2004) Inverse Scattering with Sparse Bayesian Vector Regression. *Inverse Problems*, **20**, 217-231.
<https://doi.org/10.1088/0266-5611/20/6/S13>