

# Chapter 4

## Sparse Bayesian Learning

*Sparse Bayesian Learning* [6] is a general Bayesian framework within supervised Machine Learning. It can be applied to both regression and classification tasks. The *Relevance Vector Machine*, or *RVM*, is a particular specialisation of the Sparse Bayesian Learning model which has identical functional form to the Support Vector Machine (SVM). However, the RVM comes with a number of key advantages over the SVM. The solution produced by a RVM is typically much sparser than the solution by a comparable SVM. Furthermore, the RVM is a probabilistic model and as such, allows us to estimate error bounds in its predictions.

In this chapter, we will derive the Sparse Bayesian Learning model for regression. We will summarise both the original inference algorithm [6] and also the faster “Sequential Sparse Bayesian Learning Algorithm” [7].

### 4.1 Model Specification

We are given a data set of  $N$  input vectors  $\{\mathbf{x}^{(i)}\}_{i=1}^N$  and their associated *targets*  $\{y^{(i)}\}_{i=1}^N$ . The input vectors live in  $D$ -dimensional space,  $\mathbf{x} \in \mathbb{R}^D$ . The targets are real values,  $y \in \mathbb{R}$ .<sup>1</sup>

We model the data using a linearly-weighted sum of  $M$  fixed basis functions  $\{\phi_j(\cdot)\}_{j=1}^M$  and base our predictions on the function  $f(\cdot)$  defined as

$$f(\mathbf{x}; \mathbf{w}) = \sum_{j=1}^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad (4.1)$$

---

<sup>1</sup>When using the Sparse Bayesian model for regression, we assume the targets are real-valued. It is also possible to use the model for classification in which case the targets are assumed to be discrete class labels.

where  $\mathbf{w} = [w_1, \dots, w_M]^T$  and  $\boldsymbol{\phi}(\cdot) = [\phi_1(\cdot), \dots, \phi_M(\cdot)]^T$ . Using a large number  $M$  of non-linear basis functions  $\phi_j : \mathbb{R}^D \rightarrow \mathbb{R}$  allows for a highly flexible model.

The *Relevance Vector Machine*, or RVM, is a specialisation of the Sparse Bayesian Learning model in which the basis functions take the form of *kernel functions*

$$\phi_j(\cdot) \equiv K(\cdot, \mathbf{x}^{(j)}).$$

This defines a basis function for each training data point  $\mathbf{x}^{(i)}$ . Typically, we also include an additional *bias* term  $\phi_0(\cdot) \equiv 1$ , so that  $M = N + 1$ . The RVM has identical functional form to the popular Support Vector Machine (SVM), but superior properties. It typically gives sparser solutions than the SVM and has the additional advantage of providing confidence measures for its predictions.

However, in the following derivation, we will stick to the case of general basis functions  $\phi_j : \mathbb{R}^D \rightarrow \mathbb{R}$ . Thus  $M$  need not equal  $N + 1$  and may, in fact, be a lot larger.

To train the model (4.1), i.e. find values for  $\mathbf{w}$  that are optimal in some sense, we make the standard assumption that our training data are samples from the model with additive noise:

$$\begin{aligned} y^{(i)} &= f(\mathbf{x}^{(i)}; \mathbf{w}) + \epsilon^{(i)} \\ &= \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)}) + \epsilon^{(i)} \quad i = 1, \dots, N. \end{aligned} \quad (4.2)$$

The errors  $\{\epsilon^{(i)}\}_{i=1}^N$  are assumed to be independent samples from a zero-mean Gaussian distribution with variance  $\sigma^2$

$$p(\epsilon^{(i)}) = \mathcal{N}(\epsilon^{(i)} | 0, \sigma^2) \quad i = 1, \dots, N. \quad (4.3)$$

Combining equation (4.2) with equation (4.1), we may express the model for the complete data using matrix notation:

$$\mathbf{y} = \boldsymbol{\Phi} \mathbf{w} + \boldsymbol{\epsilon} \quad (4.4)$$

where  $\boldsymbol{\epsilon} = [\epsilon^{(1)}, \dots, \epsilon^{(N)}]^T$ . The  $N \times M$  matrix  $\boldsymbol{\Phi}$  is known as the design matrix. The  $i$ th row of  $\boldsymbol{\Phi}$  is given by  $\boldsymbol{\phi}(\mathbf{x}^{(i)})^T$ . The  $j$ th column of  $\boldsymbol{\Phi}$  is given by  $\boldsymbol{\phi}_j = [\phi_j(\mathbf{x}^{(1)}), \dots, \phi_j(\mathbf{x}^{(N)})]^T$ , which is also referred to as the  $j$ th *basis vector*. Thus

$$\boldsymbol{\Phi} = [\boldsymbol{\phi}_1 \quad \dots \quad \boldsymbol{\phi}_M] = \begin{bmatrix} \boldsymbol{\phi}(\mathbf{x}^{(1)})^T \\ \vdots \\ \boldsymbol{\phi}(\mathbf{x}^{(N)})^T \end{bmatrix}$$

Combining equation (4.4) and equation (4.3), we find that the complete data likelihood function is given by

$$\begin{aligned} p(\mathbf{y} | \mathbf{w}, \sigma^2) &= \mathcal{N}(\mathbf{y} | \mathbf{w}, \sigma^2 \mathbf{I}_M) \\ &= (2\pi\sigma^2)^{-N/2} \exp\left\{-\frac{1}{2\sigma^2} \|\mathbf{y} - \Phi\mathbf{w}\|^2\right\} \end{aligned} \quad (4.5)$$

where  $\mathbf{I}_M$  is the  $M \times M$  identity matrix.

So far, we have specified the general linear regression model. To get to the sparse Bayesian formulation, we define a zero-mean Gaussian prior distribution over the parameters  $\mathbf{w}$

$$p(\mathbf{w} | \boldsymbol{\alpha}) = \prod_{j=1}^M \mathcal{N}(w_j | \alpha_j^{-1}) \quad (4.6)$$

where  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_M]^T$  is a vector of  $M$  hyperparameters. It is important to note that each hyperparameter  $\alpha_j$  is solely responsible for controlling the strength of the prior of its associated weight  $w_j$ . If  $\alpha_j$  is large, the prior over  $w_j$  is very strongly peaked at zero. This form of the prior distribution is, more than anything, responsible for the dramatic sparsity in the final model.

To complete the specification, we must define a prior over the noise parameter  $\sigma^2$  and the a hyperprior over the hyperparameters  $\boldsymbol{\alpha}$ . Following the derivation in [6], we use the following Gamma <sup>2</sup> priors

$$\begin{aligned} p(\boldsymbol{\alpha} | a, b) &= \prod_{j=1}^M \text{Gamma}(\alpha_j | a, b) \\ p(\beta | c, d) &= \text{Gamma}(\beta | c, d) \end{aligned}$$

where  $\beta \equiv \sigma^{-2}$ .

As a side note, consider the prior of  $\mathbf{w}$  after marginalising out the dependence on the hyperpriors  $\boldsymbol{\alpha}$ . Since each  $w_j$  is normally distributed with an unknown precision parameter  $\alpha_j$  and since the (hyper)prior over  $\alpha_j$  is the Gamma distribution and

---

<sup>2</sup> The Gamma distribution is defined by

$$\text{Gamma}(z | a, b) = \Gamma(a)^{-1} b^a z^{a-1} \exp(-bz) \quad z, a, b > 0$$

where  $\Gamma(\cdot)$  is the Gamma function defined by

$$\Gamma(z) = \int_0^\infty t^{z-1} \exp(-t) dt.$$

therefore conjugate to  $p(w_j | \alpha_j)$ , it follows that the resulting integral can be evaluated analytically

$$\begin{aligned} p(\mathbf{w} | a, b) &= \int p(\mathbf{w} | \boldsymbol{\alpha}) p(\boldsymbol{\alpha} | a, b) d\boldsymbol{\alpha} \\ &= \prod_{j=1}^M \int \mathcal{N}(w_j | 0, \alpha_j^{-1}) \text{Gamma}(\alpha_j | a, b) d\alpha_j \\ &= \prod_{j=1}^M \frac{b^a \Gamma(a + \frac{1}{2})}{(2\pi)^{\frac{1}{2}} \Gamma(a)} \left( b + \frac{w_j^2}{2} \right)^{-(a + \frac{1}{2})}. \end{aligned}$$

This corresponds to a product of independent Student-t density functions over the weights  $w_j$ . The choice  $a = b = 0$  implies that  $p(\mathbf{w} | a, b) \propto \prod_{j=1}^M 1/|w_j|$ . As discussed in [6], it is this hierarchical formulation of the weight prior that is ultimately responsible for encouraging sparse solutions.

## 4.2 Model Inference

We have specified the likelihood model for the data and a prior distribution over the model parameters. The next step in Bayesian inference is to compute the posterior distribution of the parameters. We begin by setting up Bayes' Rule

$$p(\mathbf{w}, \boldsymbol{\alpha}, \sigma^2 | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{w}, \boldsymbol{\alpha}, \sigma^2) p(\mathbf{w}, \boldsymbol{\alpha}, \sigma^2)}{\int p(\mathbf{y} | \mathbf{w}, \boldsymbol{\alpha}, \sigma^2) p(\mathbf{w}, \boldsymbol{\alpha}, \sigma^2) d\mathbf{w} d\boldsymbol{\alpha} d\sigma^2} \quad (4.7)$$

The integral in the denominator of (4.7) is computationally intractable and we must resort to an alternative strategy. First, we decompose the left-hand-side of equation (4.7) as

$$p(\mathbf{w}, \boldsymbol{\alpha}, \sigma^2 | \mathbf{y}) = p(\mathbf{w} | \mathbf{y}, \boldsymbol{\alpha}, \sigma^2) p(\boldsymbol{\alpha}, \sigma^2 | \mathbf{y}). \quad (4.8)$$

Next, we use Bayes' Rule to compute the posterior distribution of the weights given  $\boldsymbol{\alpha}$  and  $\sigma^2$

$$p(\mathbf{w} | \mathbf{y}, \boldsymbol{\alpha}, \sigma^2) = \frac{p(\mathbf{y} | \mathbf{w}, \sigma^2) p(\mathbf{w} | \boldsymbol{\alpha})}{p(\mathbf{y} | \boldsymbol{\alpha}, \sigma^2)} \quad (4.9)$$

The denominator of the right-hand-side is known as the *marginal likelihood* and given by

$$p(\mathbf{y} | \boldsymbol{\alpha}, \sigma^2) = \int p(\mathbf{y} | \mathbf{w}, \sigma^2) p(\mathbf{w} | \boldsymbol{\alpha}) d\mathbf{w} \quad (4.10)$$

Since  $\boldsymbol{\alpha}$  and  $\sigma^2$  are treated as fixed quantities in equation (4.9), the Gaussian density  $p(\mathbf{w} | \boldsymbol{\alpha})$  is the conjugate prior to the Gaussian likelihood function  $p(\mathbf{y} | \mathbf{w}, \sigma^2)$ . Thus,

the integral in equation (4.10) is a convolution of two Gaussians and therefore equal to another Gaussian:

$$\begin{aligned} p(\mathbf{y} \mid \boldsymbol{\alpha}, \sigma^2) &= \mathcal{N}(\mathbf{y} \mid \mathbf{0}, \mathbf{C}) \\ &= (2\pi)^{-N/2} |\mathbf{C}|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} \right\} \end{aligned} \quad (4.11)$$

where

$$\mathbf{C} = \sigma^2 \mathbf{I}_N + \boldsymbol{\Phi} \mathbf{A}^{-1} \boldsymbol{\Phi}^T. \quad (4.12)$$

The posterior distribution for  $\mathbf{w}$  is also a Gaussian,  $p(\mathbf{w} \mid \mathbf{y}, \boldsymbol{\alpha}, \sigma^2) = \mathcal{N}(\mathbf{w} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Its mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$  are given by

$$\boldsymbol{\Sigma} = \left( \sigma^{-2} \boldsymbol{\Phi}^T \boldsymbol{\Phi} + \mathbf{A} \right)^{-1} \quad (4.13)$$

$$\boldsymbol{\mu} = \sigma^{-2} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{y} \quad (4.14)$$

with  $\mathbf{A} = \text{diag}(\boldsymbol{\alpha})$ .

Finally, we need to find the posterior of the hyperparameters,  $p(\boldsymbol{\alpha}, \sigma^2 \mid \mathbf{y})$ . This part is computationally intractable, so instead we approximate the posterior by a delta-function at its mode. Hence, the problem reduces to finding the values of  $\boldsymbol{\alpha}$  and  $\sigma^2$  that maximise  $p(\boldsymbol{\alpha}, \sigma^2 \mid \mathbf{y}) \propto p(\mathbf{y} \mid \boldsymbol{\alpha}, \sigma^2) p(\boldsymbol{\alpha}) p(\sigma^2)$ .

Here, we make the simplifying assumption that  $a = b = c = d = 0$ , giving us uniform (but improper) hyperpriors (see [6] for the general case). Maximising  $p(\boldsymbol{\alpha}, \sigma^2 \mid \mathbf{y})$  is then equivalent to maximising the marginal likelihood, or equivalently, its logarithm

$$\begin{aligned} \mathcal{L}(\boldsymbol{\alpha}, \sigma^2) &= \log p(\mathbf{y} \mid \boldsymbol{\alpha}, \sigma^2) = \log \mathcal{N}(\mathbf{y} \mid \mathbf{0}, \mathbf{C}) \\ &= -\frac{1}{2} \left[ N \log 2\pi + \log |\mathbf{C}| + \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} \right] \end{aligned} \quad (4.15)$$

The procedure of finding  $\boldsymbol{\alpha}$  and  $\sigma^2$  that maximise the (log) marginal likelihood (4.15) is also known as *type-II Maximum likelihood* and *evidence approximation*.

### 4.2.1 Original Training Algorithm

The original training algorithm in [6] is derived by setting the derivatives of (4.15) to zero. We obtain the following update equations for  $\boldsymbol{\alpha}$  and  $\sigma^2$ :

$$\alpha_j^{\text{new}} = \frac{\gamma_j}{\mu_j^2} \quad (4.16)$$

**Algorithm 4.1** Sparse Bayesian Learning: Original Training Algorithm

---

```

1: Choose some initial positive values for  $\sigma^2$  and  $\alpha_j$  for  $j = 1, \dots, M$ 
2: repeat
3:    $\mathbf{A} = \text{diag}(\boldsymbol{\alpha})$ 
4:    $\boldsymbol{\Sigma} = (\sigma^{-2}\boldsymbol{\Phi}^T\boldsymbol{\Phi} + \mathbf{A})^{-1}$ 
5:    $\boldsymbol{\mu} = \sigma^{-2}\boldsymbol{\Sigma}\boldsymbol{\Phi}^T\mathbf{y}$ 

6:   for  $j = 1, \dots, M$  do
7:      $\gamma_j = 1 - \alpha_j \Sigma_{jj}$ 
8:      $\alpha_j = \gamma_j / \mu_j^2$ 
9:   end for
10:   $\sigma^2 = \|\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\mu}\|^2 / (N - \sum_j \gamma_j)$ 
11: until Convergence

```

---

$$(\sigma^2)^{\text{new}} = \frac{\|\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\mu}\|^2}{N - \sum_j \gamma_j} \quad (4.17)$$

where  $\mu_j$  is the  $j$ th component of the posterior mean  $\boldsymbol{\mu}$  (4.14). The quantities  $\gamma_j$  are defined by

$$\gamma_j = 1 - \alpha_j \Sigma_{jj} \quad (4.18)$$

where  $\Sigma_{jj}$  is the  $j$ th diagonal element of the posterior covariance  $\boldsymbol{\Sigma}$  (4.13).

To train the model, we can start by giving  $\boldsymbol{\alpha}$  and  $\sigma^2$  some initial values and evaluate the mean and covariance of the weights posterior using equations (4.14) and (4.13), respectively. Next we alternate between re-estimating the hyperparameters  $\boldsymbol{\alpha}$  and  $\sigma^2$  using (4.16) and (4.17) and updating the posterior mean and covariance parameters using (4.14) and (4.13). We continue until a relevant convergence criterion is met. For example, we may choose to stop if the change in the marginal likelihood - or, alternatively, the change in the parameter values - between two iterations is below a certain pre-defined threshold.

This procedure is summarised in Algorithm (4.1).

During training, it is typically observed that many of the hyperparameters  $\alpha_j$  tend to infinity. Equations (4.14) and (4.13) imply that the weights  $w_j$  corresponding to these hyperparameters have a posterior distribution where the mean and the variance are both zero, meaning their posterior is infinitely peaked at zero. As a consequence, the corresponding basis functions  $\phi_j(\cdot)$  are effectively removed from the model and we achieve sparsity.

In the case of the RVM, where  $\phi_j(\cdot) \equiv K(\cdot, \mathbf{x}^{(j)})$ , the input vectors  $\mathbf{x}^{(j)}$  corresponding to the remaining non-zero weights are known as the *relevance vectors* of the model.

### 4.2.2 Sequential Sparse Bayesian Learning Algorithm

A central drawback of the training algorithm discussed in the previous section is its speed. The computational complexity scales with the cube of the number of basis functions. During training, as basis functions are pruned from the model, the algorithm accelerates. Nevertheless, if  $M$  is very large, the procedure can be very expensive to run.

An alternative strategy of maximising the marginal likelihood (4.15) was developed by [7], resulting in a highly accelerated training algorithm: the *Sequential Sparse Bayesian Learning Algorithm*. It starts with a single basis function and maximises the marginal likelihood by sequentially adding and deleting candidate basis functions.

To derive the algorithm, we follow the analysis in [5] and consider the dependence of the marginal likelihood  $\mathcal{L}(\boldsymbol{\alpha}, \sigma^2)$  on a single hyperparameter  $\alpha_j$ . First, we decompose the matrix  $\mathbf{C}$ , defined in (4.12), as follows:

$$\begin{aligned}\mathbf{C} &= \sigma^2 \mathbf{I}_N + \sum_{m \neq j} \alpha_m^{-1} \boldsymbol{\phi}_m \boldsymbol{\phi}_m^T + \alpha_j^{-1} \boldsymbol{\phi}_j \boldsymbol{\phi}_j^T \\ &= \mathbf{C}_{-j} + \alpha_j^{-1} \boldsymbol{\phi}_j \boldsymbol{\phi}_j^T\end{aligned}$$

where  $\mathbf{C}_{-j} \equiv \sigma^2 \mathbf{I}_N + \sum_{m \neq j} \alpha_m^{-1} \boldsymbol{\phi}_m \boldsymbol{\phi}_m^T$  is  $\mathbf{C}$  without the contribution of the  $j$ th basis vector  $\boldsymbol{\phi}_j$ . Making use of standard identities [WHICH ONES?] for matrix inverses and determinants, we can express  $|\mathbf{C}|$  and  $\mathbf{C}^{-1}$  as

$$\begin{aligned}\mathbf{C}^{-1} &= \mathbf{C}_{-j}^{-1} - \frac{\mathbf{C}_{-j}^{-1} \boldsymbol{\phi}_j \boldsymbol{\phi}_j^T \mathbf{C}_{-j}^{-1}}{\alpha_j + \boldsymbol{\phi}_j^T \mathbf{C}_{-j}^{-1} \boldsymbol{\phi}_j} \\ |\mathbf{C}| &= |\mathbf{C}_{-j}| \left| 1 + \alpha_j^{-1} \boldsymbol{\phi}_j^T \mathbf{C}_{-j}^{-1} \boldsymbol{\phi}_j \right|\end{aligned}$$

This allows us to decompose the marginal likelihood:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\alpha}, \sigma^2) &= \mathcal{L}(\boldsymbol{\alpha}_{-j}, \sigma^2) + \frac{1}{2} \left[ \log \alpha_j - \log(\alpha_j + s_j) + \frac{q_j^2}{\alpha_j + s_j} \right] \\ &\equiv \mathcal{L}(\boldsymbol{\alpha}_{-j}, \sigma^2) + \ell(\alpha_j, \sigma^2)\end{aligned}\tag{4.19}$$

This conveniently separates terms in  $\alpha_j$  in  $\ell(\alpha_j, \sigma^2)$  from the remaining terms in  $\mathcal{L}(\boldsymbol{\alpha}_{-j}, \sigma^2)$ , which is the (log) marginal likelihood with the basis vector  $\boldsymbol{\phi}_j$  excluded.

The quantity  $s_j$  is the *sparsity factor*, defined as

$$s_j = \boldsymbol{\phi}_j^T \mathbf{C}_{-j}^{-1} \boldsymbol{\phi}_j.$$

It serves as a measure of how much the marginal likelihood would decrease if we added  $\boldsymbol{\phi}_j$  to the model. The quantity  $q_j$ , on the other hand, is known as the *quality factor*. It is defined as

$$q_j = \boldsymbol{\phi}_j^T \mathbf{C}_{-j}^{-1} \mathbf{y}$$

and measures the extent to which  $\boldsymbol{\phi}_j$  increases  $\mathcal{L}(\boldsymbol{\alpha}, \sigma^2)$  by helping to explain the data  $\mathbf{y}$ . Thus, a particular basis vector  $\boldsymbol{\phi}_j$  should not be included in the model if its sparsity factor  $s_j$  is large, unless it is offset by a large quality factor  $q_j$ .

We can see this more explicitly if we consider the first derivative of  $\ell(\alpha_j, \sigma^2)$  with respect to  $\alpha_j$  [5]

$$\frac{\partial \ell(\alpha_j, \sigma^2)}{\partial \alpha_j} = \frac{\alpha_j^{-1} s_j^2 - (q_j^2 - s_j)}{2(\alpha_j + s_j)^2}$$

Equating it to zero (and noting that  $\alpha_j$  is an inverse-variance and therefore positive), we obtain the following solution for  $\alpha_j$

$$\alpha_j = \begin{cases} s_j^2 / (q_j^2 - s_j) & \text{if } q_j^2 > s_j \\ +\infty & \text{otherwise} \end{cases}. \quad (4.20)$$

## 4.3 Making Predictions

sparse mu predictive distribution

For details on the RVM and its implementation see [3, 6, 7].



# References

- [1] Candès, E. J. and Wakin, M. B. (2008). An introduction to compressive sampling. *Signal Processing Magazine, IEEE*, 25(2):21–30.
- [2] DeVore, R. A., Jawerth, B., and Lucier, B. J. (1992). Image compression through wavelet transform coding. *Information Theory, IEEE Transactions on*, 38(2):719–746.
- [3] Pilikos, G. (2014). Signal reconstruction using compressive sensing. MPhil thesis, University of Cambridge.
- [4] Stollnitz, E. J., DeRose, T. D., and Salesin, D. H. (1995). Wavelets for computer graphics: a primer. 1. *Computer Graphics and Applications, IEEE*, 15(3):76–84.
- [5] Tipping, A. and Faul, A. (2002). Analysis of sparse bayesian learning. *Advances in neural information processing systems*, 14:383–389.
- [6] Tipping, M. E. (2001). Sparse bayesian learning and the relevance vector machine. *The journal of machine learning research*, 1:211–244.
- [7] Tipping, M. E., Faul, A. C., et al. (2003). Fast marginal likelihood maximisation for sparse bayesian models. In *AISTATS*.