

---

# Analysis of Sparse Bayesian Learning

---

Anita C. Faul    Michael E. Tipping  
Microsoft Research  
St George House, 1 Guildhall St  
Cambridge CB2 3NH, U.K.

## Abstract

The recent introduction of the ‘relevance vector machine’ has effectively demonstrated how sparsity may be obtained in generalised linear models within a Bayesian framework. Using a particular form of Gaussian parameter prior, ‘learning’ is the maximisation, with respect to hyperparameters, of the *marginal likelihood* of the data. This paper studies the properties of that objective function, and demonstrates that conditioned on an individual hyperparameter, the marginal likelihood has a unique maximum which is computable in closed form. It is further shown that if a derived ‘sparsity criterion’ is satisfied, this maximum is exactly equivalent to ‘pruning’ the corresponding parameter from the model.

## 1 Introduction

We consider the approximation, from a training sample, of real-valued functions, a task variously referred to as prediction, regression, interpolation or function approximation. Given a set of data  $\{\mathbf{x}_n, t_n\}_{n=1}^N$  the ‘target’ samples  $t_n = f(\mathbf{x}_n) + \epsilon_n$  are conventionally considered to be realisations of a deterministic function  $f$ , potentially corrupted by some additive noise process. This function  $f$  will be modelled by a linearly-weighted sum of  $M$  fixed basis functions  $\{\phi_m(\mathbf{x})\}_{m=1}^M$ :

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

and the objective is to infer values of the parameters/weights  $\{w_m\}_{m=1}^M$  such that  $\hat{f}$  is a ‘good’ approximation of  $f$ .

While accuracy in function approximation is generally universally valued, there has been significant recent interest [2, 9, 3, 5]) in the notion of *sparsity*, a consequence of learning algorithms which set significant numbers of the parameters  $w_m$  to zero.

A methodology which effectively combines both these measures of merit is that of ‘sparse Bayesian learning’, briefly reviewed in Section 2, and which was the basis for the recent introduction of the *relevance vector machine* (RVM) and related models [6, 1, 7]. This model exhibits some very compelling properties, in particular a dramatic degree of sparseness even in the case of highly overcomplete basis sets

( $M \gg N$ ). The sparse Bayesian learning algorithm essentially involves the maximisation of a marginalised likelihood function with respect to hyperparameters in the model prior. In the RVM, this was achieved through re-estimation equations, the behaviour of which was not fully understood. In this paper we present further relevant theoretical analysis of the properties of the marginal likelihood which gives a much fuller picture of the nature of the model and its associated learning procedure. This is detailed in Section 3, and we close with a summary of our findings and discussion of their implications in Section 4 (and which, to avoid repetition here, the reader may wish to preview at this point).

## 2 Sparse Bayesian Learning

We now very briefly review the methodology of sparse Bayesian learning, more comprehensively described elsewhere [6]. To simplify and generalise the exposition, we omit to notate any functional dependence on the inputs  $\mathbf{x}$  and combine quantities defined over the training set and basis set within  $N$ - and  $M$ -vectors respectively. Using this representation, we first write the generative model as:

$$\mathbf{t} = \mathbf{f} + \boldsymbol{\epsilon}, \quad (2)$$

where  $\mathbf{t} = (t_1, \dots, t_N)^\top$ ,  $\mathbf{f} = (f_1, \dots, f_N)^\top$  and  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_N)^\top$ . The approximator is then written as:

$$\hat{\mathbf{f}} = \Phi \mathbf{w}, \quad (3)$$

where  $\Phi = [\phi_1 \dots \phi_M]$  is a general  $N \times M$  design matrix with column vectors  $\phi_m$  and  $\mathbf{w} = (w_1, \dots, w_M)^\top$ . Recall that in the context of (1),  $\Phi_{nm} = \phi_m(\mathbf{x}_n)$  and  $\hat{\mathbf{f}} = \{\hat{f}(\mathbf{x}_1), \dots, \hat{f}(\mathbf{x}_N)\}^\top$ .

The sparse Bayesian framework assumes an independent zero-mean Gaussian noise model, with variance  $\sigma^2$ , giving a multivariate Gaussian likelihood of the target vector  $\mathbf{t}$ :

$$p(\mathbf{t}|\mathbf{w}, \sigma^2) = (2\pi)^{-N/2} \sigma^{-N} \exp \left\{ -\frac{\|\mathbf{t} - \hat{\mathbf{f}}\|^2}{2\sigma^2} \right\}. \quad (4)$$

The prior over the parameters is mean-zero Gaussian:

$$p(\mathbf{w}|\boldsymbol{\alpha}) = (2\pi)^{-M/2} \prod_{m=1}^M \alpha_m^{1/2} \exp \left( -\frac{\alpha_m w_m^2}{2} \right), \quad (5)$$

where the key to the model sparsity is the use of  $M$  independent hyperparameters  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_M)^\top$ , one per weight (or basis vector), which moderate the strength of the prior. Given  $\boldsymbol{\alpha}$ , the *posterior* parameter distribution is Gaussian and given via Bayes' rule as  $p(\mathbf{w}|\mathbf{t}, \boldsymbol{\alpha}) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with

$$\boldsymbol{\Sigma} = (\mathbf{A} + \sigma^{-2} \Phi^\top \Phi)^{-1} \quad \boldsymbol{\mu} = \sigma^{-2} \boldsymbol{\Sigma} \Phi^\top \mathbf{t}, \quad (6)$$

and  $\mathbf{A}$  defined as  $\text{diag}(\alpha_1, \dots, \alpha_M)$ . Sparse Bayesian learning can then be formulated as a *type-II maximum likelihood* procedure, in that objective is to maximise the *marginal likelihood*, or equivalently, its logarithm  $\mathcal{L}(\boldsymbol{\alpha})$  with respect to the hyperparameters  $\boldsymbol{\alpha}$ :

$$\mathcal{L}(\boldsymbol{\alpha}) = \log p(\mathbf{t}|\boldsymbol{\alpha}, \sigma^2) = \log \int_{-\infty}^{\infty} p(\mathbf{t}|\mathbf{w}, \sigma^2) p(\mathbf{w}|\boldsymbol{\alpha}) d\mathbf{w}, \quad (7)$$

$$= -\frac{1}{2} [N \log 2\pi + \log |\mathbf{C}| + \mathbf{t}^\top \mathbf{C}^{-1} \mathbf{t}], \quad (8)$$

with  $\mathbf{C} = \sigma^2 \mathbf{I} + \Phi \mathbf{A}^{-1} \Phi^\top$ .

Once most-probable values  $\alpha_{\text{MP}}$  have been found<sup>1</sup>, in practice they can be plugged into (6) to give a posterior mean (most probable) point estimate for the parameters  $\mu_{\text{MP}}$  and from that a mean final approximator:  $\hat{\mathbf{f}}_{\text{MP}} = \Phi \mu_{\text{MP}}$ .

Empirically, the local maximisation of the marginal likelihood (8) with respect to  $\alpha$  has been seen to work highly effectively [6, 1, 7]. Accurate predictors may be realised, which are typically highly sparse as a result of the maximising values of many hyperparameters being infinite. From (6) this leads to a parameter posterior infinitely peaked at zero for many weights  $w_m$  with the consequence that  $\mu_{\text{MP}}$  correspondingly comprises very few non-zero elements.

However, the learning procedure in [6] relied upon heuristic re-estimation equations for the hyperparameters, the behaviour of which was not well characterised. Also, little was known regarding the properties of (8), the validity of the local maximisation thereof and importantly, and perhaps most interestingly, the conditions under which  $\alpha$ -values would become infinite. We now give, through a judicious re-writing of (8), a more detailed analysis of the sparse Bayesian learning procedure.

### 3 Properties of the Marginal Likelihood $\mathcal{L}(\alpha)$

#### 3.1 A convenient re-writing

We re-write  $\mathbf{C}$  from (8) in a convenient form to analyse the dependence on a single hyperparameter  $\alpha_i$ :

$$\begin{aligned} \mathbf{C} &= \sigma^2 \mathbf{I} + \sum_m \alpha_m \phi_m \phi_m^T, = \sigma^2 \mathbf{I} + \sum_{m \neq i} \alpha_m^{-1} \phi_m \phi_m^T + \alpha_i^{-1} \phi_i \phi_i^T, \\ &= \mathbf{C}_{-i} + \alpha_i^{-1} \phi_i \phi_i^T, \end{aligned} \quad (9)$$

where we have defined  $\mathbf{C}_{-i} = \sigma^2 \mathbf{I} + \sum_{m \neq i} \alpha_m^{-1} \phi_m \phi_m^T$  as the covariance matrix with the influence of basis vector  $\phi_i$  removed, equivalent also to  $\alpha_i = \infty$ .

Using established matrix determinant and inverse identities, (9) allows us to write the terms of interest in  $\mathcal{L}(\alpha)$  as:

$$|\mathbf{C}| = |\mathbf{C}_{-i}| |1 + \alpha_i^{-1} \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i|, \quad (10)$$

$$\mathbf{C}^{-1} = \mathbf{C}_{-i}^{-1} - \frac{\mathbf{C}_{-i}^{-1} \phi_i \phi_i^T \mathbf{C}_{-i}^{-1}}{\alpha_i + \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i}, \quad (11)$$

which gives

$$\begin{aligned} \mathcal{L}(\alpha) &= -\frac{1}{2} \left[ N \log(2\pi) + \log |\mathbf{C}_{-i}| + \mathbf{t}^T \mathbf{C}_{-i}^{-1} \mathbf{t} \right. \\ &\quad \left. - \log \alpha_i + \log(\alpha_i + \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i) - \frac{(\phi_i^T \mathbf{C}_{-i}^{-1} \mathbf{t})^2}{\alpha_i + \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i} \right], \\ &= \mathcal{L}(\alpha_{-i}) + \frac{1}{2} \left[ \log \alpha_i - \log(\alpha_i + \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i) + \frac{(\phi_i^T \mathbf{C}_{-i}^{-1} \mathbf{t})^2}{\alpha_i + \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i} \right], \\ &= \mathcal{L}(\alpha_{-i}) + \ell(\alpha_i), \end{aligned} \quad (12)$$

where  $\mathcal{L}(\alpha_{-i})$  is the log marginal likelihood with  $\alpha_i$  (and thus  $w_i$  and  $\phi_i$ ) removed from the model and we have now isolated the terms in  $\alpha_i$  in the function  $\ell(\alpha_i)$ .

<sup>1</sup>The most-probable noise variance  $\sigma_{\text{MP}}^2$  can also be directly and successfully estimated from the data [6], but for clarity in this paper, we assume without prejudice to our results that its value is fixed.

### 3.2 First derivatives of $\mathcal{L}(\boldsymbol{\alpha})$

**Previous results.** In [6], based on earlier results from [4], the gradient of the marginal likelihood was computed as:

$$\frac{\partial \mathcal{L}(\boldsymbol{\alpha})}{\partial \alpha_i} = \frac{1}{2} \left[ \frac{1}{\alpha_i} - \mu_i^2 - \Sigma_{ii} \right], \quad (13)$$

with  $\mu_i$  the  $i$ -th element of  $\boldsymbol{\mu}$  and  $\Sigma_{ii}$  the  $i$ -th diagonal element of  $\boldsymbol{\Sigma}$ . This then leads to re-estimation updates for  $\alpha_i$  in terms of  $\mu_i$  and  $\Sigma_{ii}$  where, disadvantageously, these latter terms are themselves functions of  $\alpha_i$ .

**A new, simplified, expression.** In fact, by instead differentiating (12) directly, (13) can be seen to be equivalent to:

$$\frac{\partial \mathcal{L}(\boldsymbol{\alpha})}{\partial \alpha_i} = \frac{\partial \ell(\alpha_i)}{\partial \alpha_i} = \frac{1}{2} \left[ \frac{1}{\alpha_i} - \frac{1}{\alpha_i + \boldsymbol{\phi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\phi}_i} - \frac{(\boldsymbol{\phi}_i^T \mathbf{C}_{-i}^{-1} \mathbf{t})^2}{(\alpha_i + \boldsymbol{\phi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\phi}_i)^2} \right], \quad (14)$$

where advantageously,  $\alpha_i$  now occurs only explicitly since  $\mathbf{C}_{-i}$  is independent of  $\alpha_i$ . For convenience, we combine terms and re-write (14) as:

$$\frac{\partial \mathcal{L}(\boldsymbol{\alpha})}{\partial \alpha_i} = \frac{\alpha_i^{-1} S_i^2 - (Q_i^2 - S_i)}{2(\alpha_i + S_i)^2}, \quad (15)$$

where, for simplification of this and forthcoming expressions, we have defined:

$$Q_i \triangleq \boldsymbol{\phi}_i^T \mathbf{C}_{-i}^{-1} \mathbf{t}, \quad S_i \triangleq \boldsymbol{\phi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\phi}_i. \quad (16)$$

The term  $Q_i$  can be interpreted as a ‘quality’ factor: a measure of how well  $\boldsymbol{\phi}_i$  increases  $\mathcal{L}(\boldsymbol{\alpha})$  by helping to explain the data, while  $S_i$  is a ‘sparsity’ factor which measures how much the inclusion of  $\boldsymbol{\phi}_i$  serves to decrease  $\mathcal{L}(\boldsymbol{\alpha})$  through ‘inflating’  $\mathbf{C}$  (*i.e.* adding to the normalising factor).

### 3.3 Stationary points of $\mathcal{L}(\boldsymbol{\alpha})$

Equating (15) to zero indicates that stationary points of the marginal likelihood occur both at  $\alpha_i = +\infty$  (note that, being an inverse variance,  $\alpha_i$  must be positive) and for:

$$\alpha_i = \frac{S_i^2}{Q_i^2 - S_i}, \quad (17)$$

subject to  $Q_i^2 > S_i$  as a consequence again of  $\alpha_i > 0$ .

Since the right-hand-side of (17) is independent of  $\alpha_i$ , we may find the stationary points of  $\ell(\alpha_i)$  analytically without iterative re-estimation. To find the nature of those stationary points, we consider the second derivatives.

### 3.4 Second derivatives of $\mathcal{L}(\boldsymbol{\alpha})$

#### 3.4.1 With respect to $\alpha_i$

Differentiating (15) a second time with respect to  $\alpha_i$  gives:

$$\frac{\partial^2 \mathcal{L}(\boldsymbol{\alpha})}{\partial \alpha_i^2} = \frac{-\alpha_i^{-2} S_i^2 (\alpha_i + S_i)^2 - 2(\alpha_i + S_i) [\alpha_i^{-1} S_i^2 - (Q_i^2 - S_i)]}{2(\alpha_i + S_i)^4}, \quad (18)$$

and we now consider (18) for both finite- and infinite- $\alpha_i$  stationary points.

**Finite  $\alpha$ .** In this case, for stationary points given by (17), we note that the second term in the numerator in (18) is zero, giving:

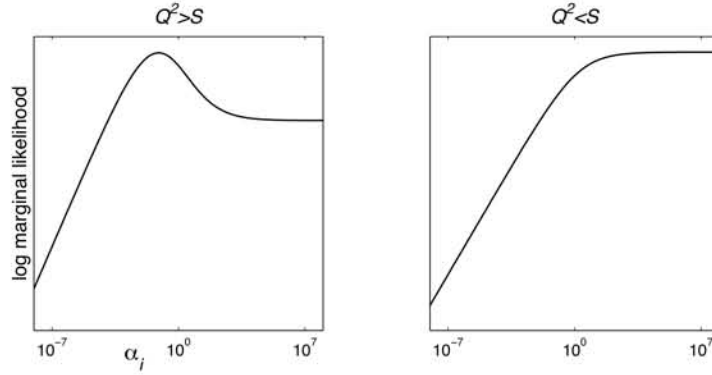
$$\left. \frac{\partial^2 \mathcal{L}(\boldsymbol{\alpha})}{\partial \alpha_i^2} \right|_{\alpha_i = \frac{S_i^2}{Q_i^2 - S_i}} = \frac{-S_i^2}{2\alpha_i^2(\alpha_i + S_i)^2}. \quad (19)$$

We see that (19) is always negative, and therefore  $\ell(\alpha_i)$  has a *maximum*, which must be unique, for  $Q_i^2 - S_i > 0$  and  $\alpha_i$  given by (17).

**Infinite  $\alpha$ .** For this case, (18) and indeed, all further derivatives, are uninformatively zero at  $\alpha_i = \infty$ , but from (15) we can see that as  $\alpha_i \rightarrow \infty$ , the sign of the gradient is given by the sign of  $-(Q_i^2 - S_i)$ .

If  $Q_i^2 - S_i > 0$ , then the gradient at  $\alpha_i = \infty$  is negative so as  $\alpha_i$  decreases  $\ell(\alpha_i)$  must increase to its unique maximum given by (17). It follows that  $\alpha_i = \infty$  is thus a minimum. Conversely, if  $Q_i^2 - S_i < 0$ ,  $\alpha_i = \infty$  is the unique maximum of  $\ell(\alpha_i)$ . If  $Q_i^2 - S_i = 0$ , then this maximum and that given by (17) coincide.

We now have a full characterisation of the marginal likelihood as a function of a single hyperparameter, which is illustrated in Figure 1.



**Figure 1:** Example plots of  $\ell(\alpha_i)$  against  $\alpha_i$  (on a log scale) for  $Q_i^2 > S_i$  (left), showing the single maximum at finite  $\alpha_i$ , and  $Q_i^2 < S_i$  (right), showing the maximum at  $\alpha_i = \infty$ .

### 3.4.2 With respect to $\alpha_j$ , $j \neq i$

To obtain the off-diagonal terms of the second derivative (Hessian) matrix, it is convenient to manipulate (15) to express it in terms of  $\mathbf{C}$ . From (11) we see that

$$\boldsymbol{\phi}_i^T \mathbf{C}^{-1} \mathbf{t} = \frac{\alpha_i Q_i}{\alpha_i + S_i}, \quad \text{and} \quad \boldsymbol{\phi}_i^T \mathbf{C}^{-1} \boldsymbol{\phi}_i = \frac{\alpha_i S_i}{\alpha_i + S_i}. \quad (20)$$

Utilising these identities in (15) gives:

$$\frac{\partial \mathcal{L}(\boldsymbol{\alpha})}{\partial \alpha_i} = \frac{1}{2\alpha_i^2} [\boldsymbol{\phi}_i^T \mathbf{C}^{-1} \boldsymbol{\phi}_i - (\boldsymbol{\phi}_i^T \mathbf{C}^{-1} \mathbf{t})^2]. \quad (21)$$

We now write:

$$\frac{\partial^2 \mathcal{L}(\boldsymbol{\alpha})}{\partial \alpha_i \partial \alpha_j} = \nabla_{ij}^2 + \delta_{ij} \nabla_{ii}^2, \quad (22)$$

where  $\delta_{ij}$  is the Kronecker ‘delta’ function, allowing us to separate out the additional (diagonal) term that appears only when  $i = j$ .

Writing, similarly to (9) earlier,  $\mathbf{C} = \mathbf{C}_{-j} + \alpha_j^{-1} \phi_j \phi_j^T$ , substituting into (21) and differentiating with respect to  $\alpha_j$  gives:

$$\begin{aligned} \nabla_{ij}^2 &= \frac{1}{2\alpha_i^2} \left[ \frac{(\phi_i^T \mathbf{C}_{-j}^{-1} \phi_j)^2}{(\alpha_j + \phi_j^T \mathbf{C}_{-j}^{-1} \phi_j)^2} - 2(\phi_i^T \mathbf{C}^{-1} \mathbf{t}) \frac{(\phi_i^T \mathbf{C}_{-j}^{-1} \phi_j)(\phi_j^T \mathbf{C}_{-j}^{-1} \mathbf{t})}{(\alpha_j + \phi_j^T \mathbf{C}_{-j}^{-1} \phi_j)^2} \right], \\ &= \frac{\phi_i^T \mathbf{C}_{-j}^{-1} \phi_j}{2\alpha_i^2 (\alpha_j + \phi_j^T \mathbf{C}_{-j}^{-1} \phi_j)^2} \left[ \phi_i^T \mathbf{C}_{-j}^{-1} \phi_j - 2(\phi_i^T \mathbf{C}^{-1} \mathbf{t})(\phi_j^T \mathbf{C}_{-j}^{-1} \mathbf{t}) \right], \\ &= \frac{\phi_i^T \mathbf{C}^{-1} \phi_j}{2\alpha_i^2 \alpha_j^2} \left[ \phi_i^T \mathbf{C}^{-1} \phi_j - 2(\phi_i^T \mathbf{C}^{-1} \mathbf{t})(\phi_j^T \mathbf{C}^{-1} \mathbf{t}) \right], \end{aligned} \quad (23)$$

while we have

$$\nabla_{ii}^2 = -\frac{1}{\alpha_i^3} \left[ \phi_i^T \mathbf{C}^{-1} \phi_i - (\phi_i^T \mathbf{C}^{-1} \mathbf{t})^2 \right]. \quad (24)$$

If all hyperparameters  $\alpha_i$  are individually set to their maximising values, *i.e.*  $\boldsymbol{\alpha} = \boldsymbol{\alpha}_{\text{MP}}$  such that all  $\partial \mathcal{L}(\boldsymbol{\alpha}) / \partial \alpha_i = 0$ , then even if all  $\partial^2 \mathcal{L}(\boldsymbol{\alpha}) / \partial \alpha_i^2$  are negative, there may still be a non-axial direction in which the likelihood could be increasing. We now rule out this possibility by showing that the Hessian is negative semi-definite.

First, we note from (21) that if  $\partial \mathcal{L}(\boldsymbol{\alpha}) / \partial \alpha_i = 0$ ,  $\nabla_{ii}^2 = 0$ . Then, if  $\mathbf{v}$  is a generic nonzero direction vector:

$$\begin{aligned} \mathbf{v}^T \left[ \frac{\partial^2 \mathcal{L}(\boldsymbol{\alpha})}{\partial \alpha_i \partial \alpha_j} \right] \mathbf{v} &= \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \frac{v_i v_j}{\alpha_i^2 \alpha_j^2} \left[ (\mathbf{C}^{-1/2} \phi_i)^T (\mathbf{C}^{-1/2} \phi_j) \right]^2 \\ &\quad - \left( \sum_{i=1}^M \frac{|v_i|}{\alpha_i^2} |\phi_i^T \mathbf{C}^{-1} \mathbf{t}| \|\mathbf{C}^{-1/2} \phi_i\| \right)^2, \\ &\leq \frac{1}{2} \left( \sum_{i=1}^M \frac{|v_i|}{\alpha_i^2} \|\mathbf{C}^{-1/2} \phi_i\|^2 \right)^2 \\ &\quad - \left( \sum_{i=1}^M \frac{|v_i|}{\alpha_i^2} |\phi_i^T \mathbf{C}^{-1} \mathbf{t}| \|\mathbf{C}^{-1/2} \phi_i\| \right)^2, \end{aligned} \quad (25)$$

where we use the Cauchy-Schwarz inequality. If the gradient vanishes, then for all  $i = 1, \dots, M$  either  $\alpha_i = \infty$ , or from (21),  $\phi_i^T \mathbf{C}^{-1} \phi_i = (\phi_i^T \mathbf{C}^{-1} \mathbf{t})^2$ . It follows directly from (25) that the Hessian is negative semi-definite, with (25) only zero where  $\mathbf{v}$  is orthogonal to all finite  $\alpha$  values.

## 4 Summary

Sparse Bayesian learning proposes the iterative maximisation of the marginal likelihood function  $\mathcal{L}(\boldsymbol{\alpha})$  with respect to the hyperparameters  $\boldsymbol{\alpha}$ . Our analysis has shown the following:

- I. As a function of an individual hyperparameter  $\alpha_i$ ,  $\mathcal{L}(\boldsymbol{\alpha})$  has a unique maximum computable in closed-form. (This maximum is, of course, dependent on the values of all other hyperparameters.)

- II. If the criterion  $Q_i^2 - S_i$  (defined in Section 3.2) is negative, this maximum occurs at  $\alpha_i = \infty$ , equivalent to the removal of basis function  $i$  from the model.
- III. The point where all individual marginal likelihood functions  $\ell(\alpha_i)$  are maximised is a joint maximum (not necessarily unique) over all  $\alpha_i$ .

These results imply the following consequences.

- From I, we see that if we update, in any arbitrary order, the  $\alpha_i$  parameters using (17), we are guaranteed to increase the marginal likelihood at each step, unless already at a maximum. Furthermore, we would expect these updates to be more efficient than those given in [6], which individually only increase, not maximise,  $\ell(\alpha_i)$ .
- Result III indicates that sequential optimisation of individual  $\alpha_i$  cannot lead to a stationary point from which a joint maximisation over all  $\alpha$  may have escaped. (*i.e.* the stationary point is not a saddle point.)
- The result II confirms the qualitative argument and empirical observation that many  $\alpha_i \rightarrow \infty$  as a result of the optimisation procedure in [6]. The inevitable implication of finite numerical precision prevented the genuine sparsity of the model being verified in those earlier simulations.
- We conclude by noting that the maximising hyperparameter solution (17) remains valid if  $\alpha_i$  is already infinite. This means that *basis functions not even in the model* can be assessed and their corresponding hyperparameters updated if desired. So as well as the facility to increase  $\mathcal{L}(\alpha)$  through the ‘pruning’ of basis functions if  $Q_i^2 - S_i \leq 0$ , new basis functions can be introduced if  $\alpha_i = \infty$  but  $Q_i^2 - S_i > 0$ . This has highly desirable computational consequences which we are exploiting to obtain a powerful ‘constructive’ approximation algorithm [8].

## References

- [1] C. M. Bishop and M. E. Tipping. Variational relevance vector machines. In C. Bouillier and M. Goldszmidt, editors, *Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence*, pages 46–53. Morgan Kaufmann, 2000.
- [2] S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. Technical Report 479, Department of Statistics, Stanford University, 1995.
- [3] Y. Grandvalet. Least absolute shrinkage is equivalent to quadratic penalisation. In L. Niklasson, M. Bodén, and T. Ziemke, editors, *Proceedings of the Eighth International Conference on Artificial Neural Networks (ICANN98)*, pages 201–206. Springer, 1998.
- [4] D. J. C. MacKay. Bayesian interpolation. *Neural Computation*, 4(3):415–447, 1992.
- [5] A. J. Smola, B. Schölkopf, and G. Rätsch. Linear programs for automatic accuracy control in regression. In *Proceedings of the Ninth International Conference on Artificial Neural Networks*, pages 575–580, 1999.
- [6] M. E. Tipping. The Relevance Vector Machine. In S. A. Solla, T. K. Leen, and K.-R. Müller, editors, *Advances in Neural Information Processing Systems 12*, pages 652–658. MIT Press, 2000.
- [7] M. E. Tipping. Sparse kernel principal component analysis. In *Advances in Neural Information Processing Systems 13*. MIT Press, 2001.
- [8] M. E. Tipping and A. C. Faul. Bayesian pursuit. Submitted to NIPS\*01.
- [9] V. N. Vapnik. *Statistical Learning Theory*. Wiley, 1998.