

# ON THE RELATIONSHIP BETWEEN BAYESIAN ERROR BARS AND THE INPUT DATA DENSITY

C K I Williams, C Qazaz, C M Bishop and H Zhu

Neural Computing Research Group, Aston University, UK.

## ABSTRACT

We investigate the dependence of Bayesian error bars on the distribution of data in input space. For generalized linear regression models we derive an upper bound on the error bars which shows that, in the neighbourhood of the data points, the error bars are substantially reduced from their prior values. For regions of high data density we also show that the contribution to the output variance due to the uncertainty in the weights can exhibit an approximate inverse proportionality to the probability density. Empirical results support these conclusions.

## 1 INTRODUCTION

When given a prediction, it is also very useful to be given some idea of the ‘‘error bars’’ associated with that prediction. Error bars arise naturally in a Bayesian treatment of neural networks and are made up of two terms, one due to the posterior weight uncertainty, and the other due to the intrinsic noise in the data<sup>1</sup>. As the two contributions are independent, we have

$$\sigma_y^2(\mathbf{x}) = \sigma_w^2(\mathbf{x}) + \sigma_v^2(\mathbf{x}) \quad (1)$$

where  $\sigma_w^2(\mathbf{x})$  is the variance of the output due to weight uncertainty and  $\sigma_v^2(\mathbf{x})$  is the variance of the intrinsic noise.

Under the assumption that the posterior in weight space can be approximated by a Gaussian (MacKay (1)), we have

$$\sigma_w^2(\mathbf{x}) = \mathbf{g}^T(\mathbf{x})\mathbf{A}^{-1}\mathbf{g}(\mathbf{x}) \quad (2)$$

where  $\mathbf{A}$  is the Hessian matrix of the model and  $\mathbf{g}(\mathbf{x}) = \partial y(\mathbf{x}; \mathbf{w})/\partial \mathbf{w}$  is the vector of the derivatives of the output with respect to the weight parameters in the network.  $\mathbf{A}$  contains contributions from both the prior distribution on the weights and the effect of the training data.

Although the weight uncertainty component of the error bar is given by equation 2, the dependence of this quantity on the location of the training points is not at all obvious. Intuitively we would expect

<sup>1</sup>If the network used is not the correct generative model for the data there will be a third component due to model mis-specification; we do not discuss this further in this paper.

the error bars from the prior (i.e. before any data is seen) to be quite large, and that the effect of the training data would be to reduce the magnitude of the error bars for those regions of the input space close to the data points, while leaving large error bars further away. The purpose of this paper is to provide theoretical insights to support this intuition. In particular, our analysis focusses on generalized linear regression (such as radial basis function networks with fixed basis function parameters) and allows us to quantify the extent of the reduction and the length scale over which it occurs.

We also show that the relationship  $\sigma_w^2(\mathbf{x}) \simeq \sigma_v^2[Np(\mathbf{x})V(\mathbf{x})]^{-1}$  holds approximately, where  $p(\mathbf{x})$  is the density of the data in the input space,  $N$  is the number of data points in the training set and  $V(\mathbf{x})$  is a function of  $\mathbf{x}$  that measures a volume in the input space. This relationship pertains to the ‘‘high-data’’ limit where the effect of the data overwhelms the prior in the Hessian.

## 2 GENERALIZED LINEAR REGRESSION

Consider a generalized linear regression (GLR) model of the form

$$y(\mathbf{x}) = \boldsymbol{\phi}^T(\mathbf{x})\mathbf{w} = \sum_{j=1}^m w_j \phi_j(\mathbf{x}) \quad (3)$$

where  $j = 1, \dots, m$  labels the basis functions  $\{\phi_j\}$  of the model. Given a data set  $D = ((\mathbf{x}_1, t_1), (\mathbf{x}_2, t_2), \dots, (\mathbf{x}_N, t_N))$ , a squared error function with noise variance<sup>2</sup>  $\sigma_v^2 = \beta^{-1}$  and a regularizer of the form  $\alpha \mathbf{w}^T \mathbf{S} \mathbf{w}/2$ , the posterior mean value of the weights  $\hat{\mathbf{w}}$  is the choice of  $\mathbf{w}$  that minimizes the quadratic form

$$E = \frac{1}{2\sigma_v^2} \sum_i^N \left( t_i - \sum_j w_j \phi_j(\mathbf{x}_i) \right)^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{S} \mathbf{w} \quad (4)$$

so that  $\hat{\mathbf{w}}$  is the solution of

$$(\beta \mathbf{B} + \alpha \mathbf{S})\hat{\mathbf{w}} = \beta \boldsymbol{\Phi}^T \mathbf{t} \quad (5)$$

<sup>2</sup>In this section we assume that  $\sigma_v^2$  is independent of  $\mathbf{x}$ . This assumption can be easily relaxed, but at the expense of somewhat more complicated notation.

where  $\Phi$  is the  $n \times m$  design matrix

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \cdots & \phi_m(\mathbf{x}_1) \\ \phi_1(\mathbf{x}_2) & \phi_2(\mathbf{x}_2) & \cdots & \phi_m(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_n) & \phi_2(\mathbf{x}_n) & \cdots & \phi_m(\mathbf{x}_n) \end{pmatrix} \quad (6)$$

$\mathbf{B} = \Phi^T \Phi$  and  $\mathbf{t}$  is the vector of targets. Writing  $\mathbf{A} = \beta \mathbf{B} + \alpha \mathbf{S}$ , we find

$$\hat{y}(\mathbf{x}) = \phi^T(\mathbf{x}) \hat{\mathbf{w}} = \beta \phi^T(\mathbf{x}) \mathbf{A}^{-1} \Phi^T \mathbf{t} \stackrel{def}{=} \mathbf{k}^T(\mathbf{x}) \mathbf{t} \quad (7)$$

where  $\hat{y}(\mathbf{x})$  is the function obtained from equation 3 using  $\hat{\mathbf{w}}$  as the weight vector. Equation 7 defines the effective kernel  $\mathbf{k}(\mathbf{x})$  and makes it clear that  $\hat{y}(\mathbf{x})$  can be written as a linear combination of the target values, i.e. it is a *linear smoother* (see, e.g. Hastie and Tibshirani (2)).

The contribution of the uncertainty of the weights to the variance of the prediction is given from equation 2 by

$$\sigma_w^2(\mathbf{x}) = \phi^T(\mathbf{x}) \mathbf{A}^{-1} \phi(\mathbf{x}) \quad (8)$$

Note that for generalized linear regression this expression is exact, and that the error bars (given  $\sigma_w^2$ ) are independent of the targets.

### 3 ERROR BARS FOR GLR

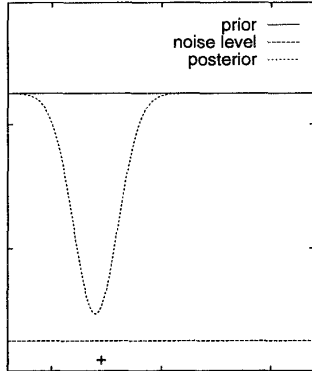


Figure 1: A schematic illustration of the effect of one data point on  $\sigma_w^2(\mathbf{x})$ . The posterior variance is reduced from its prior level in the neighbourhood of the data point (+), but remains above the noise level.

In this section we analyze the response of the prior variance to the addition of the data points. In particular we show that the effect of a single data point is to pull the  $\sigma_w^2(\mathbf{x})$  surface down to a value less than  $2\sigma_w^2(\mathbf{x})^3$  at and nearby to the data point, and that

<sup>3</sup>The analysis in this section permits the noise level to vary as a function of  $\mathbf{x}$ .

the length scale over which this effect operates is determined by the prior covariance function

$$C(\mathbf{x}, \mathbf{x}') = \phi^T(\mathbf{x}) \mathbf{A}_0^{-1} \phi(\mathbf{x}') \quad (9)$$

where  $\mathbf{A}_0 = \alpha \mathbf{S}$ .

The main tool used in this analysis is the effect of adding just one data point. A schematic illustration of this effect is shown in Figure 1. The variance due to the prior is quite large (and roughly constant over  $\mathbf{x}$ -space). Adding a single data point pulls down the variance in its neighbourhood (but not as far as the  $\sigma_w^2$  limit).

Figure 1 is relevant because we can show (see Appendix A.1) that  $\sigma_w^2(\mathbf{x})$ , when all data points are used to compute the Hessian, is never greater than  $\sigma_w^2(\mathbf{x})$  when any subset of the data points are used, and hence the surface pertaining to any particular data point is an upper bound on the overall surface.

To obtain a bound on the depth of the dip, consider the case when there is only one data point (at  $\mathbf{x} = \mathbf{x}_i$ ), so that the Hessian is given by  $\mathbf{A} = \mathbf{A}_0 + \beta(\mathbf{x}_i) \phi(\mathbf{x}_i) \phi^T(\mathbf{x}_i)$ . Using the identity

$$(\mathbf{M} + \mathbf{v}\mathbf{v}^T)^{-1} = \mathbf{M}^{-1} - \frac{(\mathbf{M}^{-1}\mathbf{v})(\mathbf{v}^T\mathbf{M}^{-1})}{1 + \mathbf{v}^T\mathbf{M}^{-1}\mathbf{v}} \quad (10)$$

it is easy to show that

$$\sigma_w^2|_{\mathbf{x}_i}(\mathbf{x}_i) = \sigma_w^2(\mathbf{x}_i) \frac{r_i}{1 + r_i} \quad (11)$$

where  $\sigma_w^2|_{\mathbf{x}_i}$  denotes the posterior weight uncertainty surface due to a data point at  $\mathbf{x}_i$  and

$$r_i = \frac{\phi^T(\mathbf{x}_i) \mathbf{A}_0^{-1} \phi(\mathbf{x}_i)}{\sigma_w^2(\mathbf{x}_i)} \quad (12)$$

i.e.  $r_i$  is the ratio of the prior to noise variances at the point  $\mathbf{x}_i$ . For any positive value of  $z$ , the function  $z/(1+z)$  lies between 0 and 1, hence we see that the  $\sigma_w^2$  contribution to the error bars must always be less than  $\sigma_w^2(\mathbf{x})$  at a data point. Typically the noise variance is much smaller than the prior variance, so  $r_i \gg 1$ .

Further evidence that  $\sigma_w^2$  at any data point is of the order of  $\sigma_w^2(\mathbf{x})$  is provided by the calculation in appendix A.2 which shows that the average of  $\sigma_w^2(\mathbf{x}_i)$  at the data points is less than  $\sigma_w^2 m/N$ , where  $m$  is the number of weights in the model and  $N$  is the number of data points.

For a single data point at  $\mathbf{x}_i$ , we can use equation 10 to show that

$$\sigma_w^2|_{\mathbf{x}_i}(\mathbf{x}) = C(\mathbf{x}, \mathbf{x}) - \frac{(C(\mathbf{x}, \mathbf{x}_i))^2}{1 + C(\mathbf{x}_i, \mathbf{x}_i)} \quad (13)$$

Hence the width of the depression in the variance surface is related to the characteristic length scale of the

prior covariance function  $C(\mathbf{x}, \mathbf{x}')$ . It is also possible to show that if a test point  $\mathbf{x}$  has zero covariance  $C(\mathbf{x}, \mathbf{x}_i)$  with all of the training points  $\{\mathbf{x}_i\}$ , then its posterior variance will be equal its prior variance. We are currently exploring the properties of  $C(\mathbf{x}, \mathbf{x}')$  for different weight priors and choices of basis functions. However, we note that a simple diagonal prior  $\mathbf{S} = \mathbf{I}$  as used by some authors is not in general a very sensible prior, because if the type of basis functions used (e.g. Gaussians, tanh functions etc.) is changed, then the covariance structure of the prior also changes. More sensibly, the weight prior should be chosen so as to approximate some desired prior covariance function  $C(\mathbf{x}, \mathbf{x}')$ .

#### 4 DENSITY DEPENDENCE OF $\sigma_w^2(\mathbf{x})$

As we have already noted, error bars on network predictions would be expected to be relatively large in regions of input space for which there is little data, and smaller in regions of high data density. In this section, we establish an approximate proportionality between the variance due to weight uncertainty and the inverse of the probability density of training data, valid in regions of high data density. A relationship of this kind was conjectured in Bishop (3).

We first consider a special case of the class of generalized linear models where the basis functions are non-overlapping bin (or ‘‘top-hat’’) activation functions<sup>4</sup>. Let the  $i^{\text{th}}$  basis function have height  $h_i$  and a  $d$ -dimensional ‘‘base area’’ of  $V_i$ , where  $d$  is the dimensionality of  $\mathbf{x}$ . If we choose a diagonal prior ( $\mathbf{S} = \mathbf{I}$ ) then the Hessian is diagonal and thus easy to invert.

$$\mathbf{B}_{ij} = \sum_{q=1}^N \phi_i(\mathbf{x}_q) \phi_j(\mathbf{x}_q) = \begin{cases} n_i h_i^2 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases} \quad (14)$$

where  $n_i$  is the number of data points falling in bin  $i$ . From equation (8) the error bars associated with a point  $\mathbf{x}$  which falls into the  $i^{\text{th}}$  bin are given by

$$\sigma_w^2(\mathbf{x}) = \frac{1}{\alpha/h_i^2 + \beta n_i} \quad (15)$$

As usual, the effect of the prior is to reduce the size of error bar compared to the case where it is not present. In the limit of  $\alpha \rightarrow 0$  we have

$$\sigma_w^2(\mathbf{x}) = \frac{\sigma_v^2}{n_i} = \frac{\sigma_v^2}{NV_i \hat{p}(\mathbf{x})} \quad (16)$$

where  $N$  is the total number of data points and  $\hat{p}(\mathbf{x})$  is the histogram estimate of the density inside the bin containing  $\mathbf{x}$ . Equation (16) demonstrates that for

<sup>4</sup>This analysis can easily be extended to arbitrary non-constant basis functions as long as they do not overlap.

this kind of model the error bars are inversely proportional to the input density and to the volume factor  $V_i$ . It also shows that we can understand the reduction in the variance  $\sigma_y^2(\mathbf{x})$  in regions of high density as the  $1/n_i$  effect for the variance of the mean of  $n_i$  (iid) Gaussian variables each of which has variance  $\sigma_v^2$ .

The aim of the remainder of this section is to show how results similar to those for the bin basis functions can be obtained, in certain circumstances, for generalized linear regression models, i.e. that the error bars will be inversely proportional to  $p(\mathbf{x})$  and an area factor  $V(\mathbf{x})$ . The key idea needed is that of an *effective kernel*, which we now describe.

As noted in equation 7, we can write  $\hat{y}(\mathbf{x}) = \mathbf{k}^T(\mathbf{x})\mathbf{t}$ , where  $\mathbf{k}(\mathbf{x})$  is the effective kernel. To take this analysis further it is helpful to think of  $\hat{y}(\mathbf{x}) = \mathbf{k}^T(\mathbf{x})\mathbf{t} = \sum_i k_i t_i$  as an approximation to the integral  $\hat{y}(\mathbf{x}) = \int K(\mathbf{z}; \mathbf{x}) t(\mathbf{z}) d\mathbf{z}$ , where  $K(\mathbf{z}; \mathbf{x})$  (regarded as a function of  $\mathbf{z}$ ) is the effective kernel for the point  $\mathbf{x}$  and  $t(\mathbf{z})$  is a ‘‘target function’’.

Following similar reasoning we obtain

$$\mathbf{B} = \sum_{q=1}^N \phi(\mathbf{x}_q) \phi^T(\mathbf{x}_q) \simeq N \int p(\mathbf{x}) \phi(\mathbf{x}) \phi^T(\mathbf{x}) d\mathbf{x} \quad (17)$$

If the original basis functions  $\phi$  are linearly combined to produce a new set  $\tilde{\phi} = \mathbf{C}\phi$ , then the matrix  $\mathbf{C}$  can be chosen so that  $\int p(\mathbf{x}) \tilde{\phi}_i(\mathbf{x}) \tilde{\phi}_j(\mathbf{x}) d\mathbf{x} = \delta_{ij}$ , where  $\delta_{ij}$  is the Kronecker delta. From now on it is assumed that we are working with the orthonormal basis functions (i.e. the tildes are omitted) and that  $\mathbf{B} = N\mathbf{I}$ . Ignoring the weight prior we obtain

$$\hat{\mathbf{w}} = \beta \mathbf{A}^{-1} \Phi^T \mathbf{t} = \frac{1}{N} \Phi^T \mathbf{t} \quad (18)$$

However,

$$(\Phi^T \mathbf{t})_i = \sum_{q=1}^N \phi_i(\mathbf{x}^q) t(\mathbf{x}^q) \simeq N \int \phi_i(\mathbf{z}) p(\mathbf{z}) t(\mathbf{z}) d\mathbf{z} \quad (19)$$

and so

$$\hat{y}(\mathbf{x}) = \frac{1}{N} \phi^T(\mathbf{x}) \Phi^T \mathbf{t} \quad (20)$$

$$= \int \left\{ \sum_i \phi_i(\mathbf{x}) \phi_i(\mathbf{z}) p(\mathbf{z}) \right\} t(\mathbf{z}) d\mathbf{z} \quad (21)$$

$$= \int K(\mathbf{z}; \mathbf{x}) t(\mathbf{z}) d\mathbf{z} \quad (22)$$

We can also show that  $K(\mathbf{z}; \mathbf{x})$  is the projection of the delta function onto the basis space  $\{\psi_i\}$ , where  $\psi_i(\mathbf{x}) = \phi_i(\mathbf{x}) p(\mathbf{x})$ , and that if a constant (bias) function is one of the original basis functions (before orthonormalization), then  $\int K(\mathbf{z}; \mathbf{x}) d\mathbf{z} = 1$ . The fact that  $K(\mathbf{z}; \mathbf{x})$  is an approximation to the delta

function suggests that as the number of basis functions increases the effective kernel should become more tightly peaked and concentrated around  $\mathbf{x}$ .

We now turn to the variance of the generalized linear model. Using orthonormal basis functions, the error bar at  $\mathbf{x}$  is given by  $\sigma_y^2(\mathbf{x}) = \frac{\sigma_v^2}{N} \phi^T(\mathbf{x}) \phi(\mathbf{x})^5$ . However, this can be rewritten in terms of the effective kernel

$$\sigma_w^2(\mathbf{x}) = \frac{\sigma_v^2}{N} \int \frac{K^2(\mathbf{z}; \mathbf{x})}{p(\mathbf{z})} d\mathbf{z} \quad (23)$$

using the orthonormality properties. If  $K(\mathbf{z}; \mathbf{x})$  is sharply peaked around  $\mathbf{x}$  (i.e. it looks something like a Gaussian) then the  $p(\mathbf{z})$  in the denominator can be pulled through the integral sign as  $p(\mathbf{x})$ . Also,  $\int K^2(\mathbf{z}; \mathbf{x}) d\mathbf{z}$  measures the inverse base area of  $K(\mathbf{z}; \mathbf{x})$ ; for example, for a one dimensional Gaussian with standard deviation  $\sigma$  centered at  $\mathbf{x}$  we find that  $\int K^2(\mathbf{z}; \mathbf{x}) d\mathbf{z} = 1/(2\sigma\sqrt{\pi})$ . Defining

$$\int K^2(\mathbf{z}; \mathbf{x}) d\mathbf{z} \stackrel{def}{=} \frac{1}{V(\mathbf{x})} \quad (24)$$

we can write

$$\sigma_w^2(\mathbf{x}) \simeq \frac{\sigma_v^2}{N p(\mathbf{x}) V(\mathbf{x})} \quad (25)$$

By extending the analysis of appendix A.2 to the continuous case we obtain

$$\int \sigma_w^2(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \frac{\sigma_v^2}{N} \gamma \quad (26)$$

where  $\gamma$  is the effective number of parameters in the model (1), showing that we would expect  $\sigma_w^2(\mathbf{x})$  to be larger for a model with more parameters.

Under the assumption that  $K(\mathbf{z}; \mathbf{x})$  is sharply peaked about  $\mathbf{x}$  we have obtained a result in equation 25 similar to equation 16 for the bin basis functions. We will now present evidence to show that this relationship holds experimentally.

The first experiment has a one dimensional input space. The probability density from which the data was drawn is shown in Figure 2(A). Figure 2(B) shows that for a range of GLR models (and for a two-layer perceptron) there is a close relationship between  $1/\sigma_w^2(\mathbf{x})$  and the density, indicating that  $V(\mathbf{x})$  is roughly constant in the high density regions for these models. This conclusion is backed up by Figure 4, which plots  $1/V(\mathbf{x}) = \int K^2(\mathbf{z}; \mathbf{x}) d\mathbf{z}$  against  $\mathbf{x}$ . The log-log plot in Figure 3 also indicates that the relationship  $\sigma_w^2(\mathbf{x}) \propto p^{-1}(\mathbf{x})$  holds quite reliably, especially for areas with high data density.

<sup>5</sup>It is interesting to note that the error bar  $\sigma_w^2(\mathbf{x})$  can also be obtained from the finite-dimensional effective kernel defined by  $\hat{y}(x) = k^T(x)t$ . Using the assumption that each  $t_i$  has independent, zero-mean noise of variance  $\sigma_t^2$ , we find that the variance of the linear combination  $\hat{y}(x)$  is  $\sigma_w^2(x) = \sigma_t^2 k^T k$ , which can easily be shown to be equivalent to  $\sigma_w^2(x) = \sigma_v^2 \phi^T \phi / N$  for  $\alpha = 0$ .

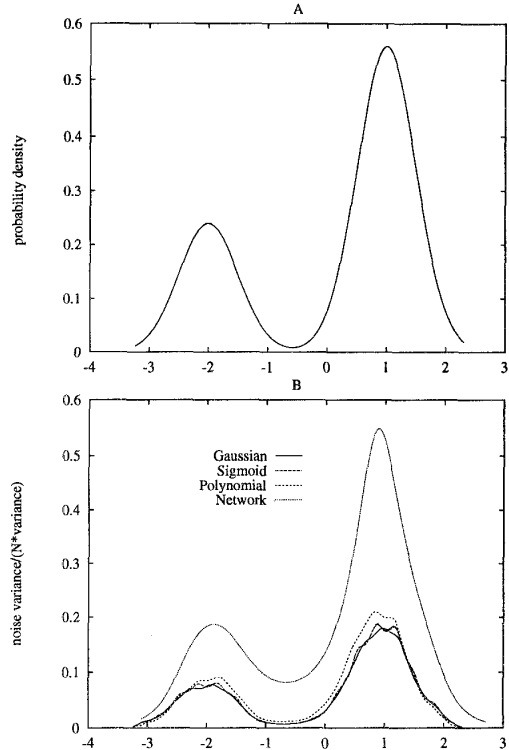


Figure 2: (A) A mixture of two Gaussian densities, from which data points were drawn for the experiments. (B) shows the (scaled) inverse variance against  $x$  for three generalized linear regression (GLR) models and a neural network. The GLR models used Gaussian, sigmoid and polynomial basis functions respectively, and each model consisted of 16 basis functions and a bias and was trained on 1000 data points. (B) also shows the inverse variance for a two layer perceptron with two hidden units. The net was trained on a data set consisting of 200 data points with inputs drawn from the density shown in panel (A) and targets generated from  $\sin(x)$  with added zero-mean Gaussian noise of standard deviation 0.1. For all four models the similarity between the inverse variance for these models and the plot of the density is striking.

Figure 2(B) also shows that the dependence of the overall magnitude of  $\sigma_w^2$  on the number of effective parameters described in equation 26 holds; the two-layer perceptron, which has only seven weights compared to the 16 in the GLR models, has a correspondingly larger inverse variance.

Some effective kernels for the GLR model with a bias and 16 Gaussian basis functions of standard deviation 0.5, spaced equally between  $-5.0$  and  $4.0$  are shown in Figures 5 and 6<sup>6</sup>. The kernels in Figure 5 correspond to areas of high density and show a strong, narrow single peak. For regions of low density Figure 6 shows that the kernels are much wider and more oscillatory, indicating that target values from a wide range of  $x$  values are used to compute  $\hat{y}(x)$ . As the widths of the kernels in the low density regions

<sup>6</sup>Similar  $\{\phi_i\}$  and kernels are obtained for sigmoidal and polynomial basis functions.

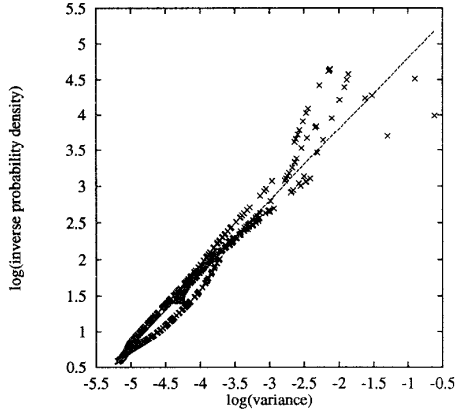


Figure 3: Plot of the log inverse density of the input data against the log of  $\sigma_w^2(x)$  for a generalized linear model with 16 Gaussian basis functions. Note that the points lie close to the line with slope 1, indicating that  $\sigma_w^2(x) \propto p^{-1}(x)$ .

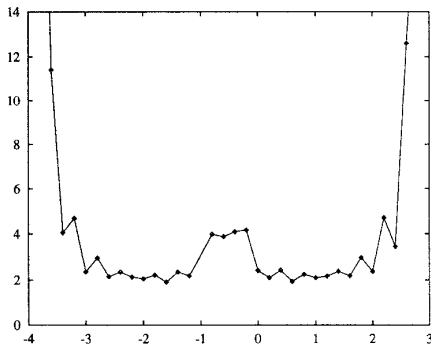


Figure 4: Plot of  $1/V(x) = \int K^2(z; x) dz$  against  $x$  for a GLR model with 16 Gaussian basis functions spaced equally between  $-5.0$  and  $4.0$ , and a bias. Note that the plot is roughly constant in regions of high density.

are greater than the length scale of the variation of the density, we would expect the approximation used in equation 25 to break down at this point.

We have conducted several other experiments with one and two dimensional input spaces which produce similar results to those shown in the log-log plot, Figure 3, including a two-layer perceptron which learned to approximate a function of two inputs.

While this relationship between  $\sigma_w^2(\mathbf{x})$  and the input data density is interesting, it should be noted that its validity is limited at best to regions of high data density. Furthermore, in such regions the contribution to the error bars from  $\sigma_w^2(\mathbf{x})$  is dwarfed by that from the noise term  $\sigma_\nu^2$ . This can be seen in the case of non-overlapping basis functions from equation 16. More generally we can consider the extension of the result 11 to the case of  $n_i$  data points all located at  $\mathbf{x}_i$ . This leads to

$$\sigma_{w|n_i \mathbf{x}_i}^2(\mathbf{x}_i) = \sigma_\nu^2 \frac{r_i}{1 + n_i r_i} \simeq \sigma_\nu^2 / n_i \quad (27)$$

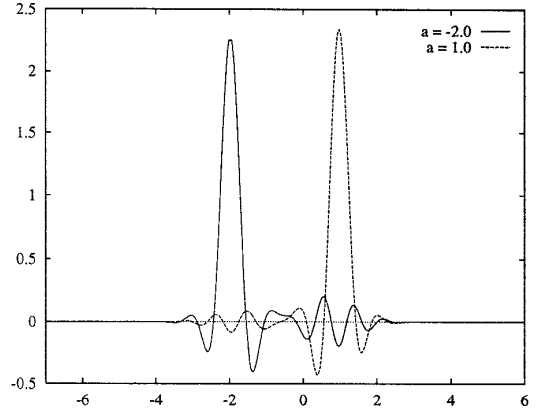


Figure 5: Effective kernels at  $x = 1.0$  and  $-2.0$ , corresponding to high density regions, as shown in figure 2(A). See text for further discussion.

again indicating that for regions of high data density the noise term will dominate.

## 5 DISCUSSION

In this paper we have analyzed the behaviour of the Bayesian error bars for generalized linear regression models. For the case of a single isolated data point we have shown that the error bar is pulled down close to the noise level, and that the length scale over which this effect occurs is characterized by the prior covariance function. We have also shown theoretically that, in regions of high data density, the contribution to the output variance due to the uncertainty in the weights can exhibit an approximate inverse proportionality to the data density. These findings have been supported by numerical simulation. Also, we have noted that, in such high-density regions, this contribution to the variance will be insignificant compared to the contribution arising from the noise term.

Although much of the theoretical analysis has been performed for generalized linear regression models, there is empirical evidence that similar results hold also for multi-layer networks. Furthermore, if the outputs of the network have linear activation functions, then under least-squares training it is effectively a generalized linear regression model with adaptive basis functions. It is therefore a linear smoother with  $\hat{y}(\mathbf{x}) = \mathbf{k}^T(\mathbf{x})\mathbf{t}$ , and hence the result that  $\sigma_w^2 = \mathbf{k}^T \mathbf{k} \sigma_\nu^2$  will still hold. Other results, including the expression 36 derived in Appendix A.2, also hold for general non-linear networks, provided we make the usual Gaussian approximation for the posterior weight distribution, and the outer-product approximation to the Hessian.

One potentially important limitation of the models

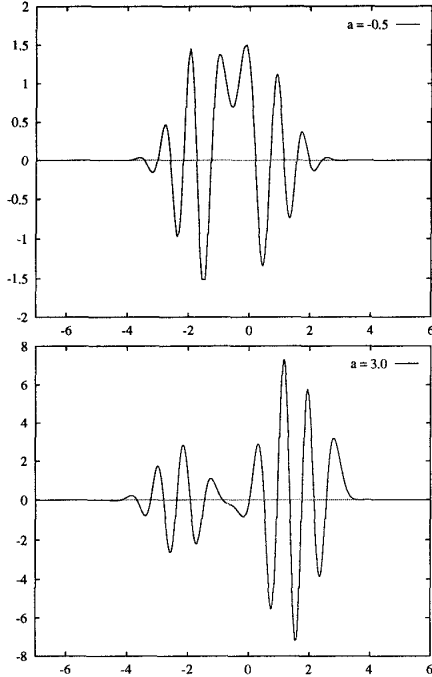


Figure 6: Effective kernels at  $x = -0.5$  and  $3.0$  corresponding to low density regions of the input space, as shown in figure 2(A). Note that the density function seems to define an “envelope” for the lower kernel; even though  $x$  may be in a low density region, the magnitude of  $K(z; x)$  is largest in the high density regions. See text for further discussion.

considered in this paper (and indeed of the models considered by most authors) is that the noise variance  $\sigma_v^2$  is assumed to be a constant, independent of  $\mathbf{x}$ . To understand why this assumption may be particularly restrictive, consider the situation in which there is a lot of data in one region of input space and a single data point in another region. The estimate of the noise variance, which we shall assume to be relatively small, will be dominated by the high density region. However, as we have seen, the error bar will be pulled down to less than  $2\sigma_v^2$  in the neighbourhood of the isolated data point. The model is therefore highly confident of the regression function (i.e. the most probable interpolant) in this region even though there is only a single data point present! If, however, we relax the assumption of a constant  $\sigma_v^2$  then we see that there in the neighbourhood of the isolated data point there is little evidence to suggest a small value of  $\sigma_v^2$  and so we would expect much larger error bars. We are currently investigating models in which  $\sigma_v^2(\mathbf{x})$  is adapted to the data.

#### Acknowledgements

This work was mainly supported by EPSRC grant GR/J75425 (CW) and by an EPSRC post-graduate scholarship to CQ.

## 6 APPENDICES

### A.1

In this appendix we show that for generalized linear regression,  $\sigma_{y|D}^2(x) \leq \sigma_{y|T}^2(x)$ , where  $D$  is the full data set  $((x_1, t_1), \dots, (x_N, t_N))$  and  $T$  is a subset of this data set.

We first note that as  $\sigma_v^2(x)$  is equal in both cases, we are only concerned about the relative contributions from the weight uncertainty to the overall variance. The key to the proof is to decompose the Hessian  $A$  into two parts,  $A_1$  and  $A_2$ , where

$$A_1 = A_0 + \sum_{q \in T} \beta_q \phi_q \phi_q^T \quad A_2 = \sum_{q \notin T} \beta_q \phi_q \phi_q^T \quad (28)$$

and  $A_0 = \alpha S$ . Note that  $A_1$  and  $A_2$  are symmetric non-negative definite, and hence  $A_1^{-1}$  and  $A_2^{-1}$  are also (using the Moore-Penrose pseudo-inverse if necessary). The matrix identity

$$(A_1 + A_2)^{-1} = A_1^{-1} - A_1^{-1}(A_1^{-1} + A_2^{-1})^{-1}A_1^{-1} \quad (29)$$

implies that for any vector  $v$

$$v^T (A_1 + A_2)^{-1} v = v^T A_1^{-1} v - (v^T A_1^{-1} v) (A_1^{-1} + A_2^{-1})^{-1} (A_1^{-1} v) \quad (30)$$

From non-negative definite condition we see that the second term in equation 31 is always non-negative, and hence

$$v^T A_1^{-1} v \geq v^T (A_1 + A_2)^{-1} v \quad (31)$$

Substituting  $\phi(x)$  for  $v$  completes the proof.

### A.2

In this appendix we show that  $\langle \sigma_w^2 \rangle$ , the average value of  $\sigma_w^2(x)$  evaluated at the data points, is equal to  $\sigma_v^2 \gamma / N$ , where  $\gamma$  ( $\leq m$ ) is the *effective number of parameters* in the model (1).

$$\langle \sigma_w^2 \rangle = \frac{1}{N} \sum_i \sigma_w^2(x_i) \quad (32)$$

$$= \frac{1}{N} \sum_i \phi_i^T A^{-1} \phi_i \quad (33)$$

$$= \frac{1}{N} \text{tr}[(\sum_i \phi_i \phi_i^T) A^{-1}] \quad (34)$$

$$= \frac{\sigma_v^2}{N} \text{tr}(\beta B A^{-1}) \quad (35)$$

$$= \frac{\sigma_v^2}{N} \gamma \quad (36)$$

## REFERENCES

1. MacKay D. J. C., 1992, “Bayesian Interpolation”, *Neural Computation*, **4**(3), 415–447.
2. Hastie T. J. and Tibshirani R. J., 1990, “Generalized Additive Models”, Chapman and Hall.
3. Bishop C. M., 1994, “Novelty detection and neural network validation”, *IEE Proceedings: Vision, Image and Signal Processing*, **141**, 217–222.