



---

# THEORETICAL REPORT

## IFS-TR-029

---

*THE USE OF BAYESIAN INFERENCE IN THE TRAINING  
AND INTERPRETATION OF NEURAL NETWORK  
PREDICTORS*

---

### ABSTRACT

The most common method used to train neural networks is the minimisation of the misfit error for a training set of examples using gradient descent based techniques. In the statistical sense, we can view this as a process by which the parameters of the neural network (the weights) are estimated from a sample of data (the training set) using the method of Maximum Likelihood.

However, as was discussed in Theoretical Reports IFS-TR-027 and IFS-TR-028, purely minimising the data misfit is not usually sufficient to ensure that the network is able to generalise well on unseen data. For this reason, regularisation constraints such as weight decay are usually added to the data misfit error during training to control the complexity of the mapping performed by the neural network.

In this report, we demonstrate how such regularisation techniques can be shown to have a Bayesian interpretation. In particular, we show how the method of weight decay (IFS-TR-028) is consistent with the assumption of a Gaussian prior on the weights of the network and a Bayesian approach to training. We also show how a Bayesian interpretation of the weight estimation process leads to the possibility of assigning error bars to the predictions made by a neural network.

We conclude this report with a discussion of some of the practicalities involved in the application of Bayesian approaches to neural network training and review the material contained in other Theoretical Reports that deal with such matters in more detail.

## THE BAYESIAN APPROACH TO NEURAL NETWORKS

We shall consider the use of Bayesian techniques and their application to neural network predictions in two contexts. Firstly we shall examine the Bayesian interpretation of the training process of the neural networks. We shall show how minimisation of the data misfit error is in fact a **special case** of the Bayesian Framework.

Secondly, we shall examine the implication of the Bayesian model of neural network training with respect to the interpretation of the outputs of neural networks. In particular, we shall show how the Hessian of the error derivatives with respect to the weights may be used to assign uncertainties (error bars) on the predictions. We also show, how the Bayesian interpretation lends support to the idea of using committees of neural networks.

### Framework

Consider the task of training a fixed neural network architecture,  $\mathbf{N}$  on a set of training examples

$$D = \{\mathbf{x}_i, t_i, i = 1 \dots N\} \quad (1)$$

where  $\mathbf{x}_i$  is the  $i^{\text{th}}$  training input vector,  $t_i$  is the target output and  $N$  is the number of training patterns. We have assumed, for simplicity, that the target outputs are scalar, but they could equally well be multi-dimensional. In what follows, extension to multi dimensional targets is trivial.

In the standard method of training a neural network we try to find a set of weights  $\mathbf{w}$ , such that the network function  $f(\mathbf{x}, \mathbf{w})$  minimises a least square error term  $E_D$  defined

$$E_D = \frac{1}{2} \sum_{i=1}^N (f(\mathbf{x}_i, \mathbf{w}_i) - t_i)^2 \quad (2)$$

If we assume that the target values  $t_i$  are in fact generated by a process with additive Gaussian noise then the set of weights that minimises the above will correspond to the Maximum Likelihood estimate of the weights,  $\mathbf{w}_{ML}$ .

In the Bayesian approach to network training, rather than finding a single *best* set of weights we attempt to find the probability distribution of the weights given the training data  $D$ , i.e. we wish to find

$$p(\mathbf{w} | D) = \frac{p(D | \mathbf{w})p(\mathbf{w})}{p(D)} \quad (3)$$

The denominator of the above is a weight independent normalisation factor and can be written in terms of the two densities contained in the numerator,

$$p(D) = \int p(D | \mathcal{W}) p(\mathcal{W}) d\mathcal{W} \quad (4)$$

So, in order to determine the density  $p(\mathcal{W} | D)$ , we need to obtain expressions for the densities  $p(\mathcal{W})$  and  $p(D | \mathcal{W})$ .

### Evaluating $p(D | \mathcal{W})$

Let us assume that the target functions are generated by an unknown process with additive, zero mean Gaussian noise with a variance of  $\sigma_n^2$ . We can then define the probability of observing a particular target  $t_i$  given an input vector  $\mathcal{X}_i$  and particular set of network weights,  $\mathcal{W}$  as

$$p(t_i | \mathcal{X}_i, \mathcal{W}) = \frac{1}{\sqrt{2\pi} \sigma_n} \exp\left(-\frac{(f(\mathcal{X}_i, \mathcal{W}) - t_i)^2}{2\sigma_n^2}\right) \quad (5)$$

We can therefore write the likelihood of the entire data set (assuming that the training examples are independent) as

$$\begin{aligned} p(D | \mathcal{W}) &= \prod_{i=1}^N p(t_i | \mathcal{X}_i, \mathcal{W}) \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} \exp\left(-\sum_{i=1}^N \frac{(f(\mathcal{X}_i, \mathcal{W}) - t_i)^2}{2\sigma_n^2}\right) \end{aligned} \quad (6)$$

Substituting Equation (2) into the above and using the standard notational conveniences of writing

$$\beta = \frac{1}{\sigma_n^2}, \quad (7)$$

and

$$Z_D(\beta) = \left(\frac{2\pi}{\beta}\right)^{\frac{N}{2}}, \quad (8)$$

we can re-write the Equation (6) as

$$p(D | \vec{\omega}) = \frac{1}{Z_D} \exp(-\beta E_D) \quad (9)$$

Clearly, for the above the maximum for the likelihood is obtained for a minimum value of  $E_D$ . Thus, standard neural network training techniques that minimise data misfit error,  $E_D$  are equivalent to maximum likelihood estimation.

### Evaluating $p(\vec{\omega})$

This density represents a prior on the values of the weights in the network. It has been empirically shown by a number of researchers that the final weights of an unregularised, trained network tend to exhibit Gaussian distributions. Let us assume that the prior for the weights may be given by a zero-mean Gaussian with a standard deviation of  $\sigma_w$ . We can then write the density of the weights as

$$p(\vec{\omega}) = \frac{1}{\sqrt{2\pi}\sigma_w} \exp\left(-\sum_{i=1}^W \frac{\omega_i^2}{2\sigma_w^2}\right) \quad (10)$$

where  $W$  is the total number of weights contained in the neural network and  $\omega_i$  is the value of the  $i^{\text{th}}$  network weight. Using similar substitutions to those introduced above in the previous Section,

$$\alpha = \frac{1}{\sigma_w^2} \quad (11)$$

and

$$Z_w(\alpha) = \left(\frac{2\pi}{\alpha}\right)^{\frac{W}{2}} \quad (12)$$

we can re-write this as

$$p(\vec{\omega}) = \frac{1}{Z_w(\alpha)} \exp(-\alpha E_w) \quad (13)$$

where we have defined the *weight error*,  $E_w$  as the sum of the squares of the weights,

$$E_w = \frac{1}{2} \sum_{i=1}^w \omega_i^2 \quad (14)$$

## Combining the Densities

Combining expressions for the two densities above we can now write the posterior distribution for the weights,

$$p(\mathcal{W} | D) = \frac{\exp(-M(\mathcal{W}))}{\int \exp(-M(\mathcal{W})) d\mathcal{W}} \quad (15)$$

where  $M(\mathcal{W})$  is defined

$$M(\mathcal{W}) = \beta E_D + \alpha E_W \quad (16)$$

Now that we have an expression for this posterior distribution, we can predict the target value,  $t$  of a new input vector,  $\mathcal{X}$  by using the integral

$$p(t | \mathcal{X}) = \int p(t | \mathcal{W}, D) p(\mathcal{W} | D) d\mathcal{W} \quad (17)$$

The problem with the above expression is that it implies that we evaluate the target for **every possible set of weights**,  $\mathcal{W}$ . Networks will usually contain 100 plus weights, so clearly, it is impractical to evaluate the above integral, as it would involve performing a 100 plus dimensional integration!

What is usually done instead is to assume that the quantity  $p(\mathcal{W} | D)$  is sharply peaked about some maximum value,  $\mathcal{W}_{MAP}$ . We can therefore approximate the above integral by evaluating it only about the immediate vicinity of the *maximum a posteriori* (MAP) estimate. We will discuss how to perform this integration in the next Section. First of all we shall discuss methods for finding the *maximum a posteriori* weights,  $\mathcal{W}_{MAP}$ .

## Finding the MAP Weights

From Equation (15) we see that in order to find the *maximum a posteriori probability*  $P(\mathcal{W} | D)$  of the weights  $\mathcal{W}$ , we must minimise the function

$$M(\mathcal{W}) = \beta E_D + \alpha E_W \quad (18)$$

with respect to the network weights,  $\mathcal{W}$ . Expanding the original terms of  $E_D$  and  $E_W$  we obtain

$$M(\mathcal{W}) = \frac{\beta}{2} \sum_{i=1}^N (f(\mathcal{X}_i, \mathcal{W}) - t_i)^2 + \frac{\alpha}{2} \sum_{i=1}^W \omega_i^2 \quad (19)$$

This function is familiar as the error function used to train neural networks with weight decay (IFS-TR-028). The weight decay parameter,  $\lambda$  is in this case given by

$$\lambda = \frac{\alpha}{\beta} \quad (20)$$

Thus, training a neural network using weight decay is, in fact, equivalent to attempting to find the MAP weights assuming a Gaussian prior on the weights and additive Gaussian noise on the target values. Furthermore, we are given an interpretation of the weight decay,  $\lambda$  as being the ratio of the expected variances of the noise on the targets of the training data and the prior of the weights.

## Using the MAP Estimate to Make Predictions and Estimate Error-Bars

Suppose we have trained a neural network to minimise the function  $M(\mathcal{W})$  as discussed above. At the end of training we hope to have found the maximum a posteriori set of weights  $\mathcal{W}_{MAP}$ . We now want to use this set of weights to provide predictions about a new input vector,  $\mathcal{X}$ . Recall that, in order to do this we wish to evaluate the integral,

$$p(t | \mathcal{X}, D) = \int p(t | \mathcal{X}, \mathcal{W}, D) p(\mathcal{W} | D) d\mathcal{W} \quad (21)$$

One option would be to assume that the value of  $p(\mathcal{W} | D)$  is 1.0 at the value  $\mathcal{W} = \mathcal{W}_{MAP}$  and 0.0 for all other values of  $\mathcal{W}$  (i.e. it is a delta function). We would therefore simply use the output of the final trained network as an approximation to the above integral. This approach corresponds to the usual non-Bayesian approach to training neural networks using simple weight decay.

A more sophisticated approach is to assume that about some local neighbourhood of  $\mathcal{W}_{MAP}$  we can make a quadratic approximation to the function  $M(\mathcal{W})$ , i.e.

$$M(\mathcal{W}) = M(\mathcal{W}_{MAP}) + \frac{1}{2} (\mathcal{W} - \mathcal{W}_{MAP})^T H (\mathcal{W} - \mathcal{W}_{MAP}) \quad (22)$$

where  $H$  is the Hessian of the regularised error function with respect to the weights. This can be written in terms of the Hessian of the unregularised error function,

$$H = \beta \nabla \nabla E_D + \alpha I \quad (23)$$

A number of techniques exist in the literature for the determination of the Hessian of the unregularised error function. For a summary, see IFS-TR-039.

Substituting the above expansion into Equation (21) and making a linear approximation for the output of the network about the MAP weights, i.e.

$$f(\mathbf{x}, \mathbf{w}_{MAP} + \Delta \mathbf{w}) = f(\mathbf{x}, \mathbf{w}_{MAP}) + \mathbf{g}^T \Delta \mathbf{w} \quad (24)$$

where  $\mathbf{g}$  is the standard error gradient evaluated about  $\mathbf{w}_{MAP}$

$$\mathbf{g} = \nabla_{\mathbf{w}} y |_{\mathbf{w}_{MAP}} \quad (25)$$

we obtain the final expression for a prediction of the target,  $t$  for a new input vector  $\mathbf{x}$ ,

$$p(t | \mathbf{x}, D) = \frac{1}{\sqrt{2\pi}\sigma_t} \exp\left(-\frac{(t - f(\mathbf{x}, \mathbf{w}_{MAP}))^2}{2\sigma_t^2}\right) \quad (26)$$

i.e. the probability of the predicted output is a Gaussian with expected mean equal to the output of the network with the MAP weights and variance given by

$$\sigma_t^2 = \frac{1}{\beta} + \mathbf{g}^T H^{-1} \mathbf{g} \quad (27)$$

So, in practical terms we use the network with the MAP weights to estimate the *expected* value of the target in the usual manner, but we can also associate with this estimate a measure of the uncertainty associated with this prediction.

## Finding Suitable Values for $\alpha$ and $\beta$

So far, we have not discussed how to obtain suitable estimates for the so called *hyper parameters*  $\alpha$  and  $\beta$ . We examine two possible techniques for this in Theoretical Report IFS-TR-031.

## Summary

In this report, the following has been shown

- Training using the standard weight decay is consistent with training a neural network using Bayesian techniques with a Gaussian prior assumed for the distribution of the weights of the network.
- The Bayesian analysis allows us to understand more precisely the role of the regularisation constant,  $\lambda$  in the weight decay algorithm. Specifically it allows us to view  $\lambda$  as representing the ratio of the expected variances of the target outputs and the weight priors.



- 
- It is possible, through the consideration of the Hessian matrix to calculate error bars or *uncertainties* in the output responses of the neural networks. This is particularly valuable in the case of financial time series prediction as it allows us to estimate the volatility of the predicted time series.

*Author: Darren Toulson*

*Revision: v 1.00 7th March, 1997*

*See ALSO: Theoretical Reports IFS-TR-027, IFS-TR-028, IFS-TR-031, IFS-TR-039*