



A Noninformative Prior for Neural Networks

HERBERT K.H. LEE

herbie@soe.ucsc.edu

*Department of Applied Math and Statistics, University of California at Santa Cruz, 1156 High Street,
Santa Cruz, CA 95064, USA*

Editors: Nando de Freitas, Christophe Andrieu, Arnaud Doucet

Abstract. While many implementations of Bayesian neural networks use large, complex hierarchical priors, in much of modern Bayesian statistics, noninformative (flat) priors are very common. This paper introduces a noninformative prior for feed-forward neural networks, describing several theoretical and practical advantages of this approach. In particular, a simpler prior allows for a simpler Markov chain Monte Carlo algorithm. Details of MCMC implementation are included.

Keywords: Bayesian statistics, improper prior, Markov chain Monte Carlo

1. Introduction

Much of the existing Bayesian neural network literature relies upon complex hierarchical priors for the parameters of the network, for example, MacKay (1992), Neal (1996), Rios Insua and Müller (1998), and Andrieu, de Freitas, and Doucet (2001). This paper proposes a fully noninformative prior with several philosophical and practical advantages. It also discusses how to implement model fitting via Markov chain Monte Carlo with this prior.

In the Bayesian context, one's choice of prior is meant to reflect either knowledge from previous experience or data, or to reflect personal and subjective beliefs about the problem (or both). In order to accurately choose a prior, one needs to understand how to interpret the parameters of a neural network model. Unfortunately, few people seem to realize how difficult a problem this is. Section 2 will discuss this issue in more detail. The bottom line is that if one cannot interpret the parameters effectively, it may be unwise to pretend so by using a proper prior. A noninformative prior is a way to acknowledge this high degree of uncertainty about the meaning of the parameters and to let the data more fully guide the model. Such priors were first fully established by Jeffreys (1961). There is now an extensive literature on such reference priors, and Kass and Wasserman (1996) provide a thorough review of this literature.

This noninformative prior is much simpler in structure, and it thus leads to a much simpler final model. This has advantages in reducing the dimension of the parameter space with the attendant benefits of shorter MCMC programming and run times. Since the noninformative prior is flat with respect to most parameters, the posterior modes of the parameters will be the same as would be obtained under least-squares fitting. This has the advantage that one's results are directly comparable to those of many other researchers, and that one can directly relate to those doing non-Bayesian analysis. However, one still gains all of the benefits of

the Bayesian approach. For example, it is straightforward to get estimates of variability or error probabilities from MCMC output, and one can justifiably interpret the results in terms of probabilities. This approach also allows one to do model averaging (Kass & Raftery, 1995) over several different fitted networks to improve predictive power.

Noninformative priors, including the one of this paper, often have many invariance properties (Jeffreys, 1961; Hartigan, 1964). For example, Jeffreys's prior is invariant with respect to all differentiable transformations of the input variables. The prior of this paper is easily seen to be invariant with respect to a shift in location or scale. Thus one need not be concerned about rescaling the variables (except, perhaps, for computational reasons).

This noninformative prior, as well as many other priors for neural networks, have been shown to be asymptotically consistent for the posterior, in the sense that the posterior probability accumulates in a Hellinger neighborhood of the truth (Lee, 2000a).

In many problems, including this one, it is important to consider the issue of model selection. How many hidden nodes are best? Are all of the inputs necessary? Noninformative priors, being not informative, do not provide the same shrinkage as some other priors (or methods such as weight decay), and so the user must keep this in mind. It is advisable to use some model selection technique, and these are discussed in Section 6.

This paper focuses exclusively on feed-forward neural networks with a single hidden layer of units with logistic activation functions, without direct connections from the inputs to the outputs. For regression problems, linear output units are used, while for classification, the softmax approach is taken. However, all of these methods and results are fully generalizable to any variety of feed-forward networks. The primary focus of this paper is the case of regression, although some classification examples are also presented. The particular form of the model in this paper in the case of regression is:

$$y_i = \beta_o + \sum_{j=1}^k \beta_j \Psi(\gamma'_j \mathbf{x}_i) + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

where x_i is a (possibly vector-valued) input, y_i is the target, γ are the weights from the inputs to the hidden layer, β are the weights from the hidden layer to the outputs, and the activation function is

$$\Psi(z) = \frac{1}{1 + \exp(-z)}.$$

2. Interpretation of the parameters

This section will demonstrate the infeasibility of trying to interpret the parameters of a neural network in the general case. Before getting to the general case, let us first look at some cases where the parameters can be interpreted, so that we can see what happens to cause a loss of interpretability.

The case of a single hidden node is straightforward. With a single input variable, the model for the fitted values, \hat{y}_i , is

$$\hat{y}_i = \beta_0 + \frac{\beta_1}{1 + \exp(-\gamma_0 - \gamma_1 x_i)}.$$

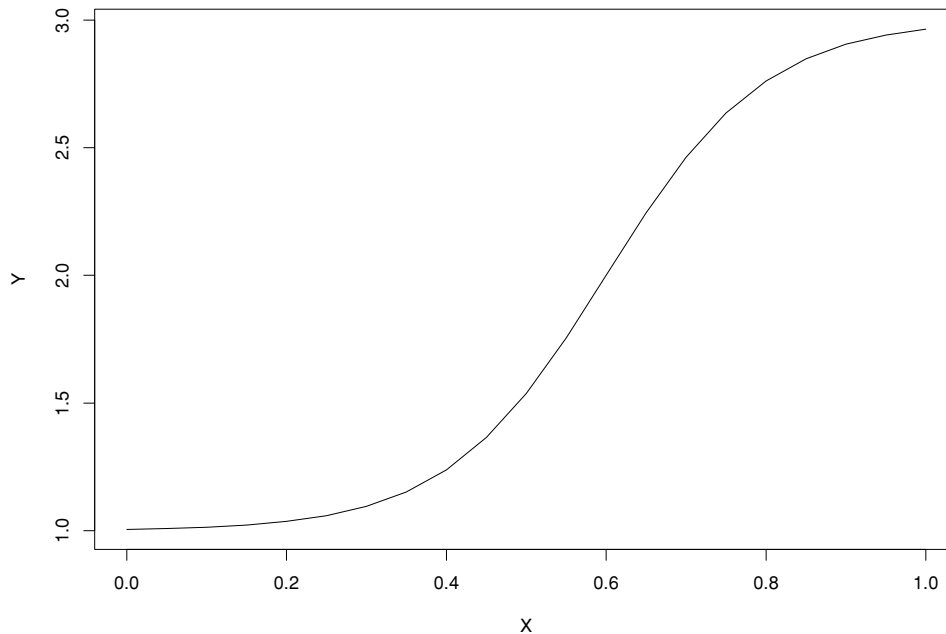


Figure 1. A one-node function.

Figure 1 shows this fitted function for $\beta_0 = 1$, $\beta_1 = 2$, $\gamma_0 = -6$, and $\gamma_1 = 10$ over x in the unit interval. In the one-hidden node case, the parameters are easily interpreted. β_0 represents the overall location of y , a sort of intercept, in that $y = \beta_0$ when the logistic function is close to zero (which it is here for values of x near or below zero). β_1 is the overall scale factor for y , in that the logistic function ranges from 0 to 1, so β_1 is like the range of y , which in this case is $3 - 1 = 2$. The γ parameters control the location and scale of the logistic function. The center of the logistic function occurs at $-\frac{\gamma_0}{\gamma_1}$, here 0.6. The larger the value of γ_1 , the steeper the increase in y as one moves away from the center. If γ_1 is positive, then the logistic rises from left to right. If γ_1 is negative, then the logistic will decrease as x increases.

This interpretation works fine for a single hidden node. For well-separated nodes (in terms of their centers and slopes, so that for any particular value of x , at most one node is near its center and the rest are far enough from their centers that a small change in x does not produce a noticeable change in the logistic function), one can continue to use this interpretation. Two nodes with centers near each other, but with opposite signs on γ_1 , can combine to produce a peak, similar in spirit to a kernel or a radial basis function. However, when the nodes effectively overlap, which is typically the case in a real problem, the parameters can no longer be interpreted as above.

For example, the left plot of figure 2 results from maximum likelihood (least-squares) fitting of a two-node network to the motorcycle accident data of Silverman (1985). The x -axis is the time in milliseconds after the crash, and the y -axis the acceleration force on the head of the rider. The two-node model does a good job of capturing the main features of

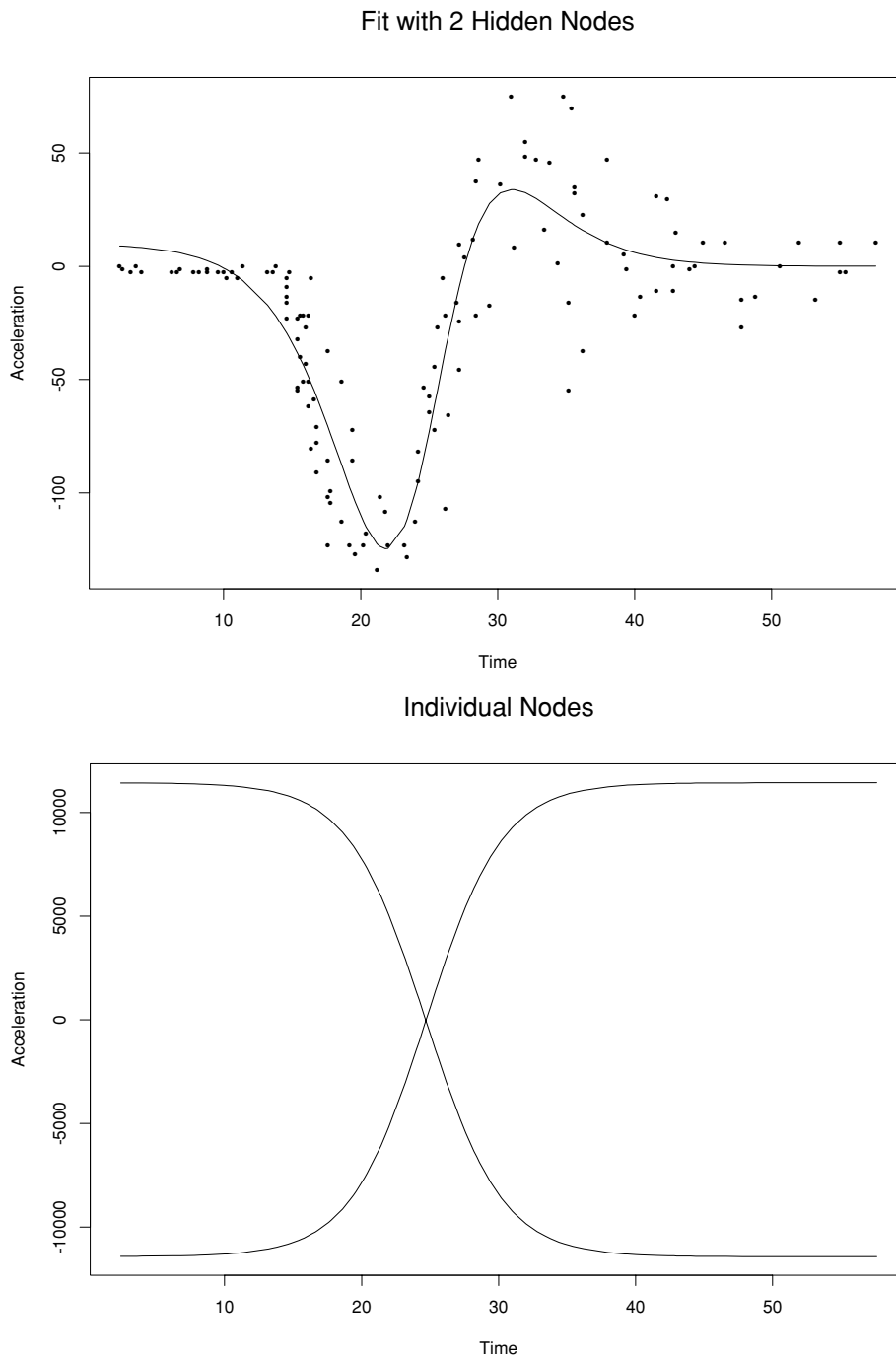


Figure 2. Logistic basis functions.

the data, although perhaps it over-smooths for small x values. Keep in mind that this fit is achieved with only two hidden nodes. It is impossible to use any of the above interpretations of the nodes for this fit, because there are more points of inflection in the fit than there are hidden nodes. The two individual nodes are shown in the right plot of the figure, and they have very similar centers and slopes, but combine to produce a highly non-linear fit. Note that the scale of the basis functions is two orders of magnitude larger than the original data. Looking at only the data, it is not at all clear that the particular basis functions of figure 2 would work so well. This is an extremely simple case, with only two hidden nodes and a single input variable, and one can easily imagine how much worse things can get with more hidden nodes.

Alternative models in the neural networks literature require the selection of many hyperparameters for their hierarchical priors. There is little guidance on how the selection of these hyperparameters affects the posterior in specific cases. Typically such models are meant to impose some shrinkage in the fitted model, but how much shrinkage can be allowed in the example of figure 2 before the least-squares solution becomes infeasible under the posterior, because the scale of the least-squares fit is so different than the scale of the data? In most cases, even a small amount of shrinkage will eliminate the least-squares fit of this figure. Yet a shrinkage model won't be easily able to fit the data as well, requiring more hidden nodes and increasing the computational complexity. A noninformative prior allows the data to drive the model, rather than the prior.

At this point, two comments should be made. First, in the limiting case of infinitely many hidden nodes, Neal (1996) has shown that the model converges to a Gaussian process model, and demonstrates how the prior parameters affect the posterior. This is indeed an elegant solution to understanding the prior. Yet in practice, many users of neural networks opt for a relatively small number of hidden nodes, and it is these finite models that this paper is meant to address. Second, prior shrinkage does have its benefits, in particular, it can help guard against overfitting (analogously to weight decay). One must take more care with a noninformative prior, and some suggestions are in Section 6.

3. An improper prior

The prior I propose for feed-forward neural networks is

$$\pi(\beta, \gamma, \sigma^2) \propto \frac{1}{\sigma^2}. \quad (1)$$

Notice that this prior is only defined up to a constant. Since it has an infinite integral, it does not matter what the constant is. Priors with an infinite integral, such as this one, are called *improper*. While the thought of using an improper prior may sound odd, improper priors have found to be very useful in practice. Many of Jeffreys's (1961) original reference priors were improper. The key is to ensure that the posterior is proper (finite). Below are some details to enforce propriety for this prior.

This choice of prior is well-established in least-squares regression (see for example, Gelman et al., 1995). It places uniform prior density over the real line for all parameters

except the variance of the error. For the variance, it is uniform with respect to the log of the variance, since variances must be positive and are right-skewed. By using a prior such as this one, we make no assumptions about the location or scale of any of the parameters in the model. Gelman et al. (1995) present some theoretical advantages of this family of priors.

If we condition on the hidden nodes (i.e., temporarily think of all hidden node parameters as fixed), then fitting the linear output weights is exactly a linear regression problem. It makes sense to thus apply results from the regression literature. The hidden nodes can be seen as basis functions, and a noninformative prior is then noninformative about the fitted function, allowing the data to drive the fit.

In the case of simple least-squares regression, the posterior is necessarily proper (assuming at least as many data points as parameters). In the case of a neural network, propriety of the posterior is not guaranteed. Hence we need to make a few minor modifications to the prior. The basic idea is to guarantee that the logistic basis functions of the hidden layer are linearly independent. We do this by putting certain restrictions on the parameters during the MCMC fitting process, a method now common in the mixture model literature (Diebolt & Robert, 1994; Wasserman, 2000).

First, a piece notation is helpful. Denote the outputs of the hidden layer as

$$z_{ij} = \left[1 + \exp \left(-\gamma_{j0} - \sum_{h=1}^p \gamma_{jh} x_{ih} \right) \right]^{-1}$$

and let Z be the matrix with elements (z_{ij}) . The fitting of the vector β is merely a least-squares regression on the design matrix Z . What is necessary is a restriction on Z to ensure the linear independence of the columns of Z . A sufficient condition for linear independence is that the determinant of $Z^T Z$ is larger than some positive value C_n . We can let $C_n \rightarrow 0$ as $n \rightarrow \infty$ as long as $C_n > 0$ for all n (n is the sample size). To ensure propriety of the posterior, we also need to bound the individual γ 's such that $|\gamma_{jh}| < D_n$ for all j, h . D_n is a bound which is allowed to grow with n , as discussed in the next paragraph. A full theoretical justification of these restrictions follows in Section 5.

It is computationally useful to also bound the other parameters (namely $|\beta_j| < D_n$; σ does not need to be bounded). This can help with numerical stability in the MCMC fitting. It is also theoretically useful in that this prior can be shown to be asymptotically consistent for the posterior (Lee, 2000a). The consistency proof requires that $D_n = o(\exp(n^r))$ for all $r > 0$.

The basic idea of the constraining values C_n and D_n is that as the sample size n increases, C_n and D_n can become more extreme (closer to zero for C_n , arbitrarily large for D_n). Asymptotically, the rate does not matter for C_n , and the rate for D_n required by the consistency theorem is exponentially fast. In practice, one can just pick a large constant value for D_n such as 100,000, and a small constant value for C_n such as 0.0001. What matters is the numerical stability of one's computational routines, and numbers such as these usually suffice for double precision arithmetic.

Thus the restricted prior is

$$\pi_n(\beta, \gamma, \sigma^2) = \pi(\beta, \gamma, \sigma^2) \mathbf{I}_{\{\theta \in \Omega_n\}} \propto \frac{1}{\sigma^2} \mathbf{I}_{\{\theta \in \Omega_n\}}. \quad (2)$$

where $I_{\{\}} is an indicator function (equal to one when its argument is true, zero otherwise), and Ω_n is the parameter space restricted such that:$

1. $|Z^T Z| > C_n$
2. $|\gamma_{jh}| < D_n$ for all j, h
3. $|\beta_j| < D_n$ for all j

4. Fitting the model with MCMC

Fitting the model is straightforward using Markov chain Monte Carlo. The full posterior is proportional to the likelihood times the prior:

$$f(\beta, \gamma, \sigma^2 | y) \propto \frac{1}{\sigma^2} (2\pi\sigma^2)^{-n/2} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^k \beta_j z_{ij}\right)^2\right] I_{\{\theta \in \Omega_n\}}.$$

The complete conditional distributions for σ^2 and β given the other parameters and the data are inverse-gamma and normal, respectively:

$$\sigma^2 | \gamma, y \sim \Gamma^{-1}\left(\frac{n-k-1}{2}, \frac{1}{2} \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \sum_{j=1}^k \hat{\beta}_j z_{ij}\right)^2\right) \quad (3)$$

$$\beta | \gamma, \sigma^2, y \sim N((Z^T Z)^{-1} Z^T y, (Z^T Z)^{-1} \sigma^2), \quad (4)$$

where $\hat{\beta}$ are the fitted coefficients from a least-squares regression of y on Z . Thus the MCMC algorithm is as follows:

1. Start with arbitrary initial values for γ .
2. Compute Z .
3. Draw σ^2 via Eq. (3).
4. Draw β via Eq. (4).
5. For each j from $1, \dots, k$, do the following Metropolis step:
 - (a) Generate a candidate $\tilde{\gamma}_j \sim N(\gamma_j, 0.05^2)$.
 - (b) Re-compute Z with $\tilde{\gamma}_j$ and compute $|Z^T Z|$.
 - (c) If $|Z^T Z| > 0.0001$ and $|\gamma_{jh}| < 100,000$ for all j, h , accept $\tilde{\gamma}_j$ with probability $\min(1, \frac{f(\beta, \tilde{\gamma}, \sigma^2 | y)}{f(\beta, \gamma, \sigma^2 | y)})$; otherwise, reject the candidate and keep the current value of γ_j .
6. Repeat steps 2 though 5 for the required number of iterations.

The resulting draws will be a sample from the joint posterior distribution of the parameters. An explanation of why the determinant restriction keeps the posterior from being improper follows in the next section. Note that the choice of $C_n = 0.0001$ and $D_n = 100,000$ are arbitrary and can be adjusted to match individual computational requirements. They generally do not affect the actual fit of the model, since the logistic function tails off very quickly (i.e., $1/(1 + \exp(-x))$ is numerically equivalent to 1.0 for $x > 40$.) The choice of

0.05 in the proposal function of step 5(a) should be treated as a tuning parameter, and can be adjusted to aim for an acceptance probability of approximately forty percent.

Many diagnostics exist for checking convergence, or at least lack thereof, of the chain. It is important to use only those values drawn from the stationary distribution of the chain, with the initial burn-in runs discarded. The reader is referred to Gilks, Richardson, and Spiegelhalter (1996) for a further discussion of convergence diagnostics, which includes several examples and additional references.

Note that there exist many more sophisticated methods of generating Metropolis-Hastings proposals to achieve higher acceptance probabilities and better mixing of the chain, such as Hybrid Monte Carlo (Neal, 1996). Such algorithms can also be applied to this model with their standard benefits. However, there is a trade-off between programming time and run time, in that more complicated approaches can decrease the necessary run time, but may take much more time and effort to program. Some concern has also been voiced about trying to make the MCMC iterations as independent as possible from each other. However, it is not normally necessary to use independent samples when estimating posterior quantities, such as means or predictions (see for example, Gilks, Richardson, & Spiegelhalter, 1996). The approach of this paper is to use a simpler algorithm and just let it run for more iterations, which often takes a similar amount of computing time in the end. In some cases, using a much simpler model leads to much faster run times (see the abalone example of Section 7.2). However, this is merely an issue of implementation and this noninformative prior can be used with any choice of fitting algorithm.

5. Theoretical justification for the restrictions

This section contains first a heuristic explanation of the restrictions, followed by more rigorous theory. To understand why the restrictions force propriety on the posterior, first look at a simplified scenario. Consider a single input x with n observations. Instead of using logistic functions, suppose the activation functions are step (Heaviside or indicator) functions of the form

$$\Psi_j(x) = I_{\{x \leq a_j\}} \quad \text{or} \quad \Psi_j(x) = I_{\{x > a_j\}}.$$

Suppose for now that network contains only two hidden nodes. Then $z_j = (z_{1j}, \dots, z_{nj}) = (1, \dots, 1, 0, \dots, 0)$ or $(0, \dots, 0, 1, \dots, 1)$ and z_0 is a vector of ones (the bias input). Let

$$r_i = \sum_j z_{ij} \quad \text{for } i = 1, 2; \quad r_{12} = \sum_i z_{i1} z_{2i}.$$

r_1 is the number of cases with input 1 larger (or smaller depending on the direction of inequality in the indicator) than the first threshold a_1 . r_{12} is the number of cases for which both hidden units give outputs of 1. The determinant of interest is

$$|Z'Z| = \begin{vmatrix} r_1 & r_{12} & r_1 \\ r_{12} & r_2 & r_2 \\ r_1 & r_2 & n \end{vmatrix} = r_1 r_2 n + 2r_1 r_2 r_{12} - r_1^2 r_2 - r_1 r_2^2 - r_{12}^2 n.$$

This determinant will be zero if any of the following happen: $r_i = 0$, $r_i = n$, $r_1 = r_2 = r_{12}$, or $\{r_1 = n - r_2 \text{ and } r_{12} = 0\}$. Essentially one needs to ensure that the thresholds for the activation functions are all separated by data points (i.e., there is a data point x_i that is between the thresholds of the two indicators), and that no threshold occurs outside the range of the data. These are exactly the conditions that prevent linear dependence. If we choose a small enough C_n (for indicator functions, all we need is $0 < C_n < 1$), then the same conditions guarantee that $|Z^t Z| > C_n$. The same logic applies to datasets with more inputs and to networks with more hidden nodes.

For indicator function hidden nodes, the requirement on the determinant also prevents impropriety in the posterior. Consider the case of trying to fit an indicator function where the parameter is the threshold. Suppose one tries to fit a threshold smaller than the smallest data point, x_{\min} . Then the data do not provide any information for distinguishing between putting the threshold at $x_{\min} - b_1$ and $x_{\min} - b_2$ for any positive numbers b_1 and b_2 . Since this equivalence set has infinite mass under the prior, the posterior is improper. On the other hand, when fitting a threshold inside the range of the data, the data force propriety on the posterior.

The indicator functions relate to logistic functions in that indicator functions are a limiting case of logistic functions. Let

$$\Psi(x) = \frac{1}{1 + \exp(-\gamma_0 - \gamma_1 x)}$$

be a logistic function. If $\gamma_0, \gamma_1 \rightarrow \infty$ such that $\frac{\gamma_0}{\gamma_1} \rightarrow a$ for some constant a , then $\Psi(x) \rightarrow I_{\{x > -a\}}$. We would then need the determinant condition to guarantee linear independence. However, this is not the only condition we need to guarantee a proper posterior. The triangular region where $\gamma_0, \gamma_1 \rightarrow \infty$ such that $\frac{\gamma_0}{\gamma_1} \rightarrow a$ has infinite area. Over this region, as the parameters get large, the likelihood converges to some non-zero constant (in most problems, the likelihood converges to zero in the tails). Thus the posterior over this region alone is improper. For this reason, we also need to bound the individual parameters.

The logistic functions allow values between zero and one, so that two columns of $Z^t Z$ could be very similar, but not identical. In regression this is called multicollinearity. The near-linear dependence causes instability in the parameter estimates. It is computationally desirable to avoid this case, which we can do by requiring the determinant to be larger than some small positive number C_n rather than merely requiring it to be non-zero.

Now for the full proof of the propriety of the prior. Denote the likelihood by L_n . Let π be the noninformative prior of Eq. (1), π_n be the restricted prior of Eq. (2), and Ω_n the restricted parameter space described with Eq. (2). Let C_n decrease to 0 (for example, $C_n = 1/n$) and let D_n increase with n (for example, $D_n = 100,000 + n$). First the likelihood is re-written so that β can be integrated out, which involves completing the square for β .

$$\begin{aligned} L_n &= f(\beta, \gamma, \sigma | \mathbf{y}) = (2\pi\sigma^2)^{-n/2} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n \left(\sum_{j=0}^k \beta_j z_j - y_i\right)^2\right] \\ &= (2\pi\sigma^2)^{-n/2} \exp\left[-\frac{1}{2\sigma^2} (\mathbf{Z}\beta - \mathbf{Y})^t (\mathbf{Z}\beta - \mathbf{Y})\right] \quad \text{in vector notation} \end{aligned}$$

$$\begin{aligned}
&= (2\pi\sigma^2)^{-\frac{k+1}{2}} |\mathbf{Z}'\mathbf{Z}|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2}[\beta - (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{Y}]'(\mathbf{Z}'\mathbf{Z})[\beta - (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{Y}]\right\} \\
&\quad * (2\pi\sigma^2)^{-\frac{n-(k+1)}{2}} |\mathbf{Z}'\mathbf{Z}|^{\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2}[\mathbf{Y}'\mathbf{Y} - \hat{\mathbf{Y}}'\hat{\mathbf{Y}}]\right\} \\
&= f(\beta | \gamma, \sigma, \mathbf{y})f(\gamma, \sigma | \mathbf{y}),
\end{aligned}$$

where $\hat{\mathbf{Y}}$ is the vector of fitted values, $\hat{\mathbf{Y}} = E[\mathbf{Y} | \mathbf{X}]$. Note that $f(\beta | \gamma, \sigma, \mathbf{y})$ is a proper density as long as $|\mathbf{Z}'\mathbf{Z}| > 0$, which is true over Ω_n . Denote by Γ_n the subspace of Ω_n that relates to all of the γ parameters. Then the posterior is proper:

$$\begin{aligned}
\int L_n \pi_n &= \int_{\Omega_n} f(\beta | \gamma, \sigma, \mathbf{y})f(\gamma, \sigma | \mathbf{y}) \left[\frac{1}{\sigma^2}\right] d\beta d\sigma d\gamma \\
&= \int_{\Gamma_n} \int \left[\int f(\beta | \gamma, \sigma, \mathbf{y}) d\beta \right] \frac{1}{\sigma^2} f(\gamma, \sigma | \mathbf{y}) d\sigma d\gamma \\
&= \int_{\Gamma_n} \int \frac{1}{\sigma^2} (2\pi\sigma^2)^{-\frac{n-(k+1)}{2}} |\mathbf{Z}'\mathbf{Z}|^{\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2}[\mathbf{Y}'\mathbf{Y} - \hat{\mathbf{Y}}'\hat{\mathbf{Y}}]\right\} d\sigma d\gamma \\
&\leq \int_{\Gamma_n} \int \frac{1}{\sigma^2} (2\pi\sigma^2)^{-\frac{n-(k+1)}{2}} |\mathbf{Z}'\mathbf{Z}|^{\frac{1}{2}} d\sigma d\gamma \\
&= (2\pi)^{-\frac{n-(k+1)}{2}} (n-k+1)^{-1} \int_{\Gamma_n} |\mathbf{Z}'\mathbf{Z}|^{\frac{1}{2}} d\gamma
\end{aligned}$$

The last integral is finite because Γ_n is a bounded set and the integrand is finite. Thus the posterior is proper.

In addition to showing that the adjusted prior leads to a proper posterior, it is also important to show that the adjusted prior is asymptotically equivalent to the original improper prior, which can be shown in both a global and local sense. First, it is clear that, for any compact set κ ,

$$\int_{\kappa} |\pi_n - \pi| = \int_{\kappa} |\pi I_{\Omega_n} - \pi| \rightarrow 0 \text{ as } n \rightarrow \infty$$

because $|\mathbf{Z}'\mathbf{Z}|$ must be non-zero for the true function (or else it would have one fewer node), and because for a large enough n , Ω_n will contain all elements of κ that satisfy the determinant condition. This equation says that, in the limit as the sample size grows, π_n converges to π on all compact sets. In this sense, the two priors are ‘‘asymptotically globally equivalent’’ (Wasserman, 2000).

Second is a condition of ‘‘asymptotic local equivalence’’. It relates to correct second-order frequentist coverage properties (Wasserman, 2000). The key is that the original and adjusted priors have the same local properties (while the adjusted prior is better behaved in the tails). Suppose there exists a true value of the parameters, θ_0 . Then for large n ,

$$\left| \frac{\partial \log \pi_n}{\partial \theta_0} - \frac{\partial \log \pi}{\partial \theta_0} \right| = O_p\left(\frac{1}{\sqrt{n}}\right)$$

because if n is large enough, θ_0 will be contained in Ω_n .

6. Avoiding overfitting

Overfitting is a traditional problem faced by data analysts. Several methods have been developed for combating overfitting in the case of neural networks, including weight decay, early stopping, and Bayesian methods. A popular device in the Bayesian context is Automatic Relevance Determination (ARD) (MacKay, 1992; Neal, 1996), which adds another layer to the hierarchical prior to adjust the weights for the relative importance of the inputs. In a fully Bayesian paradigm, one puts a prior over the space of possible models and explores the posterior for the model space, either choosing the model of highest posterior probability or doing model averaging, i.e., predicting with a weighted average of predictions from multiple models where the weights are the posterior probabilities of the models. Some methodology for reversible-jump MCMC is presented in Müller and Rios Insua (1998) and Andrieu, de Freitas, and Doucet (2001).

An alternative Bayesian approach is to approximate the posterior probabilities directly. For a review of the difficulties in estimating these posterior probabilities for a neural network, see Lee (2002). If the model space is sufficiently small, one can simply estimate probabilities for all models exhaustively. Should the model space be too large to estimate all possible models, one can use a searching technique such as Bayesian Random Searching (BARS). A brief description BARS is provided here, and more details can be found in Lee (2000b, 2001). BARS is motivated by Markov Chain Monte Carlo Model Composition (MC³) (Raftery, Madigan, & Hoeting, 1997), which uses MCMC on the model space, estimating posterior probabilities of the models with the fraction of time the Markov chain spends visiting each model. In order to compute the transition probabilities for this chain, one needs the ratio of the posterior probabilities of the two models, or equivalently, their Bayes factor. This can be approximated using the BIC (Lee, 2002). Since these posterior probabilities are estimated as part of MC³, instead of discarding them and relying on the steady state properties of the chain, BARS simply keeps track of the probabilities of all models visited, using these probabilities as the final estimates, and using the chain merely as a stochastic method for exploring the model space. This approach is similar to that of Chipman, George, and McCulloch (1998) in the context of CART.

7. Examples

7.1. Robot arm data

The canonical dataset in the Bayesian neural network literature is the robot arm data of MacKay (1992). The idea of the data is to model the relationship between the two joint angles of the arm (x_1 and x_2) and the resulting arm position in Cartesian coordinates (denoted by y_1 and y_2). In this case, the data were actually simulated from the true functions and Gaussian noise was added. The true model is

$$y_1 = 2.0 \cos(x_1) + 1.3 \cos(x_1 + x_2) + \varepsilon_1 \quad (5)$$

$$y_2 = 2.0 \sin(x_1) + 1.3 \sin(x_1 + x_2) + \varepsilon_2, \quad (6)$$

Table 1. Comparison of methods on the robot arm data.

Method	Mean square error on test data	R^2
MacKay	0.00557	0.9966
Neal	0.00549	0.9967
Müller and Insua	0.00620	0.9963
Andrieu, de Freitas, and Doucet	0.00502	0.9970
This paper	0.00501	0.9970

where $\varepsilon_i \stackrel{iid}{\sim} N(0, 0.05^2)$. The data are divided into two groups: a training set of 200 observations, and a test set of 200 observations. The models are fit using only the training set, and then the models can be validated on the test set.

Fitting models with each of two through sixteen hidden nodes and estimating the posterior probabilities of these models via the BIC approximation (Lee, 2002) finds that a six hidden nodes network has about 98% of the posterior probability of the model space. MCMC was run with the methods of Section 4 using 20,000 burn-in iterations (far more than necessary) and 20,000 posterior samples. The MCMC output was then used to fit the model on the 200 cases in the test data set. The theoretical optimum value for the mean squared error (MSE) is 0.005, and the model of this paper has an MSE of 0.00501. The small mean square error shows that this model does indeed fit quite well.

This dataset has also been analyzed by several others in the Bayesian neural network literature, in particular MacKay (1992), Neal (1996), Müller and Rios Insua (1998a) and Andrieu, de Freitas, and Doucet (2001). Table 1 shows the MSE and R^2 achieved on the test data by the methods of each of the above competing models. All of these methods have an error rate on the same order of magnitude, although the methods of this paper do well within this group of methods.

7.2. Abalone data

This dataset, available at the UCI Machine Learning Repository (<http://www.ics.uci.edu/~mllearn/MLRepository.html>), involves predicting the age of abalone from eight physical measurements (Nash et al., 1994). Here the response is treated as a continuous variable. The data were divided into a training set of the first 3133 cases and a test set of the last 1044 cases. For comparison, I fit both the model of this paper (BARS chooses seven inputs and two hidden nodes) and the model of Neal using his neural network code (Neal, 1996, available at <http://www.cs.toronto.edu/~radford/fbm.software.html>). For Neal's code, multiple combinations of nodes and prior hyperparameter values were tried with and without ARD. The best found used 16 hidden nodes with ARD. Table 2 shows the results of predicting on the test set. Both methods achieve similar fits. However in terms of running time, Neal's code took six hours to run 250 burn-in runs and 750 posterior samples, while the algorithm of this paper took only twelve minutes to run 10,000 burn-in runs and 20,000 posterior samples on the same Compaq Alpha processor.

Table 2. Comparison of error rates on the Abalone data.

Method	Mean square error on test data	R^2
Hierarchical neural network	3.89	0.34
Neural network of this paper	3.93	0.33

Table 3. Comparison of error rates on the median house price data.

Method	Mean square error on test data	R^2
Hierarchical neural network	6.10×10^9	0.54
Neural network of this paper	3.65×10^9	0.72

7.3. Median house price data

A third regression example is the median house price by census block in California (Pace & Barry, 1997), available at Statlib (<http://lib.stat.cmu.edu/datasets/>). Eight input variables are available for 20,640 cases, which were equally divided into a training set and a test set. Results comparing the code of Neal (1996) for a network of 20 nodes with ARD to the model of this paper (using four hidden nodes) are shown in Table 3. Note that the noninformative prior model does significantly better.

7.4. Diabetes data

Ripley (1996) provides an analysis of a dataset on diabetes in Pima Indian women (the data are also available at the UCI Repository, <http://www.ics.uci.edu/~mlearn/MLRepository.html>). The idea is to predict the presence of diabetes using seven health covariates. There are 532 complete records,¹ of which 200 are used as a training set and the other 332 are used as a test set. About 33% of the population has diabetes. The methods of this paper are easily applied to classification examples, either by changing the output node from linear to logistic (or using the Softmax model for multiple categories; this approach requires Metropolis-Hastings updates for the output weights), or by using a latent added variable approach (as in Albert & Chib, 1993); additional details are in Lee (1998).

BARS finds that the best model uses only one hidden node and four explanatory variables: number of pregnancies, plasma glucose, body mass, and pedigree. The error rate on the test set with this model is 65 of 332, or 19.6%. For comparison, Table 4 reports misclassification rates as reported by Ripley (1996) for a number of other analyses, as well as the result of running Neal's neural network code (Neal, 1996). Again for Neal's code, multiple combinations of nodes and prior hyperparameter values were tried, and the model shown was the best found, which used 18 hidden nodes with ARD.

Table 4. Comparison of error rates on the Pima Indian diabetes data.

Standard linear discrimination	20.2%
Robust linear discrimination	22.9%
Logistic regression	19.9%
Multivariate adaptive regression splines (MARS)	22.6%
Projection pursuit regression (PPR)	22.6%
Multi-layer neural network	22.6%
Nearest neighbor with CV	24.7%
Classification tree	24.4%
Hierarchical neural network	20.2%
Neural network of this paper	19.6%

Table 5. Comparison of error rates on the loan applications data.

Logistic regression	35%
Nearest neighbor	36%
CART	31%
Hierarchical neural network	29%
Neural network of this paper	31%

7.5. Loan applications data

Lee (2001) describes a dataset on applications for personal loans at a bank. One can attempt to predict whether the loan was approved or denied from 23 covariates. There is a large amount of correlation between covariates (e.g., someone with a mortgage will be a homeowner; a person cannot have lived in their current residence for more years than they are old), so some model selection technique is required. 53% of the applications were denied when reviewed by a trained loan officer. The data are split into a training set of 4000 observations and a test set of the remaining 4508 observations. Using the training set, BARS finds that the optimal model uses seven explanatory variables and two hidden nodes. Table 5 compares the misclassification rates of various methods on the test data. Note that the data are very messy, and not all information available to loan officers has been coded in the dataset, so one can only expect large error rates.

8. Conclusions

Modern Bayesian methods include noninformative priors in the toolbox. This paper is meant to help bring this perspective and practice to the neural network community. Noninformative priors have both philosophical and practical advantages over both complicated hierarchical priors and over non-Bayesian methods. They simplify the choice of prior, the model, the computational details, and the interpretation of the results. Markov chain Monte Carlo

algorithms can be much simpler when a noninformative prior is used. However, it is not a free lunch—overfitting can become more of a problem with noninformative priors than might be the case in other methods, so the user should consider a model selection/shrinkage technique to complement the model (as one should in general). The methods of this paper are easily extended to other varieties of feed-forward networks. I hope that these ideas will be found useful in other areas of the machine learning community as well.

Acknowledgments

This work was partially supported by National Science Foundation grant DMS-9803433 and National Institutes of Health grant RO1 CA54852-08. The author is grateful to Larry Wasserman for all of his help, and to the editors and two anonymous referees who provided many useful comments.

Note

1. The UCI web site documentation stated that there is no missing data, however some records have zeroes that are not possible, such as a zero blood pressure.

References

- Albert, J. H., & Chib, S. (1993). Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88, 669–679.
- Andrieu, C., de Freitas, N., & Doucet, A. (2001). Robust full Bayesian learning for radial basis networks. *Neural Computation*, 13:10, 2359–2407.
- Chipman, H., George, E., & McCulloch, R. (1998). Bayesian CART model search (with discussion). *Journal of the American Statistical Association*, 93, 935–960.
- Diebolt, J., & Robert, C. (1994). Estimation of finite mixture distributions through Bayesian sampling. *Journal of the Royal Statistical Society B*, 56, 363–375.
- Gelman, A., Carlin, J. B., Stern, H. S., & Rubin, D. B. (1995). *Bayesian data analysis*. London: Chapman and Hall.
- Gilks, W. R., Richardson, S., & Spiegelhalter, D. J. (1996). *Markov chain Monte Carlo in practice*. London: Chapman and Hall.
- Hartigan, J. (1964). Invariant prior distributions. *Annals of Mathematical Statistics*, 35:2, 836–845.
- Jeffreys, H. (1961). *Theory of probability*, 3rd edn. New York: Oxford University Press.
- Kass, R. E., & Wasserman, L. (1996). The selection of prior distributions by formal rules. *Journal of the American Statistical Association*, 91:435, 1343–1370.
- Kass, R. E., & Raftery, A. E. (1995). Bayes factors. *Journal of the American Statistical Association*, 90:430, 773–795.
- Lee, H. K. H. (1998). Model selection and model averaging for neural networks. Ph.D. thesis, Carnegie Mellon University, Department of Statistics.
- Lee, H. K. H. (2000a). Consistency of posterior distributions for neural networks. *Neural Networks*, 13:6, 629–642.
- Lee, H. K. H. (2000b). A framework for nonparametric regression using neural networks. Technical Report 00-32, Duke University, ISDS.
- Lee, H. K. H. (2001). Model selection for neural network classification. *Journal of Classification*, 18:2, 227–243.
- Lee, H. K. H. (2002). Difficulties in estimating the normalizing constant of the posterior for a neural network. *Journal of Computational and Graphical Statistics*, 11:1, 222–235.

- MacKay, D. J. C. (1992). Bayesian methods for adaptive methods. Ph.D. thesis, California Institute of Technology, Program in Computation and Neural Systems.
- Nash, W. J., Sellers, T. L., Talbot, S. R., Cawthorn, A. J., & Ford, W. B. (1994). The population biology of abalone (*haliotis* species) in Tasmania. I. Blacklip abalone (*h. rubra*) from the North Coast and Islands of Bass Strait. Technical Report 48, Sea Fisheries Division, Australia.
- Neal, R. M. (1996). Bayesian learning for neural networks. New York: Springer.
- Pace, R. K. & Barry, R. (1997). Sparse spatial autoregressions. *Statistics and Probability Letters*, 33, 291–297.
- Raftery, A. E., Madigan, D., & Hoeting, J. A. (1997). Bayesian model averaging for linear regression models. *Journal of the American Statistical Association*, 437, 179–191.
- Rios Insua, D., & Müller, P. (1998). Feedforward neural networks for nonparametric regression. In D. Dey, P. Müller, & D. Sinha (Eds.), *Practical nonparametric and semiparametric bayesian statistics* (pp. 181–193). New York: Springer-Verlag.
- Ripley, B. D. (1996). *Pattern recognition and neural networks*. Cambridge: Cambridge University Press.
- Silverman, B. W. (1985). Some aspects of the spline smoothing approach to non-parametric curve fitting. *Journal of the Royal Statistical Society B*, 47, 1–52.
- Wasserman, L. (2000). Asymptotic inference for mixture models by using data-dependent priors. *Journal of the Royal Statistical Society B*, 62, 159–180.

Received June 8, 2000

Revised July 3, 2001

Accepted September 5, 2001

Final manuscript September 5, 2001