

Bayesian Neural Networks for Image Analysis

Aki Vehtari and Jouko Lampinen
Laboratory of Computational Engineering,
Helsinki University of Technology,
P.O.Box 9400, FIN-02015, HUT, Finland

Abstract

We demonstrate the advantages of using Bayesian neural networks for image analysis. The Bayesian approach provides consistent way to do inference by combining the evidence from data to prior knowledge from the problem. A practical problem with neural networks is to select the correct complexity for the model, i.e., the right number of hidden units or correct regularization parameters. The Bayesian approach offers efficient tools for avoiding overfitting even with very complex models, and facilitates estimation of the confidence intervals of the results. In this contribution we review the Bayesian methods for neural networks and present comparison results from case studies in process tomography and image segmentation. In the first case, neural networks were used to solve the inverse problem in electrical impedance tomography. The Bayesian networks provided consistently better results than other methods. In the second case, the goal was to locate trunks of trees in forest scenes. With Bayesian network it was possible to use large number of potentially useful features and prior for determining the relevance of the features automatically.

1 Introduction

A universal task in many areas of image analysis is to infer some needed piece of information from measurements that only partly determine the information.

Profound example of such problems is the perception of the three dimensional structure of view from two dimensional projection, i.e., image, with means like shape from shading, binocular stereopsis, and perception of perspective. Another example is restoration of images from blurred or noisy recordings. Also, classification and segmentation of image regions or objects based on a set of precomputed features is similar problem, as the features are often insufficient for uniquely separating the classes.

Recently Bayesian approaches have shown considerable potential in such problems. In the Bayesian approach prior information from the problem is combined to the evidence from the data, giving the posterior prob-

ability of the solutions. Predictions are made by integrating over this posterior distribution. In case of insufficient data the prior dominates the solution, and the effect of the prior diminishes with increased evidence from the data. In one of the pioneering works by Geman et.al. [1], Bayesian approach was developed for image restoration. This work also introduced the Gibbs sampling technique for computing posterior distributions. The advantages of Bayesian approaches in computational modeling of perception are discussed in [2].

In classification and non-linear function approximation neural networks have become very popular in recent years. With neural networks the main difficulty is controlling the complexity of the model. Another problem of standard neural network models is the lack of tools for analyzing the results (confidence intervals, like 10 % and 90 % quantiles, etc.).

For neural networks, MacKay introduced Bayesian approach [3] based on Gaussian approximation. Recently Neal introduced hybrid Monte Carlo method [4] that facilitates Bayesian learning for neural networks with no compromising approximations. The main advantages of Bayesian neural networks are:

- Automatic complexity control: Bayesian inference techniques allow the values of regularization coefficients to be selected using only the training data, without the need to use separate training and validation data.
- Possibility to use prior information and hierarchical models for the hyperparameters.
- Predictive distributions for outputs.

In this contribution we demonstrate the advantages of Bayesian neural networks in two case problems. In section 3 we give a review of the Bayesian methods for neural networks. In section 4 we report results on using Bayesian networks for image reconstruction in electrical impedance tomography. In section 5 we present results comparing Bayesian networks and other classification methods for classification of objects in forest scenes.

2 Multi Layer Perceptron

In this section we briefly review Multi Layer Perceptron (MLP) neural network. See [5] for thorough introduction to MLPs. We concentrate here to one hidden layer MLP networks with hyperbolic tangent (tanh) activation function, but Bayesian methods described can be used for other types of neural networks, like RBF networks, too. Basic MLP network model with k outputs is

$$f_k(x, w) = w_{k0} + \sum_{j=1}^m w_{kj} \tanh \left(w_{j0} + \sum_{i=1}^d w_{ji} x_i \right), \quad (1)$$

where x is a d -dimensional input vector, w denotes weights and indices i and j correspond to hidden and output units, respectively.

MLP is often considered as a generic semiparametric model, which means that the effective number of parameters may be less than the number of available parameters. Effective number of parameters determines the complexity of the model. For small weights the network mapping is almost linear and has low effective complexity, since the central region of sigmoidal activation function can be approximated by linear transformation. Traditionally complexity of MLP has been controlled with early stopping or weight decay [5].

In early stopping weights are initialized to very small values. Part of the training data is used to train the MLP and the other part is used to monitor the validation error. Iterative optimization algorithms used for minimizing the training error gradually take parameters in use. Training is stopped when the validation error begins to increase. Since training is stopped before a minimum of the training error, the effective number of parameters remains less than the number of available parameters.

Intuitively, the optimization algorithm first fits the model more to the underlying process and when the optimization is continued it fits more to the noise in the training set. With early stopping, optimization is tried to be stopped before too much fitting to the noise has occurred.

The basic early stopping is rather inefficient, as it is very sensitive to the initial conditions of the network and only part of the available data is used to train the model. These limitations can easily be alleviated by using a committee of early stopping networks, with different partitioning of the data to training and stopping sets for each network. When used with caution MLP early stopping committee is good baseline method for neural networks.

In weight decay penalizing term is added to the error function. Using sum of squares of weights the weights are encouraged to be small. In practice each layer in an

MLP should have different regularization parameter [5], giving the penalty term

$$\alpha_1 \sum_{j,i} w_{ji}^2 + \alpha_2 \sum_{j,k} w_{kj}^2. \quad (2)$$

Problem is how to select good values for α_i . Traditionally this has been done with cross validation (CV). Since CV gives noisy estimate for error, it does not guarantee that good values for α_i can be found. Also it becomes easily computationally prohibitive as computational expenses grow exponentially with number of parameters to be selected.

3 Bayesian Learning for MLP

Bayesian methods use probability to quantify uncertainty in inferences and the result of Bayesian learning is a probability distribution expressing our beliefs regarding how likely the different predictions are. Bayesian paradigm offers consistent way to do inference using models with even very large number of parameters. See e.g. [6] for good introduction to Bayesian methods.

3.1 Bayesian Learning

Consider a regression or classification problem involving the prediction of a noisy vector y of target variables given the value of a vector x of input variables.

The process of Bayesian learning is started by defining a model \mathcal{M} , and prior distribution $p(\theta)$ for the model parameters. Prior distribution expresses our initial beliefs about parameter values, before any data has observed.

After observing new data $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$, prior distribution is updated to the posterior distribution using Bayes' rule

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)} \propto L(\theta|D)p(\theta), \quad (3)$$

where the likelihood function $L(\theta|D)$ gives the probability of the observed data as function of the unknown model parameters.

To predict the new output $y^{(n+1)}$ for new input $x^{(n+1)}$, predictive distribution is obtained by integrating the predictions of the model with respect to the posterior distribution of the model parameters

$$p(y^{(n+1)}|D) = \int p(y^{(n+1)}|x^{(n+1)}, \theta)p(\theta|D)d\theta. \quad (4)$$

This is same as taking the average prediction of all the models weighted by their goodness.

Note that predictive distribution for $y^{(n+1)}$ is implicitly conditioned on hypotheses that hold throughout – no probability judgments can be made in vacuum [6] –

and to be more explicit notation as the following might be used

$$p(y^{(n+1)}|D, H) = \int p(y^{(n+1)}|x^{(n+1)}, \theta, H)p(\theta|D, H)d\theta, \quad (5)$$

where H refers to the set of hypotheses or assumptions used to define the model.

3.2 Models

As noted above we start from defining a model for our problem. Statistical model is defined with the likelihood function, which in case of independent and exchangeable data points is given by

$$L(\theta|D) = \prod_{i=1}^n p(y^{(i)}|x^{(i)}, \theta), \quad (6)$$

where n is the number of data points.

In the likelihood equation the term $p(y^{(i)}|x^{(i)}, \theta)$ depends on our problem. In regression problems, it is generally assumed that the distribution of target data can be described by a deterministic function of inputs, corrupted by additive Gaussian noise of a constant variance. Probability density for a target y_j is then

$$p(y_j|x, w, \sigma) = \frac{1}{\sqrt{2\pi}\sigma_j} \exp(-(y - f_j(x, w))^2/2\sigma_j^2), \quad (7)$$

where σ_j^2 is the noise variance for the target. See [4, 7] for input dependent noise models. For a two class classification (logistic regression) model, the probability that a binary-valued target, y_j , has the value 1 is

$$p(y_j = 1|x, w) = [1 + \exp(-f_j(x, w))]^{-1} \quad (8)$$

and for many class classification (softmax) model, the probability that a class target, y , has value j is

$$p(y = j|x, w) = \exp(f_j(x, w)) / \sum_k \exp(f_k(x, w)). \quad (9)$$

In equations (7), (8) and (9) function $f_j(x, w)$ is in this case an MLP network. Traditionally in many methods one of the problems has been to find a good topology for the MLP. In Bayesian approach we could use infinite number of hidden units [4]. We do not need to restrict the size of the MLP based on the size of the training set, but in practice, we will have to use finite number of hidden units due to computational limits. MCMC methods (section 3.5) produce the correct answer eventually, but it may sometimes take unreasonable amount of time[4].

3.3 Priors

Next, we have to define the prior information about our model parameters, before any data has been seen. Usual prior is that the model has some unknown complexity but the model is not constant or extremely flexible. To express this prior belief we set hierarchical model specification.

Parameters w define the model $f(x, w)$. As discussed in section 2, complexity of the MLP network can be controlled by controlling the size of the weights w . Corresponding prior to weight decay is to use Gaussian prior distribution for weights w given hyperparameter α

$$p(w|\alpha) = (2\pi)^{-m/2} \alpha^{m/2} \exp(-\alpha \sum_{i=1}^m w_i^2/2). \quad (10)$$

This prior states that smaller weights are more probable, but how much more is determined by the value of hyperparameter α . Since we do not know the correct value for hyperparameter α , we set a vague hyperprior $p(\alpha)$ expressing our belief that complexity controlled by α is unknown but the model is not constant or extremely flexible. A convenient form for this hyperprior is vague Gamma distribution with mean μ and shape parameter a

$$p(\alpha) \sim \text{Gamma}(\mu, a) \propto \alpha^{a/2-1} \exp(-\alpha a/2\mu). \quad (11)$$

In order to have prior for weights which is invariant under the linear transformations of data, separate priors (each having its own hyperparameters α_i) for different weight groups in each layer of a MLP are used.

In MLP networks, the weights from less important inputs are typically smaller than weights from more important inputs¹. Prior belief that some inputs are likely to be more relevant than others can be implemented by using different priors for weight groups from each input, and hierarchical hyperpriors for these priors. The posteriors for hyperparameters should then adjust according to relevance of the inputs. This prior is called Automatic Relevance Determination (ARD) [8, 4].

For regression models we need prior for σ in equation (7), which is conveniently specified in terms of corresponding precision, $\tau = \sigma^{-2}$. As for α , our prior information is usually quite vague, stating that σ is not zero or extremely large. This prior can be expressed with vague Gamma-distribution with mean μ and shape parameter a

$$p(\tau) \sim \text{Gamma}(\mu, a) \propto \tau^{a/2-1} \exp(-\tau a/2\mu). \quad (12)$$

¹Note that in the non-linear network the effect of an input may be small even if the weights from it are large and vice versa, but in general the size of the weights roughly reflects the relevance of the input.

3.4 Prediction

After defining the model and prior information, we combine the evidence from the data to get the posterior distribution for the parameters

$$p(w, \alpha, \tau|D) \propto L(w, \alpha, \tau|D)p(w, \alpha, \tau) . \quad (13)$$

Predictive distribution for new data is then obtained by integrating over this posterior distribution

$$p(y^{(n+1)}|x^{(n+1)}, D) = \int p(y^{(n+1)}|x^{(n+1)}, w, \alpha, \tau)p(w, \alpha, \tau|D) dw\alpha\tau \quad (14)$$

We can also evaluate expectations of various functions with respect to the posterior distribution for parameters. For example in regression we may evaluate the expectation for a component of $y^{(n+1)}$

$$\hat{y}_k^{(n+1)} = \int f_k(x^{(n+1)}, w)p(w, \alpha, \tau|D) dw\alpha\tau , \quad (15)$$

which corresponds to the best guess with squared error loss.

The posterior distribution for the parameters $p(w, \alpha, \tau|D)$ is typically very complex, with many modes. Evaluating the integral of equation (15) is therefore a difficult task.

The integral can be approximated with Gaussian approximations to modes. Then predictive distribution is approximated by the corresponding integral with respect to the Gaussian [3, 9]. Or we can use Monte Carlo methods, described next, to numerically approximate the integral.

3.5 Markov Chain Monte Carlo method

Recently, Neal has introduced implementation of Bayesian learning for neural networks in which the difficult integration of equation (15) is performed using Markov Chain Monte Carlo (MCMC) methods [4]. In [10] there is a good introduction to basic MCMC methods and many applications in statistical data analysis.

MCMC methods make no assumptions about the form of the posterior distribution. They may in some circumstances require a very long time to converge to the desired distribution.

The integral of equation (15) is the expectation of function $f_k(x^{(n+1)}, w)$ with respect to the posterior distribution of the parameters. This and other expectations can be approximated by Monte Carlo method, using a sample of values $w^{(t)}$ drawn from the posterior distribution of parameters

$$\hat{y}_k^{(n+1)} \approx \frac{1}{N} \sum_{t=1}^N f_k(x^{(n+1)}, w^{(t)}) . \quad (16)$$

Note that samples from the posterior distribution are drawn during the learning phase and predictions for new data can be calculated quickly using the same samples and equation (16).

In the MCMC, samples are generated using a Markov chain that has the desired posterior distribution as its stationary distribution. Difficult part is to create Markov chain which converges rapidly and in which states visited after convergence are not highly dependent.

Neal has used the hybrid Monte Carlo (HMC) algorithm [11] for parameters and Gibbs sampling for hyperparameters. HMC is an elaborate Monte Carlo method which makes efficient use of gradient information to reduce random walk behavior. The gradient indicates in which direction one should go to find states with high probability. Use of Gibbs sampling for hyperparameters helps to minimize the amount of tuning that is needed to obtain good performance in HMC.

When the amount of data increases, the evidence from data causes the probability mass to concentrate to the smaller area and we need less samples from the posterior distribution. Also less samples are needed to evaluate the mean of the predictive distribution than the tail-quantiles like, 10% and 90% quantiles. So depending on the problem 10–20 samples may be enough (given that samples are not too highly dependent).

In our examples (sections 4, 5) of Bayesian learning for neural networks with MCMC we have used Flexible Bayesian Modeling (FBM) software². The methods implemented in software are described in [4].

4 Case I: Inverse Problem in Electrical Impedance Tomography

In this section we report results on using Bayesian neural networks for solving the ill-posed inverse problem in electrical impedance tomography, EIT. The full report of the proposed approach is presented in [12]. Here we review the approach shortly and report comparison results that show that the Bayesian neural networks perform consistently better than other types of networks.

The aim in EIT is to recover the internal structure of an object from surface measurements. Number of electrodes are attached to the surface of the object and current patterns are injected from through the electrodes and the resulting potentials are measured. The inverse problem in EIT, estimating the conductivity distribution from the surface potentials, is known to be severely ill-posed, thus some regularization methods must be used to obtain feasible results [13].

²<URL:<http://www.cs.toronto.edu/~radford/fbm.software.html>>

Figure 1 shows a simulated example of the EIT problem. The volume bounded by the circles in the image represent gas bubble floating in liquid. The conductance of the gas is much lower than that of the liquid, producing the equipotential curves shown in the figure. The simulation was computed with FEM (Finite Element Method) using Matlab PDE-toolbox. Figure 2 shows

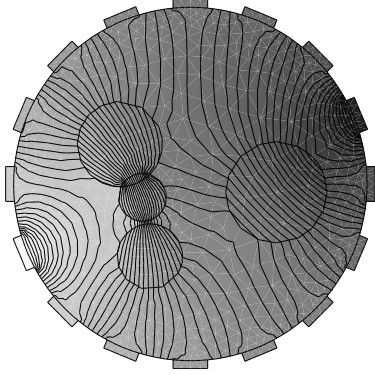


Figure 1: Example of the EIT measurement. The simulated bubble formation is bounded by the circles. The current is injected from the electrode with the lightest color and the opposite electrode is grounded. The gray level and the contour curves show the resulting potential field.

the resulting potential signals, from which the image is to be recovered.

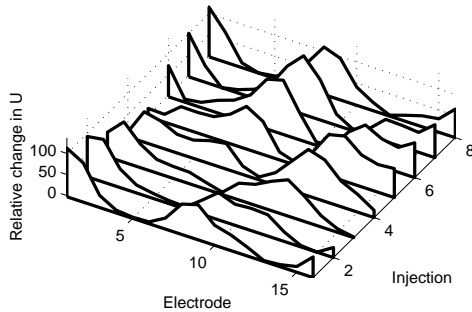


Figure 2: Relative changes in potentials compared to homogenous background. The eight curves correspond to injections from eight different electrodes.

In [12] we proposed a novel feedforward solution for the reconstruction problem. The approach is based on computing the principal component decomposition for the potential signals and the eigenimages of the bubble distribution from the autocorrelation model of the bubbles. The input to the neural network is the projection of the potential signals to the first principal components,

and the network gives the coefficients for reconstructing the image as weighted sum of the eigenimages.

The projection of the potentials and the images to the eigenspace reduces correlations from the input and the output data of the network and detaches the actual inverse problem from the representation of the potential signals and image data. For example, the resolution of the reconstructed images can be changed afterwards, independently of the inverse computation, by recomputing the eigenimages from the autocorrelation model with desired accuracy.

The reconstruction was based on 20 principal components of the 128 dimensional potential signal and 30 eigenimages with resolution 41×41 pixels. The training data consisted of 500 simulated bubble formations with one to ten overlapping circular bubbles in each image. To compute the reconstructions MLP networks containing 30 hidden units (20-30-30 network) with total of about 1500 parameters were used. MLP models tested were

MLP-ESC (NNTB3 defaults) : Early stopping committee of 20 MLP networks, with different division of data to training and stopping sets for each member. The networks were initialized with the Matlab Neural Network Toolbox 3.0 default procedure (Nguyen-Widrow algorithm).

MLP-ESC (decent defaults) : Similar committee to the previous, but the networks were initialized to near zero weights to guarantee that the mapping is smooth in the beginning.

MLP-ESC (mlp-bgd-1) : Early stopping committee used in [14] for benchmarks.

Bayesian MLP : Bayesian neural network with FBM-software, using vague priors and MCMC-run specifications similar as used in [14]. 20 networks from the posterior distribution of network parameters were used.

Figure 3 shows examples of the bubble images reconstructed with Bayesian MLP. The average number of pixels that were erroneously classified to bubble or background was 3.96 % in the test set of 500 bubble formations. Figure 4 shows the goodness of the image reconstructions with different network models for one example image.

Table 1 shows the quality of the image reconstructions with different network models, measured by error in the void fraction and percentage of erroneous pixels in the segmentation.

An important goal in the studied process tomography application was to estimate the void fraction, which is the proportion of gas and liquid in the image. With the proposed approach such goal variables can be estimated directly without explicit reconstruction of the

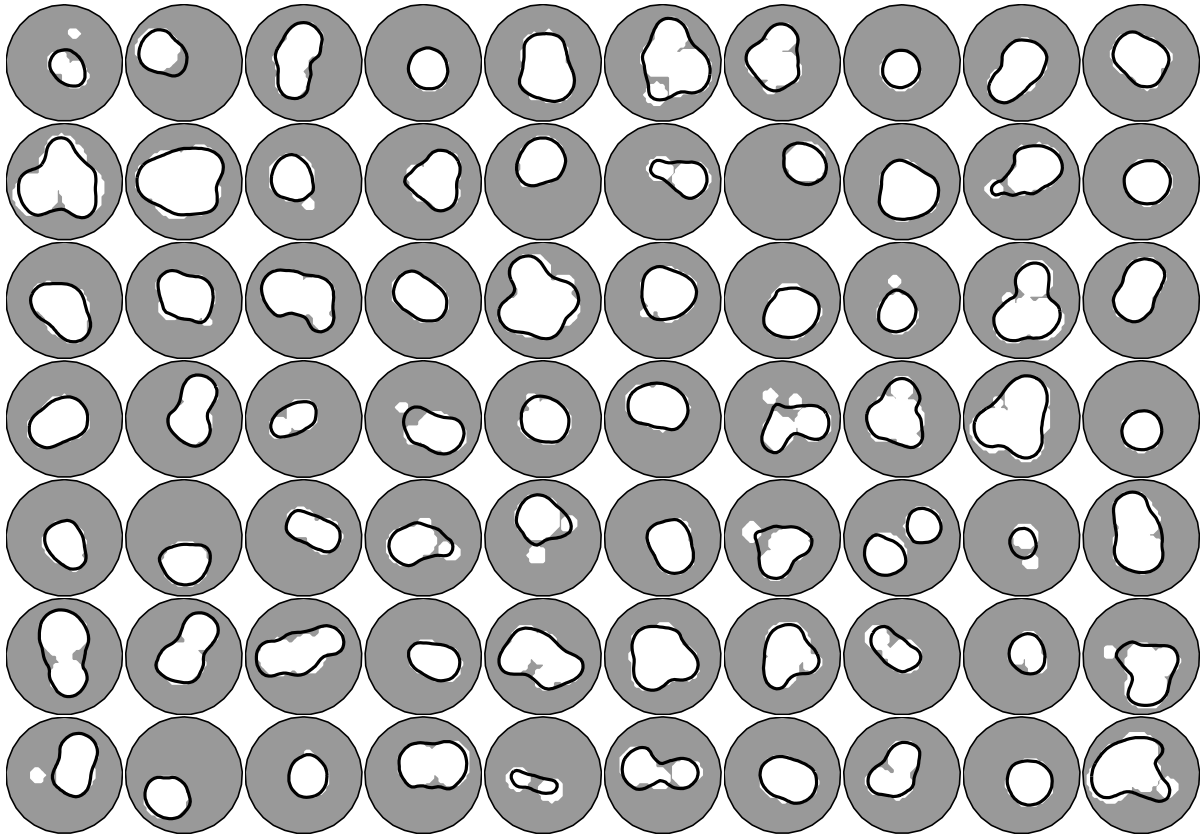


Figure 3: Examples of bubble formations reconstructed with Bayesian MLP. The white blobs show the actual simulated bubbles and the black lines show the contours of the reconstructed bubbles.

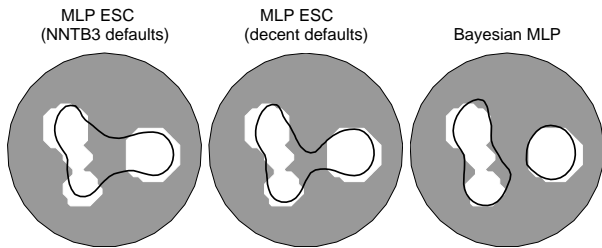


Figure 4: Example of the image reconstruction with Bayesian MLP and early stopping committees. See text for explanation of the models.

Table 1: Errors in reconstructing the bubble shape and estimating the void fraction from the reconstructed images. See text for explanation of the different models.

Method	Classification errors %	Relative error in void fraction %
MLP ESC (NNTB3 def)	4.7	16.2
MLP ESC (decent def)	4.5	15.7
Bayesian MLP	3.8	6.0

image. Table 2 shows the relative absolute errors in estimating the void fraction directly from the projections of the potential signals.

Figure 5 shows the scatter plot of the void fraction versus the estimate by the Bayesian neural network. The 10% and 90% quantiles are computed directly from the posterior distribution of the model output.

See [12] for results for effect of additive Gaussian noise to the performance of the method.

5 Case II: Forest Scene Analysis

In this section we report results of using Bayesian neural networks for classification of forest scenes, to accurately recognize and locate the trees from any background. Potential applications include forest inventory (estimation of the volume and growth rate of the trees) and autonomous forest harvester (navigation and tree manipulation tasks).

Forest scene classification task is demanding due to

Table 2: Relative errors in estimating the void fraction directly. See text for explanation of the different models. Error mean and 90% interval estimated from 4 runs with different random seeds.

Method	Relative test error, %
MLP-ESC (NNTB3 defaults)	8.6 ± 1.2
MLP-ESC (mlp-bgd-1)	6.42 ± 0.04
MLP-ESC (decent defaults)	4.10 ± 0.03
Bayesian MLP	3.16 ± 0.02

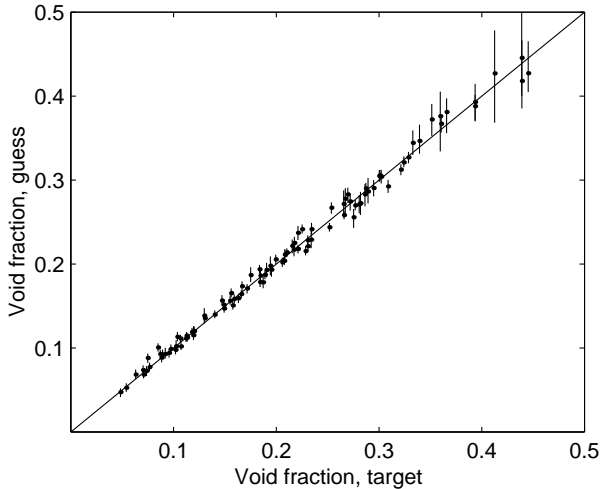


Figure 5: Scatterplot of the void fraction estimate with 10% and 90% quantiles.

the texture richness of the trees, occlusions of the forest scene objects and diverse lighting conditions under operation. This makes it difficult to determine which are optimal image features for the classification. A natural way to proceed is to extract many different types of potentially suitable features.

In [15] we extracted total of 84 statistical and Gabor features over different sized windows at each spectral channel. Due to great number of features used, many classifier methods would suffer from the curse of dimensionality, but Bayesian neural networks manage well in high dimensional problems.

The image data for teaching and testing of the classifiers was collected by using an ordinary digital camera in varying weather conditions. Ideal weather conditions were not searched, as the aim was to test the viability and the robustness of the methods. Total of 48 images were taken.

Based on the above image data a suitable dataset was prepared for the classification study. The labeling of the image data was done by hand via identifying many types of tree and background image blocks with different tex-

tures and lighting conditions. In this study only pines were considered.

To estimate classification errors of different methods we used eight folded cross-validation (CV) error estimate, i.e., 42 of 48 pictures were used for training and the six left out for error evaluation, and this scheme was repeated eight times. The models tested were

KNN LOOCV : K-nearest-neighbor, where K is chosen by leave-one-out cross-validation on the training set.

CART : Classification And Regression Tree [16].

MLP ESC : MLP early stopping committee with different division of data to training and stopping sets for each member of committee.

Bayesian MLP : Bayesian neural network with FBM-software, using vague priors and MCMC-run specifications similar as used in [14].

Bayesian MLP +ARD : Bayesian neural network with FBM-software, using vague priors, Automatic Relevance Determination prior and MCMC-run specifications similar as used in [14].

MLP models contained 20 hidden units (84-20-1 network) with total of about 1700 parameters and e.g. Bayesian MLP with ARD had in addition total of 88 hyperparameters.

We also tested Principal Component Analysis (PCA) for dimension reduction. With PCA we selected first components describing 99% of variance in training data, which were first 16 to 20 principal components depending on training set. With PCA feature MLP models had total of about 400 parameters.

CV error estimates are collected in table 3. Figure 6 shows example image classified with different methods.

Table 3: CV error estimates for forest scene classification. See text for explanation of the different models.

	Error %, all 84 features	Error %, 16-20 pca features
KNN LOOCV	20	24
CART	30	30
MLP ESC	13	19
Bayesian MLP	12	19
Bayesian MLP +ARD	11	19

6 Summary discussion

Above case problems in image analysis illustrate the advantages of using Bayesian neural networks. The approach contains automatic complexity control as the

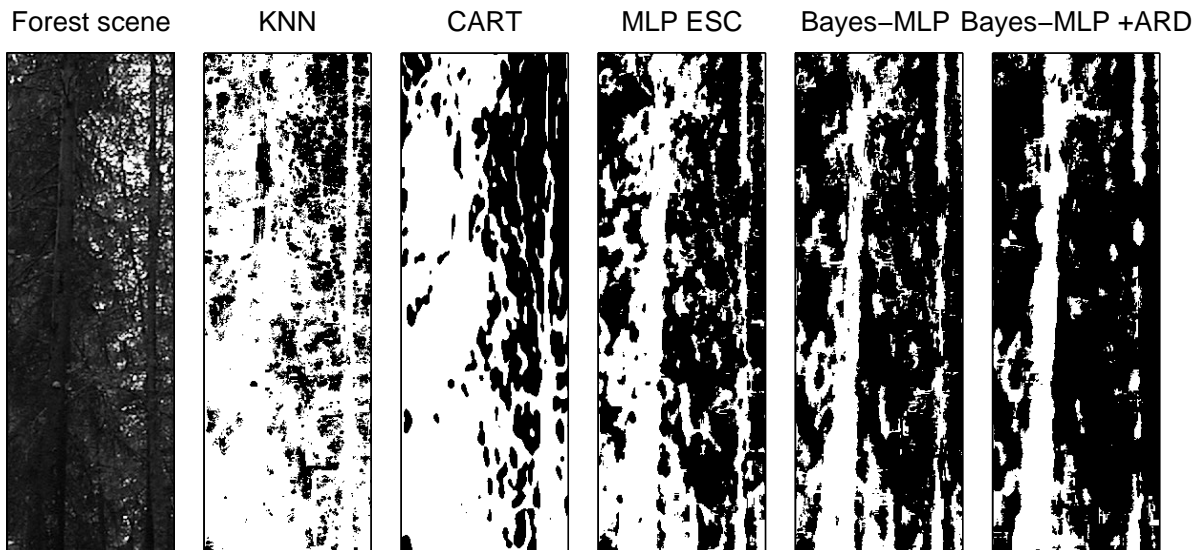


Figure 6: Examples of classified forest scene. See text for explanation of the different models.

Bayesian inference techniques allow the values of regularization coefficients to be selected using only the training data, without the need to use separate training and validation data. As we don't need to fear overfitting, we can use large number of inputs and there is no need to search for minimal set of sufficient inputs. It is possible to use prior information, like ARD. The Bayesian approach gives the predictive distributions for outputs, which can be used to estimate reliability of the predictions.

Acknowledgments

This study was partly funded by TEKES Grant 40888/97 (Project *PROMISE*, *Applications of Probabilistic Modeling and Search*).

References

- [1] S. Geman and D. Geman, "Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images," *IEEE Transactions on Pattern Analysis and Machine Intelligence* **6**, pp. 721–741, 1984.
- [2] D. C. Knull and W. Richards, eds., *Perception as Bayesian Inference*, Cambridge University Press, 1996.
- [3] D. J. C. MacKay, "A practical Bayesian framework for backpropagation networks," *Neural Computation* **4**(3), pp. 448–472, 1992.
- [4] R. M. Neal, *Bayesian Learning for Neural Networks*, vol. 118 of *Lecture Notes in Statistics*, Springer-Verlag, July 1996.
- [5] C. M. Bishop, *Neural Networks for Pattern Recognition*, Oxford University Press, 1995.
- [6] A. Gelman, J. B. Carlin, H. S. Stern, and D. R. Rubin, *Bayesian Data Analysis*, Texts in Statistical Science, Chapman & Hall, 1995.
- [7] C. M. Bishop and C. S. Qazaz, "Regression with input-dependent noise: A Bayesian treatment," in *Advances in Neural Information Processing Systems*, M. C. Mozer, M. I. Jordan, and T. Petsche, eds., vol. 9, MIT Press, (Cambridge, MA), 1997.
- [8] D. J. C. MacKay, "Bayesian non-linear modelling for the prediction competition," in *ASHRAE Transactions, V.100, Pt.2*, pp. 1053–1062, ASHRAE, (Atlanta Georgia), 1994.
- [9] D. J. C. MacKay, "The evidence framework applied to classification networks," *Neural Computation* **4**(5), pp. 720–736, 1992.
- [10] W. Gilks, S. Richardson, and D. Spiegelhalter, eds., *Markov Chain Monte Carlo in Practice*, Chapman & Hall, 1996.
- [11] S. Duane, A. D. Kennedy, B. J. Pendleton, and D. Roweth, "Hybrid Monte Carlo," *Physics Letters B* **195**, pp. 216–222, 1987.
- [12] J. Lampinen, A. Vehtari, and K. Leinonen, "Using Bayesian neural network to solve the inverse problem in electrical impedance tomography," in *Proceedings of 11th Scandinavian Conference on Image Analysis SCIA '99*, (Kangerlussuaq, Greenland), June 1999.
- [13] M. Vauhkonen, J. P. Kaipio, E. Somersalo, and P. A. Karjalainen, "Electrical impedance tomography with basis constraints," *Inverse Problems* **13**(2), pp. 523–530, 1997.
- [14] R. M. Neal, "Assessing relevance determination methods using DELVE," in *Neural Networks and Machine Learning*, C. M. Bishop, ed., vol. 168 of *NATO ASI Series F: Computer and Systems Sciences*, Springer-Verlag, 1998.
- [15] A. Vehtari, J. Heikkonen, J. Lampinen, and J. Juujärvi, "Using Bayesian neural networks to classify forest scenes," in *Intelligent Robots and Computer Vision XVII: Algorithms, Techniques, and Active Vision*, D. P. Casasent, ed., vol. 3522 of *Proceedings of SPIE*, pp. 66–73, (Boston, MA, USA), November 1998.
- [16] L. Breiman, J. Friedman, R. Olshen, and C. Stone, *Classification and regression trees*, Chapman and Hall, 1984.