

Bayesian Training of Neural Networks Using Genetic Programming

Tsili Marwala *Member, IEEE*

Abstract—Bayesian neural networks trained using Markov chain Monte Carlo (MCMC) and genetic programming in binary space within Metropolis framework is proposed. It is tested and compared to classical MCMC method and is observed to give better results than classical approach.

I INTRODUCTION

THE use of Bayesian framework to train neural networks has been the subject of research during the previous decade. Some of the techniques that have been applied so far to train neural networks using Bayesian framework include Markov chain Monte Carlo (MCMC) method [1] and the hybrid Monte Carlo method [2]. Both these methods have been applied within the framework of Metropolis et al. algorithm [3]. Markov chain Monte Carlo method has been applied to improve the abilities of mathematical models to predict the dynamics and reliability of structures [4]. MCMC was used for Bayesian curve fitting and applied to signal segmentation [5] to estimate regularization parameters for satellite image restoration [6] and for inference of stochastic volatility models of the S-P index [7].

All these applications that have been described above have one aspect in common and that is that they have been applied without paying particular attention to the issue of achieving a global optimum posterior distribution function. Kenall and Montana [8] have noted that inside every Markov chain with measurable transition density there is a discrete state space Markov chain struggling to escape from some local optimum distribution. This is in essence indicates that the issue of global posterior distribution must not be taken for granted. Several techniques have been implemented to achieve global optimum distributions such as simulated annealing [9] and evolutionary computing to identify global solutions in optimization problems [10]. This paper therefore proposes the use of genetic programming to sample a posterior probability distribution of the neural network weights. The procedure proposed operates in binary space and conventional evolutionary concepts of mutation, crossover and reproduction are applied. This procedure is then compared to the classical Markov chain Monte Carlo method that samples states in floating point space.

II NEURAL NETWORKS

This section describes multi layer perceptron (MLP) neural networks which are parameterised graphs that make probabilistic assumptions about data. The network architecture implemented in this paper contains i input units, n hidden units and o output units. The relationship between output y and input x may be written as follows [11]:

$$y_k = f_{outer} \left(\sum_{j=1}^M w_{kj}^{(2)} f_{inner} \left(\sum_{i=1}^d w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right) \quad (1)$$

Here $w_{ji}^{(1)}$ and $w_{ji}^{(2)}$ indicate weights in the first and second layers respectively going from input i to hidden unit j while $w_{j0}^{(1)}$ indicates the bias for the hidden unit j . Here M is the number of hidden units, d is the number of input units and k is the index for the output units. In this paper the function $f_{outer}(\bullet)$ is linear while $f_{inner}(\bullet)$ is a hyperbolic tangent function [12]. The Bayesian method identifies the distribution of weights in (1) that look probable in the light of data.

III BAYESIAN FRAMEWORK

In this section a method of identifying the network weights described in (1) is outlined. The problem of identifying the network weights is posed in Bayesian form as follows [13]:

$$P(w|D) = \frac{P(D|w)P(w)}{P(D)} \quad (2)$$

where $P(w)$ is the probability distribution function of the weight space in the absence of any data also known as the prior distribution function and $D \equiv (y_1, \dots, y_N)$ is a matrix containing the output data. The quantity $P(w|D)$ is the posterior probability distribution function after the data have been seen. $P(D|w)$ is the likelihood distribution function and $P(D)$ is the evidence and its function is to normalize the posterior probability distribution function. Equation (2) may be expanded to give [14]:

$$P(w|D) = \frac{1}{Z_s} \exp \left(-\beta \sum_{n=1}^N \sum_{k=1}^K \{t_{nk} - y_{nk}\} - \alpha \sum_{j=1}^M w_j \right) \quad (3)$$

where

Manuscript received December 5, 2005.
 T. Marwala is a professor of Electrical Engineering at the School of Electrical and Information Engineering, University of the Witwatersrand, Johannesburg, South Africa. phone: 7777777 fax: 7777777
 e-mail: t.marwala@ee.wits.ac.za

$$Z_s(\alpha, \beta) = \int \exp \left(-\beta \sum_n \sum_k \{t_{nk} - y_{nk}\}^2 - \frac{\alpha}{2} \sum_j w_j^2 \right) dw \quad (4)$$

$$= \left(\frac{2\pi}{\beta} \right)^{N/2} + \left(\frac{2\pi}{\alpha} \right)^{W/2}$$

In (4) the first term in the exponent is the likelihood function and the second term is the prior information n is the index for the training pattern β is the data contribution to the error k is the index for the output units and α is the coefficient of the prior information. The second term in (4) is the regularization parameter. Training the network using Bayesian approach automatically penalizes overly complex models and also gives a probability distribution of the output of the networks.

IV MCMC VIA METROPOLIS ALGORITHM

The application of Bayesian approach to neural networks results in the probability distribution functions of the network outputs. From these distribution functions the average prediction of the neural network and the variance of that prediction can be calculated. The probability distributions of these network weights are mathematically described by (4). From (4) and by following the rules of probability theory the distribution of the output parameter y , is written as [10]:

$$p(y|x, D) = \int p(y|x, w) p(w|D) dw \quad (5)$$

Equation (5) opens on (4) and is difficult to solve analytically due to relatively high dimension of weight space. Thus the integral in (5) may be approximated as follows:

$$\tilde{y} \cong \frac{1}{L} \sum_{i=1}^{R+L-1} F(w_i) \quad (6)$$

Here F is the mathematical model that gives the output given the input \tilde{y} is the average prediction of the Bayesian neural network R is the number of initial states that are discarded in the hope of reaching a stationary posterior distribution function described in (4) and L is the number of retained states. In this paper MCMC method is implemented by sampling a stochastic process consisting of random variables $\{w_1, w_2, \dots, w_n\}$ through introducing random changes to weight vector $\{w\}$ and either accepting or rejecting the state according to Metropolis et al algorithm given the differences in posterior probabilities between two states that are in transition [11]. The algorithm ensures that states with high probability form the majority of the Markov chain. Traditionally the MCMC was conducted in floating point space and this paper introduces genetic sampling of Bayesian networks which is the subject of the next section.

V MCMC: GENETIC PROGRAMMING AND METROPOLIS ALGORITHM

Genetic programming takes features from natural evolution and uses these to computationally solve practical problems. Genetic algorithms are examples of genetic programming and a procedure that is inspired by these is introduced in this section. In this paper some of the features of genetic computing are applied to sample the posterior distribution function in (5).

Genetic algorithms were inspired by Darwin's theory of natural evolution. In natural evolution members of the population compete with each other to survive and reproduce. Evolutionary successful individuals reproduce while weaker members die. As a result the genes that are successful are likely going to spread within the population. This natural optimisation method has been successfully used to optimise complex problems [12] [14]. This procedure uses a population of binary string chromosomes and each of these strings is the discretised representation of a point in the search space and therefore as a fitness function that is given by the objective function. On generating a new population three operators are performed: (i) crossover (ii) mutation (iii) and reproduction and these operators are a option in genetic MCMC sampling. The crossover operator mixes genetic information in the population by cutting pairs of chromosomes at random points along their length and exchanging over the cut sections. This is as a potential of joining successful operators together. Crossover occurs with a certain probability. In many natural systems the probability of crossover occurring is higher than the probability of mutation occurring. Simple crossover technique is used in this paper [14]. For simple crossover one crossover point is selected binary string from beginning of chromosome to the crossover point is copied from one parent and the rest is copied from the second parent. For example when **11001** undergoes simple crossover with **111** it becomes **11001111**.

The mutation operator picks a binary digit of the chromosomes at random and inverts it. This is as a potential of introducing to the population new information. Mutation occurs with a certain probability. In many natural systems the probability of mutation is low (i.e. less than 1%) In this paper binary mutation is used [14]. When binary mutation is used a number written in binary form is chosen and its value is inverted. For an example: **1** may become **0**.

Reproduction takes successful chromosomes and reproduces them in accordance to their fitness functions. In this paper Metropolis [11] criteria is used as a reproduction method. By so doing the least fit members are therefore gradually driven out of the population of states that form a Markov chain. The schematic illustration of the MCMC method trained using genetic programming is shown in Fig. In this figure an initial sample weight vector $\{w\}_n$ is generated.

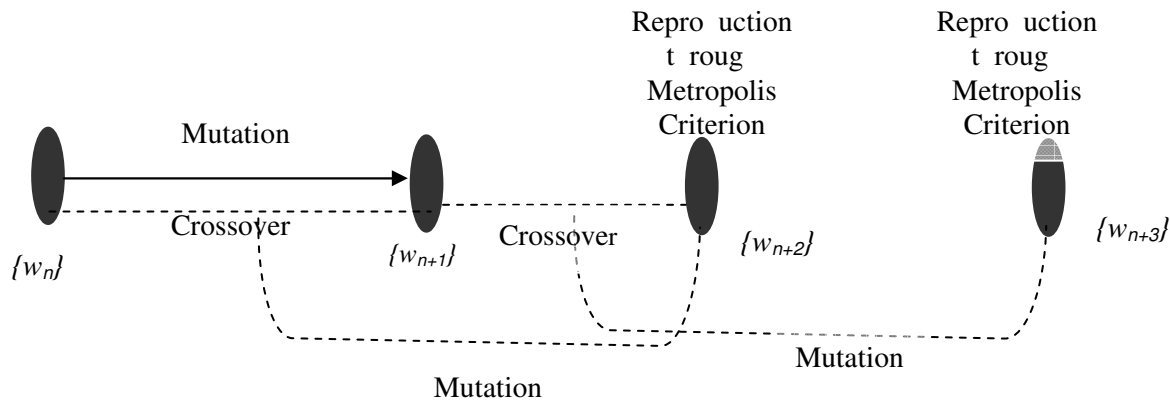


Fig. A schematic illustration of genetic sampling for the MCMC implementation

The sample is converted into binary form using Gray method [10]. The sample is then mutated to form a new sample vector $\{w\}_{n+1}$. The new weight vector $\{w\}_{n+1}$ undergoes crossover with its predecessor $\{w\}_n$ and mutates again to form a new network weight vector $\{w\}_{n+2}$. The weight vector $\{w\}_{n+2}$ is converted into floating point and its probability is calculated. This network weight vector is either accepted or rejected using Metropolis criterion [11]. Thereafter states $\{w\}_{n+2}$ and $\{w\}_{n+1}$ in binary form undergo crossover and are mutated to form $\{w\}_{n+3}$. The genetic MCMC proposed in this section is different from the traditional GA in the following nature: a) The genetic MCMC does not generate a new population of genes at any given iteration (i.e. generation in the GA framework) as is the case in GA but it generates one sample at each iteration. b) The fitness function uses Metropolis criterion while in GA this is not the case and c) The genetic MCMC has a higher mutation rate than GA. The genetic MCMC is different from a standard MCMC in the following way: a) The random walk in the classical MCMC is replaced by a procedure inspired by Darwin's theory of evolution which entails cross over, mutation and reproduction and b) It operates in floating point space.

VI. CASE STUDY : SIMULATED STUDY

In this case study Bayesian neural network that is trained using genetic MCMC is used for regression problems. The same regression problem is solved using the classical MCMC which generates states in floating point space and accept or reject the state using Metropolis et al. method. The results of the two methods are then compared. The simulated data are generated from a noisy sine function with a standard deviation of 0.1. This is the same function that was used by Nabney [5]. Twenty data points are generated around $x=0.25$. Regression analysis is conducted for the domain $x \in [0, 1]$. The MLP networks constructed have one input, five hidden units and one output unit. The optimal number of hidden units was obtained by studying the relationship between the number of hidden units and the generalization error. This was conducted by setting the number of hidden

units to fall between 8 and 18 and assessing the generalization error. The hidden layer activation functions are hyperbolic tangent functions while the output activation functions are linear functions.

On implementing Bayesian training the coefficient of the data contribution to the error β is set to 0.1 while the prior coefficient α is set to 0.01. The manner in which these parameters fit into the Bayesian framework is described by [12]. The number of retained states L is 100 while the number of discarded states R is 100. These values fit into the Bayesian framework through [13].

For the genetic part of the simulation the rate of mutation is 1% and the rate of crossover is 7%. It should be noted that the rate of mutation proposed here is higher than that of standard genetic algorithm. The proposed Bayesian method via genetic programming as a random component search and therefore may be viewed as being equivalent to the random walk that is executed in the standard Bayesian sampling. In the proposed procedure may in principle be equivalent to the standard random walk, however it takes into account of the efficient sampling in binary space which has been observed in standard genetic algorithm. It must be noted that the rate of mutation chosen here is lower than the rate of crossover which is in accordance to many natural systems. When implementing the genetic framework through genetic algorithm, bit binary numbers are used. The bounds of the magnitudes of the components of the weight vectors are $[-4, 4]$. The results obtained when Bayesian networks are trained using genetic programming are shown in Fig. 4.

The mean square error (MSE) obtained from this figure is 0.007. The results obtained when the Bayesian networks are trained using classical MCMC are shown in Figure 5.

This gives an MSE of 0.0055. When Fig. 4 is compared to Fig. 5, it is observed that genetic approach to MCMC performs better than the method that uses MCMC method that operates in floating point space because it gives lower average errors. The graph showing the errors as a function of samples accepted from state to state is shown in Fig. 4.

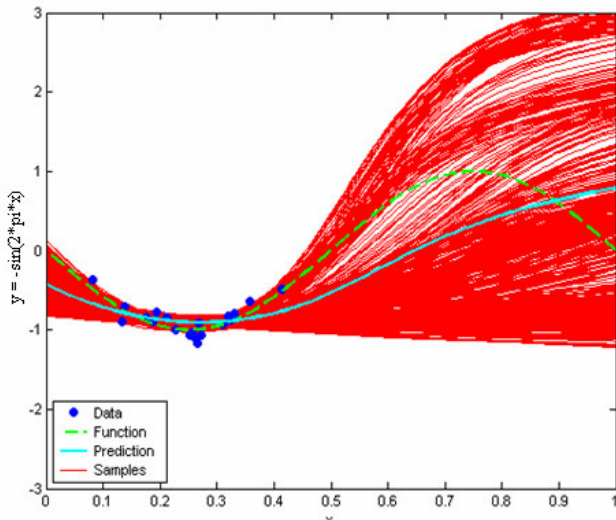


Fig 3 Results obtained when Bayesian networks are trained via genetic programming

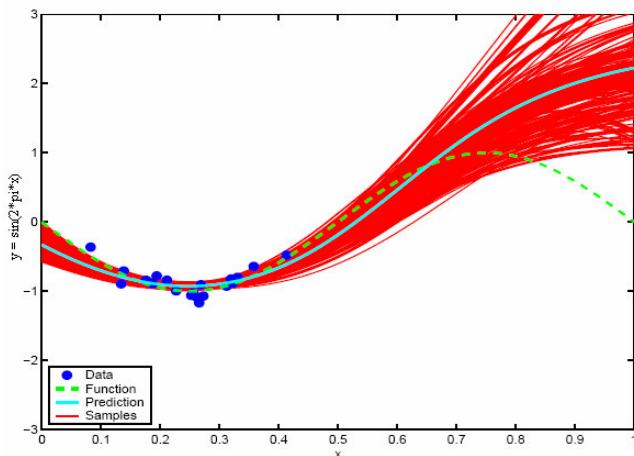


Fig 4 Results obtained when Bayesian networks are trained via traditional MCMC

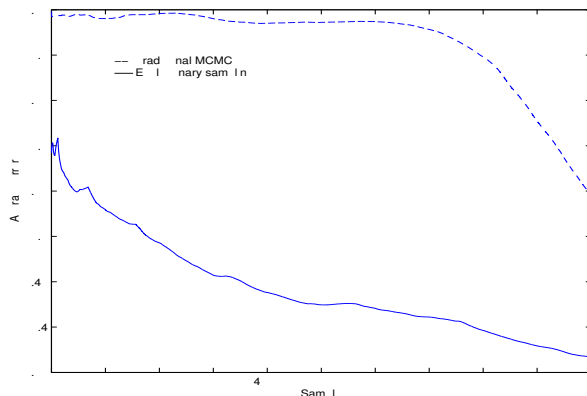


Fig 5 Prediction errors versus samples

This graph shows that the rate of convergence to a stationary posterior distribution is faster when using genetic approach than when using the standard MCMC method. This is because it has been proven that sampling through binary space is more efficient than sampling through floating point space. The reason for this is that sampling through binary space is able to explore a larger part of the weight space than if the process is conducted in floating point space. The acceptance rate of states for MCMC method that operates in floating point space is 8% whereas when using MCMC method based on genetic programming is 7%. This indicates that the MCMC method based on genetic programming is able to explore states that form a Markov chain better than the MCMC method that operates in floating point space.

VII CASE STUDY : ARTIFICIAL TASTER

Now that the simulations have been conducted we now apply the procedure to a practical problem of development of an artificial taster. Artificial taster has been the subject of research for some time. Some of the works on this subject are development of a taster based on proton transfer mass spectroscopy to successfully taste mozzarella cheese [6] a solid state electronic taster for beverage analysis [7] a taste sensor based on lipid coated crystal microbalance to evaluate beer body and smoothness [8] and wine flavour taster that uses multivariate statistics [9]. In this paper the method proposed is used to construct an artificial taster and compare to the classical approach. The artificial taster is basically an infrastructure that relates the characteristics of beer measured in the laboratory to taste score measured from a panel of professional tasters. These characteristics capture parameters that human beings are sensitive to on tasting beer and these are: alcohol level, sugar level, calorie present, extract, real extract, pH, iron, acetaldehyde, isomerization, ethyl acetate, isoamyl acetate, total igt, alcohol colour and bitterness as inputs to the artificial taster that estimates the average taste score given by a panel of professional tasters. To realise the artificial taster, MLP neural network is used. The MLP has inputs 7, hidden nodes and output 1. The genetic parameters used in the previous example are used in this case. The characteristics of beer and the corresponding taste score from an average score from a panel of professional tasters are used to construct an artificial taster. These tasters each give a taste score that ranges from 1 for a really bad beer to 10 for a good beer. When an artificial taster which is defined in this paper as

Table 1 Samples needed for convergence

Method	Number of Samples for Convergence	Average Percentage Errors
Classical Approach	5	9 %
Evolutionary	4	7 %

neural networks and measure analytical data from beer is used to predict taste scores the results obtained are shown in Table and Fig 5. The idea of giving a taste a numerical number is quite difficult to justify philosophically but psychophysically as provide appropriate measurement techniques for subjective phenomena such as taste and this is the framework that is adopted in this paper.

These results show that the genetic approach proposed in this paper converges faster than the classical MCMC. Secondly these results show that the genetic programming gives lower average errors than the classical approach. This

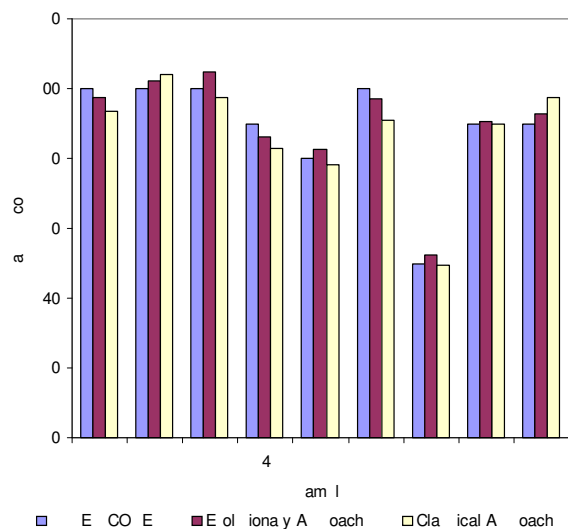


Fig 5 The taste score for a given taste sample. The top graph is for the entire 7 test data while the bottom graph is a close view of the top graph.

is because of the fact that genetic programming is able to explore a wider search space more efficiently than the classical approach.

VIII CONCLUSION

The method is tested on simulated data. The results obtained are compared to those obtained from MCMC method that operates in floating point space. It is concluded that the MCMC method based on genetic programming gives better results than MCMC method that operates in floating point space on modeling simulated data and artificial taster.

REFERENCES

[1] R. E. Kass, B. P. Carlin, A. Gelman, and M. Neal, "Markov Monte Carlo in practice: A round table discussion," *American Statistician*, vol. 5, pp. 9–19, 1998.

[2] R. M. Neal, "An improved acceptance procedure for the hybrid Monte Carlo algorithm," *J. Computational Physics*, vol. 175, pp. 969–984, 2002.

[3] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, "Equations of state calculations by fast computing machines," *Journal of Chemical Physics*, vol. 21, pp. 1087–1092, 1953.

[4] T. Marwala and S. Sibisi, "Finite element updating using Bayesian framework and modal properties," *Journal of Aircraft*, vol. 41, pp. 75–78, 2004.

[5] E. Punskey, C. An, R. A. Doucet, and W. J. Fitzgerald, "Bayesian curve fitting using MCMC with applications to signal segmentation," *IEEE Transactions on Signal Processing*, vol. 53, pp. 747–758, 2005.

[6] A. Jalobeanu, L. Blanc-Feraud, and J. Zerubia, "Hyperparameter estimation for satellite image restoration using a MCMC maximum likelihood method," *Pattern Recognition*, vol. 38, pp. 4–15, 2005.

[7] S. C. ib, F. Nar, and N. S. ep, "Markov chain Monte Carlo methods for stochastic volatility models," *Journal of Econometrics*, vol. 118, pp. 8–28, 2004.

[8] W. S. Ken, and G. Montana, "Small sets and Markov transition densities," *Stochastic Processes and their Applications*, vol. 99, pp. 77–94, 2002.

[9] R. M. Neal, "Probabilistic inference using Markov chain Monte Carlo methods," *University of Toronto Technical Report CRG-TR-93-1*, Toronto, Canada, 1993.

[10] T. Marwala, "Finite element model updating using wavelet data and genetic algorithm," *Journal of Aircraft*, vol. 39, pp. 7–9, 2002.

[11] C. M. Bishop, *Neural Networks for Pattern Recognition*. Oxford University Press, Oxford, UK, 1995.

[12] J. Holland, *Adaptation in natural and artificial systems*. University of Michigan Press, 1975.

[13] Z. Michalewicz, *Genetic algorithms + data structures = evolution programs*. Springer Verlag, 1999.

[14] D. E. Goldberg, "Genetic algorithms in search, optimization and machine learning," *Addison-Wesley Reading MA*, 1989.

[15] I. Nabney, *Netlab: algorithms for pattern recognition*. Springer, Berlin, 1994.

[16] F. Gasperi, "The mozzarella cheese flavour profile: a comparison between juice panel and proton transfer reaction mass spectroscopy," *Journal of the Science of Food and Agriculture*, vol. 8, pp. 57–64, 2004.

[17] L. Lvova, S. S. Kim, S. A. Legin, Y. Vlasov, J. S. Yang, G. S. C. and H. Nam, "All solid state electronic tongue and its application for beverage analysis," *Analytica Chimica Acta*, vol. 48, pp. 4–14, 2004.

[18] Y. Vlasov, A. Legin, and A. R. nitskaya, "Electronic tongues and their analytical application," *Analytical and Biological Chemistry*, vol. 7, pp. 4–14, 2004.

[19] A. C. Noble, and S. E. Ebeler, "Use of multivariate statistics in understanding wine flavor," *Food Reviews International*, vol. 8, pp. 1–14, 1992.