

Bootstrap Model Selection Criterion for Determining the Number of Hidden Units in Neural Network Model¹⁾

Changha Hwang and Daehak Kim²⁾

Abstract

Statistical relations between a system and empirical distribution are studied in terms of the concept of Akaike's Information Criterion. From this consideration we derive a bootstrap criterion for determining the optimal number of hidden units in neural networks.

1. Introduction

In data analysis and AI research, one of the most important applications of neural networks is modeling a system whose input and output relationship is unknown. In general, weights of the network are modified by means of a stochastic gradient decent method which minimizes a certain error function composed of input and output samples observed from the actual system.

In neural networks, an important but difficult problem is to determine the optimal number of hidden units. This problem is in principle no different from selecting regressors in a linear regression or the order of a polynomial regression. The main idea that have been developed in this area is minimizing some measure of performance over the class of possible models. An increase in the number of the parameters will lessen the output errors for given training samples but will raise additional errors for testing samples. All the candidate measures aim to predict the performance on a set of testing samples, and so to select the model with smallest generalization error. In this paper, we consider statistical relations between a system and empirical distribution in terms of the concept of Akaike's Information Criterion. From this consideration we derive a bootstrap criterion for determining the optimal number of hidden units in neural networks. Murata *et al.*(1994) proposed network information criterion using the relation between the training error and the generalization error. Liu(1995) provided a model selection criterion using the cross-validation method to estimate the generalization error based on the optimal parameter obtained by minimizing an error function. In Statistics there has

1) This research was supported in part by Catholic University of Taegu-Hyosung.

2) Dept. of Statistical Information, Catholic University of Taegu-Hyosung, Kyungbuk, Korea.

been a substantial amount of work in the problem of model selection. Chung et al.(1996) showed recently that such a bootstrap version of model selection criterion as suggested by Linhart and Zucchini(1986) has a downward bias of amount roughly equivalent to the number of parameters of approximating model.

2. Asymptotic Properties of Learning

Denote the underlying conditional probability distribution as $p(\mathbf{y}|\mathbf{x})$, the marginal probability distribution of \mathbf{x} as $p(\mathbf{x})$, and the joint probability distribution of (\mathbf{x}, \mathbf{y}) as $p(\mathbf{x}, \mathbf{y})$. Hereafter we identify a system with a probability distribution $p(\mathbf{x}, \mathbf{y})$. On the other hand, a neural network model such as a feed-forward neural network or a competitive network is denoted by a function $g(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta})$, where $\boldsymbol{\theta} \in R^m$ is a parameter vector of weights and thresholds which specify the network. It is modified according to a learning rule. We define a distance function which decides how the model $g(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta})$ is closely related to the system $p(\mathbf{x}, \mathbf{y})$.

Definition 1. A distance function or the expected loss $D(p, g(\boldsymbol{\theta}))$ between a system p and a model $g(\boldsymbol{\theta})$ is defined by

$$D(p, g(\boldsymbol{\theta})) = \int d(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}$$

where d is a loss plus a regularization term:

$$d(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) = d(p(\mathbf{x}, \mathbf{y}), g(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta})). \quad \blacksquare$$

Our purpose is to find an optimal number of hidden units and the optimal parameter vector, namely the optimal model based on minimizing the distance function given the training data set $\{(\mathbf{x}_i, \mathbf{y}_i), i = 1, \dots, n\}$. We denote an empirical distribution function as

$$p_n(\mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \delta(\mathbf{x} - \mathbf{x}_i, \mathbf{y} - \mathbf{y}_i),$$

where δ is the delta function. It is well known that p_n approximates adequately p if n is large enough. The optimal parameters for the true system and the empirical distribution are denoted by $\boldsymbol{\theta}_0$ and $\hat{\boldsymbol{\theta}}$, respectively:

$$\boldsymbol{\theta}_0 = \arg \min_{\boldsymbol{\theta}} D(p, g(\boldsymbol{\theta}))$$

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} D(p_n, g(\boldsymbol{\theta})).$$

The learning method to be described is obtained by generalizing the back-propagation method or the learning vector quantization method.

Definition 2. The parameter vector θ_t at time t is modified according to the following stochastic gradient descent method:

$$\theta_{t+1} = \theta_t - \varepsilon_t \nabla d(\mathbf{x}_t, \mathbf{y}_t, \theta_t),$$

where ∇ denotes a differential operator in respect of θ , ε_t is a positive learning coefficient, and t represents the number of modification. Here $(\mathbf{x}_t, \mathbf{y}_t)$ is an example at time t independently chosen from the training set. ■

In general, the learning coefficient ε_t may be changed depending on time t . But, when the optimal θ_0 deviates slightly, training time would be very long for adjusting the deviation because ε_t might have become too small. In this sense, it is desirable to fix ε_t at a positive small constant. After learning an estimated value $\tilde{\theta}$ is obtained:

$$\tilde{\theta} = \lim_{t \rightarrow \infty} \theta_t,$$

and it varies according to the order of samples applied at the learning period. The probability distribution of θ_t and $\tilde{\theta}$ is denoted by $\pi_t(\theta_t)$ and $\tilde{\pi}(\tilde{\theta})$, respectively. The distribution $\pi_t(\theta_t)$ converges to some distribution $\tilde{\pi}(\tilde{\theta})$. Therefore, when t is large enough, the random variable θ_t is subject to the distribution $\tilde{\pi}(\theta_t)$.

Since the learning rule is stochastic, it is not necessarily guaranteed that θ_t converges to the optimal θ_0 , even if we consider only a neighborhood of θ_0 in order to avoid the convergence to a local minimum. Amari(1967, 1993) analyzed the dynamic behavior of θ_t in the neighborhood of θ_0 and obtained the convergence speed to θ_0 and the fluctuation of θ_t around θ_0 . In this paper, we show the distribution $\tilde{\pi}(\tilde{\theta})$ approaches the normal distribution with mean θ_0 . A brief proof is given in Lemma 1.

Lemma 1. The distribution $\tilde{\pi}(\tilde{\theta})$ approaches the normal distribution

$$N(\theta_0, 2\varepsilon L_Q^{-1}G)$$

as $t \rightarrow \infty$. Here, L_Q is a linear operator defined as

$$L_Q M \equiv QM + (QM)^T, \quad M \in R^{m \times m},$$

where $G \equiv V_p[\nabla d(\mathbf{x}, \mathbf{y}, \theta_0)]$ and $Q \equiv E_p[\nabla \nabla d(\mathbf{x}, \mathbf{y}, \theta_0)]$.

Proof. After lots of algebra we get the following equation corresponding to the equation (10) of Lemma 2 in Murata *et al.*(1994).

$$\varphi_{n+1}(z) = \varphi_n(z) - \varepsilon z^t Q \frac{\partial}{\partial z} \varphi_n(z) - \varepsilon^2 z^t G z \varphi_n(z) + O(\varepsilon^2).$$

Following the next steps in the proof of Lemma 2 in Murata *et al.*(1994) gives us the desired result. ■

Murata *et al.*(1994) showed that the distribution $\tilde{\pi}(\tilde{\theta})$ approaches the normal distribution with mean $\hat{\theta}$. We introduce the result in Lemma 2.

Lemma 2. The distribution $\tilde{\pi}(\tilde{\theta})$ approaches the normal distribution

$$N(\hat{\theta}, \varepsilon L^{-1} \hat{G})$$

as $t \rightarrow \infty$. Here, $\hat{G} \equiv V_{p_n}[\nabla d(\mathbf{x}, \mathbf{y}, \hat{\theta})]$ and $\hat{Q} \equiv E_{p_n}[\nabla \nabla d(\mathbf{x}, \mathbf{y}, \hat{\theta})]$. ■

Let $\hat{\pi}(\hat{\theta})$ be the probability distribution of $\hat{\theta}$. It is known in Statistics that the distribution $\hat{\pi}(\hat{\theta})$ has the following property.

Lemma 3. $\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^d N(0, Q^{-1} G Q^{-1})$ as $n \rightarrow \infty$. ■

3. Bootstrap Model Selection Criterion

Since we do not assume that p belongs to the approximating family of models, we use nonparametric bootstrap method. But, naive nonparametric bootstrap may cause some trouble when the sample size is relatively small due to frequent redundancy. First, we choose a bootstrap sample from the empirical cumulative distribution function p_n .

$$(\mathbf{x}_1, \mathbf{y}_1)^*, \dots, (\mathbf{x}_n, \mathbf{y}_n)^* \sim^{iid} p_n.$$

Then, we can compute the bootstrap version of $\tilde{\theta}$ as the minimizer of the bootstrapped empirically discrepancy. It can be known from Chung *et al.*(1996) that bootstrap expected overall discrepancy is given by

$$E^*[D(p_n, g(\tilde{\theta}^*))].$$

In many cases where explicit forms are not easily available, it can be approximated by replacing the basic resampling procedure and taking the average of each bootstrap copies. Formally, with B bootstrap repetitions, it can be approximated as

$$B^{-1} \sum_{b=1}^B D(p_n, g(\tilde{\theta}_b^*)).$$

Proposition 1. The following approximation holds for the naive plug-in bootstrap expected

overall discrepancy:

$$E^*[D(p_n, g(\tilde{\theta}^*))] = D(p_n, g(\tilde{\theta})) + \frac{1}{n} \text{tr} G Q^{-1} + O(n^{-\frac{3}{2}}).$$

Proof. Let us consider the expansion of $D(p_n, g(\tilde{\theta}^*))$ at $\tilde{\theta}$

$$\begin{aligned} D(p_n, g(\tilde{\theta}^*)) &= D(p_n, g(\tilde{\theta})) + (\tilde{\theta}^* - \tilde{\theta})^t \{\nabla D(p_n, g(\tilde{\theta}))\} \\ &\quad + \frac{1}{2} (\tilde{\theta}^* - \tilde{\theta})^t \{\nabla \nabla D(p_n, g(\tilde{\theta}))\} (\tilde{\theta}^* - \tilde{\theta}) + \text{h.o.t.} \end{aligned}$$

The second term on the right hand side of vanishes, since $\tilde{\theta}$ is chosen that way. Therefore, the expected value of the both sides of the above equation, with respect to bootstrap distribution, turns out to be

$$E^*[D(p_n, g(\tilde{\theta}^*))] = D(p_n, g(\tilde{\theta})) + \frac{1}{2} E^*[(\tilde{\theta}^* - \tilde{\theta})^t \{\nabla \nabla D(p_n, g(\tilde{\theta}))\} (\tilde{\theta}^* - \tilde{\theta})].$$

The second term can be rewritten as

$$\begin{aligned} \frac{1}{2} E^*[(\tilde{\theta}^* - \tilde{\theta})^t \{\nabla \nabla D(p_n, g(\tilde{\theta}))\} (\tilde{\theta}^* - \tilde{\theta})] &\approx \frac{1}{2} E[(\tilde{\theta} - \theta_0)^t \{\nabla \nabla D(p_n, g(\tilde{\theta}))\} (\tilde{\theta} - \theta_0)] \\ &\approx \frac{1}{2} E_{\tilde{\pi}, \hat{\pi}}[(\tilde{\theta} - \theta_0)^t \{\nabla \nabla D(p_n, g(\tilde{\theta}))\} (\tilde{\theta} - \theta_0)]. \end{aligned}$$

The integrand on the right hand side can be rewritten as

$$\begin{aligned} &\frac{1}{2} (\tilde{\theta} - \theta_0)^t \{\nabla \nabla D(p_n, g(\tilde{\theta}))\} (\tilde{\theta} - \theta_0) \\ &= \frac{1}{2} \text{tr} [\hat{Q} \{ (\tilde{\theta} - \hat{\theta}) + (\hat{\theta} - \theta_0) \} \{ (\tilde{\theta} - \hat{\theta}) + (\hat{\theta} - \theta_0) \}^t] \\ &= \frac{1}{2} \text{tr} [\hat{Q} \{ (\tilde{\theta} - \hat{\theta})(\tilde{\theta} - \hat{\theta})^t + (\hat{\theta} - \theta_0)(\hat{\theta} - \theta_0)^t \\ &\quad + (\tilde{\theta} - \hat{\theta})(\hat{\theta} - \theta_0)^t + (\hat{\theta} - \theta_0)(\tilde{\theta} - \hat{\theta})^t \}]. \end{aligned}$$

Averaging each term subject to $\tilde{\pi}(\tilde{\theta})$ and $\hat{\pi}(\hat{\theta})$ and using

$$E_{\tilde{\pi}, \hat{\pi}}[(\tilde{\theta} - \hat{\theta})(\tilde{\theta} - \hat{\theta})^t] = \frac{\varepsilon}{2} \hat{Q}^{-1} \hat{G},$$

$$E_{\tilde{\pi}, \hat{\pi}}[(\hat{\theta} - \theta_0)(\hat{\theta} - \theta_0)^t] = \frac{1}{n} Q^{-1} G Q^{-1},$$

$$E_{\tilde{\pi}, \hat{\pi}}[(\tilde{\theta} - \hat{\theta})(\hat{\theta} - \theta_0)^t] = 0$$

we get the required result. ■

Therefore, the bootstrap model selection criterion in neural network is given by

$$D(p_n, g(\tilde{\theta})) + \frac{1}{n} \text{tr} \tilde{G} \tilde{Q}^{-1} .$$

We choose the model minimizes this criterion, which turns out to be equivalent to NIC given by Murata *et al.*(1994).

References

- [1] Amari, S. (1967). Theory of adaptive pattern classifiers, *IEEE Transactions on Electronic Computers*, EC-16, 299-307.
- [2] Amari, S. (1993). Backpropagation and stochastic gradient descent method, *Neurocomputing*, 5, 185-196.
- [3] Chung, H., Lee, K. and Koo, J. (1996). A note on bootstrap model selection criterion, *Statistics and Probability Letters*, 26, 35-41.
- [4] Linhart, H. and Zucchini, W.(1986). Model Selection, Wiley, New York.
- [5] Liu, Y. (1995). Unbiased estimate of generalization error and model selection in neural network, *Neural Networks*, 8, 215-219.
- [6] Murata, N., Yoshizawa, S. and Amari, S. (1991). A criterion for determining the number of parameters in an artificial neural network model, *Artificial Neural Networks*, T. Kohonen, *et al.* ed. Amsterdam: The Netherlands: Elsevier, 9-14.
- [7] Murata, N., Yoshizawa, S. and Amari, S. (1994). Network information criterion - Determining the number of hidden units for an artificial neural network model, *IEEE Transactions on Neural Networks*, 5, 865-872.