

# Adaptive Regularization in Neural Network Modeling

Jan Larsen<sup>1</sup>, Claus Svarer<sup>2</sup>, Lars Nonboe Andersen<sup>1</sup> and Lars Kai Hansen<sup>1</sup>

<sup>1</sup> Department of Mathematical Modelling, Building 321  
Technical University of Denmark  
DK-2800 Lyngby, Denmark  
emails: jl@imm.dtu.dk, lna@imm.dtu.dk, lkhansen@imm.dtu.dk  
www: <http://eivind.imm.dtu.dk>  
<sup>2</sup> Neurobiology Research Unit  
Department of Neurology, Building 9201  
Copenhagen University Hospital  
Blegdamsvej 9  
DK-2100 Copenhagen Ø, Denmark  
email: [csvarer@pet.rh.dk](mailto:csvarer@pet.rh.dk)  
www: <http://neuro.pet.rh.dk>

**Abstract.** In this paper we address the important problem of optimizing regularization parameters in neural network modeling. The suggested optimization scheme is an extended version of the recently presented algorithm [24]. The idea is to minimize an empirical estimate – like the cross-validation estimate – of the generalization error with respect to regularization parameters. This is done by employing a simple iterative gradient descent scheme using virtually no additional programming overhead compared to standard training. Experiments with feed-forward neural network models for time series prediction and classification tasks showed the viability and robustness of the algorithm. Moreover, we provided some simple theoretical examples in order to illustrate the potential and limitations of the proposed regularization framework.

## 1 Introduction

Neural networks are flexible tools for time series processing and pattern recognition. By increasing the number of hidden neurons in a 2-layer architecture any relevant target function can be approximated arbitrarily close [18]. The associated risk of overfitting on noisy data is of major concern in neural network design, which find expression in the ubiquitous bias-variance dilemma, see e.g., [9].

The need for regularization is two-fold: First, it remedies numerical problems in the training process by smoothing the cost function and by introducing curvature in low (possibly zero) curvature regions of cost function. Secondly, regularization is a tool for reducing variance by introducing extra bias. The overall objective of architecture optimization is to minimize the generalization error. The architecture can be optimized *directly* by stepwise selection procedures (including pruning techniques) or *indirectly* using regularization. In general, one

would prefer a hybrid scheme; however, a very flexible regularization may substitute the need for selection procedures. The numerical experiments we consider mainly hybrid pruning/adaptive regularization schemes.

The trick presented in this communication addresses the problem of adapting regularization parameters.

*The trick consists in formulating a simple iterative gradient descent scheme for adapting the regularization parameters aiming at minimizing the generalization error.*

We suggest to use an empirical estimate<sup>3</sup> of the generalization error, viz. the  $K$ -fold cross-validation [8], [38]. In [24] and [3] the proposed scheme was studied using the hold-out validation estimator.

In addition to empirical estimators for the generalization error a number of *algebraic* estimators like FPE [1], FPER [22], GEN [20], GPE [30] and NIC [32] have been developed in recent years. These estimates, however, depend on a number of statistical assumptions which can be quite hard to justify. In particular, they are  $o(1/N_t)$  estimators where  $N_t$  is the number of training examples. However, for many practical modeling set-ups it is hard to meet the large training set assumption.

In [13] properties of adaptive regularization is studied in the simple case of estimating the mean of a random variable using an algebraic estimate of the average<sup>4</sup> generalization error and [14] proposed an adaptive regularization scheme for neural networks based on an algebraic estimate. However, experiments indicate that this scheme has a drawback regarding robustness. In addition, the requirement of a large training set may not be met.

The Bayesian approach to adapt regularization parameters is to minimize the so-called evidence [5, Ch. 10], [29]. The evidence, however, does not in a simple way relate to the generalization error which is our primary object of interest.

Furthermore [2] and [37] consider the use of a validation set to tune the amount of regularization, in particular when using the early-stop technique.

Section 2 considers training and empirical generalization assessment. In Section 3 the framework for optimization of regularization parameters is presented. The experimental section 4 deals with examples of feed-forward neural networks models for classification and time series prediction. Further, in order to study the theoretical potential/limitations of the proposed framework, we include some simulations on a simple toy problem.

## 2 Training and Generalization

Suppose that the neural network is described by the vector function  $\mathbf{f}(\mathbf{x}; \mathbf{w})$  where  $\mathbf{x}$  is the input vector and  $\mathbf{w}$  is the vector of network weights and thresholds with dimensionality  $m$ . The objective is to use the network model for approximating the true conditional input-output distribution  $p(\mathbf{y}|\mathbf{x})$ , or some moments

---

<sup>3</sup> For further discussion on empirical generalization assessment, see e.g., [23].

<sup>4</sup> Average w.r.t. to different training sets.

hereof. For regression and signal processing problems we normally model the conditional expectation  $E\{\mathbf{y}|\mathbf{x}\}$ .

Assume that we have available a data set  $\mathcal{D} = \{\mathbf{x}(k); \mathbf{y}(k)\}_{k=1}^N$  of  $N$  input-output examples. In order to both train and empirically estimate the generalization performance we follow the idea of  $K$ -fold cross-validation [8], [38] and split the data set into  $K$  randomly chosen disjoint sets of approximately equal size, i.e.,  $\mathcal{D} = \cup_{j=1}^K \mathcal{V}_j$  and  $\forall i \neq j : \mathcal{V}_i \cap \mathcal{V}_j = \emptyset$ . Training and validation is replicated  $K$  times, and in the  $j$ 'th run training is done on the set  $\mathcal{T}_j = \mathcal{D} \setminus \mathcal{V}_j$  and validation is performed on  $\mathcal{V}_j$ .

The cost function,  $C_{\mathcal{T}_j}$ , for network training on  $\mathcal{T}_j$ , is supposed to be the sum of a loss function (or training error),  $S_{\mathcal{T}_j}(\mathbf{w})$ , and a regularization term  $R(\mathbf{w}, \boldsymbol{\kappa})$  parameterized by a set of regularization parameters  $\boldsymbol{\kappa}$ , i.e.,

$$C_{\mathcal{T}_j}(\mathbf{w}) = S_{\mathcal{T}_j}(\mathbf{w}) + R(\mathbf{w}, \boldsymbol{\kappa}) = \frac{1}{N_{tj}} \sum_{k=1}^{N_{tj}} \ell(\mathbf{y}(k), \hat{\mathbf{y}}(k); \mathbf{w}) + R(\mathbf{w}, \boldsymbol{\kappa}) \quad (1)$$

where  $\ell(\cdot)$  measures the distance between the output  $\mathbf{y}(k)$  and the network prediction  $\hat{\mathbf{y}}(k) = \mathbf{f}(\mathbf{x}(k); \mathbf{w})$ . In section 4 we will consider log-likelihood and the square error loss function  $\ell = |\mathbf{y} - \hat{\mathbf{y}}|^2$ .  $N_{tj} \equiv |\mathcal{T}_j|$  defines the number of training examples in  $\mathcal{T}_j$  and  $k$  is used to index the  $k$ 'th example  $[\mathbf{x}(k), \mathbf{y}(k)]$ . Training provides the estimated weight vector  $\hat{\mathbf{w}}_j = \arg \min_{\mathbf{w}} C_{\mathcal{T}_j}(\mathbf{w})$ .

The  $j$ 'th validation set  $\mathcal{V}_j$  consist of  $N_{vj} = N - N_{tj}$  examples and the validation error<sup>5</sup> of the trained network reads

$$S_{\mathcal{V}_j}(\hat{\mathbf{w}}_j) = \frac{1}{N_{vj}} \sum_{k=1}^{N_{vj}} \ell(\mathbf{y}(k), \hat{\mathbf{y}}(k); \hat{\mathbf{w}}_j) \quad (2)$$

where the sum runs over the  $N_{vj}$  validation examples.  $S_{\mathcal{V}_j}(\hat{\mathbf{w}}_j)$  is thus an estimate of the generalization error, defined as the expected loss,

$$G(\hat{\mathbf{w}}_j) = E_{\mathbf{x}, \mathbf{y}}\{\ell(\mathbf{y}, \hat{\mathbf{y}}; \hat{\mathbf{w}}_j)\} = \int \ell(\mathbf{y}, \hat{\mathbf{y}}; \hat{\mathbf{w}}_j) \cdot p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} \quad (3)$$

where  $p(\mathbf{x}, \mathbf{y})$  is the unknown joint input-output probability density. Generally,  $S_{\mathcal{V}_j}(\hat{\mathbf{w}}_j) = G(\hat{\mathbf{w}}_j) + O(1/\sqrt{N_{vj}})$  where  $O(\cdot)$  is the Landau order function<sup>6</sup>. Thus we need large  $N_{vj}$  to achieve an accurate estimate of the generalization error. On the other hand, this leaves only few data for training thus the true generalization  $G(\hat{\mathbf{w}}_j)$  increases. Consequently there exist a trade-off among the two conflicting aims which calls for finding an optimal split ratio. The optimal split ratio<sup>7</sup> is an interesting open and difficult problem since it depends on the total algorithm in which the validation error enters. Moreover, it depends on the learning curve<sup>8</sup> [16].

<sup>5</sup> That is, the loss function on the validation set.

<sup>6</sup> If  $h(x) = O(g(x))$  then  $|h(x)|/|g(x)| < \infty$  for  $x \rightarrow 0$ .

<sup>7</sup> For more elaborations on the split of data, see e.g., [2], [19], [23] and [25].

<sup>8</sup> Defined as the average generalization error as a function of the number of training examples.

The final  $K$ -fold cross-validation estimate is given by the average validation error estimates,

$$\widehat{\Gamma} = \frac{1}{K} \sum_{j=1}^K S_{\mathcal{V}_j}(\widehat{\mathbf{w}}_j). \quad (4)$$

$\widehat{\Gamma}$  is an estimate of the *average* generalization error over all possible training sets of size  $N_{t_j}$ ,

$$\Gamma = E_{\mathcal{T}}\{G(\widehat{\mathbf{w}}_j)\}. \quad (5)$$

$\widehat{\Gamma}$  is an unbiased estimate of  $\Gamma$  if the data of  $\mathcal{D}$  are independently distributed<sup>9</sup>, see e.g., [15].

The idea is now to optimize the amount of regularization by minimizing  $\widehat{\Gamma}$  w.r.t. the regularization parameters  $\boldsymbol{\kappa}$ . An algorithm for this purpose is described in Section 3. Furthermore, we might consider optimizing regularization using the hold-out validation estimate corresponding to  $K = 1$ . In this case one have to choose a split ratio. Without further ado, we will recommend a 50/50 splitting.

Suppose that we found the optimal  $\boldsymbol{\kappa}$  using the cross-validation estimate. Replications of training result in  $K$  different weight estimates  $\widehat{\mathbf{w}}_j$  which might be viewed as an ensemble of networks. In [15] we showed under certain mild conditions that when considering a  $o(1/N)$  approximation, the average generalization error of the ensemble network  $f_{\text{ens}}(\mathbf{x}) = \sum_{j=1}^K \beta_j \cdot f(\mathbf{x}, \widehat{\mathbf{w}}_j)$  equals that of the network trained on all examples in  $\mathcal{D}$  where  $\beta_j$  weights the contribution from the  $j$ 'th network and  $\sum_j \beta_j = 1$ . If  $K$  is a divisor in  $N$  then  $\forall j, \beta_j = 1/K$ , otherwise  $\beta_j = (N - N_{v_j})/N(K - 1)$ . Consequently, one might use the ensemble network to compensate for the increase in generalization error due to only training on  $N_{t_j} = N - N_{v_j}$  data. Alternatively, one might retrain on the full data set  $\mathcal{D}$  using the optimal  $\boldsymbol{\kappa}$ . We use the latter approach in the experimental section.

A minimal necessary requirement for a procedure which estimates the network parameters on the training set and optimizes the amount of regularization from a cross-validation set is: the generalization error of the regularized network should be smaller than that of the unregularized network trained on the full data set  $\mathcal{D}$ . However, this is not always the case, and is the quintessence of various “no free lunch” theorems [11], [43], [45]:

- If the regularizer is parameterized using many parameters,  $\boldsymbol{\kappa}$ , there is a potential risk of over-fitting on the cross-validation data. A natural way to avoid this situation is to limit the number of regularization parameters. Another recipe is to impose constraints on  $\boldsymbol{\kappa}$  (hyper regularization).
- The specific choice of the regularizers functional form impose prior constraints on the functions to be implemented by the network<sup>10</sup>. If the prior information is mismatched to the actual problem it might be better not to use regularization.

<sup>9</sup> That is,  $[\mathbf{x}(k_1), \mathbf{y}(k_1)]$  is independent of  $[\mathbf{x}(k_2), \mathbf{y}(k_2)]$  for all  $k_1 \neq k_2$ .

<sup>10</sup> The functional constraints are through the penalty imposed on the weights.

- The de-biasing procedure described above which compensate for training only on  $N_{tj} < N$  examples might fail to yield better performance since the weights now are optimized using all data, including those which were left out exclusively for optimizing regularization parameters.
- The split among training/validation data, and consequently the number of folds,  $K$ , may not be chosen appropriately.

These problems are further addressed in Section 4.1.

### 3 Adapting Regularization Parameters

The choice of regularizer may be motivated by

- the fact that the minimization of the cost function is normally an ill-posed task. Regularization smoothens the cost function and thereby facilitates the training. The weight decay regularizer<sup>11</sup>, originally suggested by Hinton in the neural networks literature, is a simple way to accomplish this task, see e.g., [34].
- a priori knowledge of the weights, e.g., in terms of a prior distribution (when using a Bayesian approach). In this case the regularization term normally plays the role of a log-prior distribution. Weight decay regularization may be viewed as a Gaussian prior, see e.g., [5]. Other types of priors, e.g., the Laplacian [12], [42] and soft weight sharing [33] has been considered. Moreover, priors have been developed for the purpose of restricting the number of weights (pruning), e.g., the so-called weight elimination [41].
- a desired characteristics of the functional mapping performed by the network. Typically, a smooth mapping is preferred. Regularizers which penalizes curvature of the mapping has been suggested in [4], [7], [31], [44].

In the experimental section we consider weight decay regularization and some generalizations hereof. Without further ado, weight decay regularization has proven to be useful in many neural network applications.

The standard approach for estimation of regularization parameters is more and less systematic search and evaluation of the cross-validation error. However, this is not viable for multiple regularization parameters. On the other hand, as will be demonstrated, it is possible to derive an optimization algorithm based on gradient descent.

Consider a regularization term  $R(\mathbf{w}, \boldsymbol{\kappa})$  which depends on  $q$  regularization parameters contained in the vector  $\boldsymbol{\kappa}$ . Since the estimated weights  $\hat{\mathbf{w}}_j = \arg \min_{\mathbf{w}} C_{\mathcal{T}_j}(\mathbf{w})$  are controlled by the regularization term, we may in fact consider the cross-validation error Eq. (4) as an *implicit function* of the regularization parameters, i.e.,

$$\hat{\Gamma}(\boldsymbol{\kappa}) = \frac{1}{K} \sum_{j=1}^K S_{\mathcal{V}_j}(\hat{\mathbf{w}}_j(\boldsymbol{\kappa})) \quad (6)$$

---

<sup>11</sup> Also known as ridge regression.

where  $\hat{\mathbf{w}}_j(\boldsymbol{\kappa})$  is the  $\boldsymbol{\kappa}$ -dependent vector of weights estimated from training set  $\mathcal{T}_j$ . The optimal regularization can be found by using gradient descent<sup>12</sup>,

$$\boldsymbol{\kappa}_{(n+1)} = \boldsymbol{\kappa}_{(n)} - \eta \frac{\partial \hat{\Gamma}}{\partial \boldsymbol{\kappa}}(\hat{\mathbf{w}}(\boldsymbol{\kappa}_{(n)})) \quad (7)$$

where  $\eta > 0$  is a step-size (learning rate) and  $\boldsymbol{\kappa}_{(n)}$  is the estimate of the regularization parameters in iteration  $n$ .

Suppose the regularization term is linear in the regularization parameters,

$$R(\mathbf{w}, \boldsymbol{\kappa}) = \boldsymbol{\kappa}^\top \mathbf{r}(\mathbf{w}) = \sum_{i=1}^q \kappa_i r_i(\mathbf{w}) \quad (8)$$

where  $\kappa_i$  are the regularization parameters and  $r_i(\mathbf{w})$  the associated regularization functions. Many suggested regularizers are linear in the regularization parameters, this includes the popular weight decay regularization as well as regularizers imposing smooth functions such as the Tikhonov regularizer [4], [5] and the smoothing regularizer for neural networks [31], [44]. However, there exist exceptions such as weight-elimination [41] and soft weight sharing [33]. In this case the presented method needs few modifications.

Using the results of the Appendix, the gradient of the cross-validation error equals

$$\frac{\partial \hat{\Gamma}}{\partial \boldsymbol{\kappa}}(\boldsymbol{\kappa}) = \frac{1}{K} \sum_{j=1}^K \frac{\partial S_{\mathcal{V}_j}}{\partial \boldsymbol{\kappa}}(\hat{\mathbf{w}}_j), \quad (9)$$

$$\frac{\partial S_{\mathcal{V}_j}}{\partial \boldsymbol{\kappa}}(\hat{\mathbf{w}}_j) = -\frac{\partial \mathbf{r}}{\partial \mathbf{w}^\top}(\hat{\mathbf{w}}_j) \cdot \mathbf{J}_j^{-1}(\hat{\mathbf{w}}_j) \cdot \frac{\partial S_{\mathcal{V}_j}}{\partial \mathbf{w}}(\hat{\mathbf{w}}_j). \quad (10)$$

where  $\mathbf{J}_j = \partial^2 C_{\mathcal{T}_j} / \partial \mathbf{w} \partial \mathbf{w}^\top$  is the Hessian of the cost function. As an example, consider the case of weight decay regularization with separate weight decays for two group of weights, e.g., the input-to-hidden and hidden-to output weights, i.e.,

$$R(\mathbf{w}, \boldsymbol{\kappa}) = \kappa^I \cdot |\mathbf{w}^I|^2 + \kappa^H \cdot |\mathbf{w}^H|^2 \quad (11)$$

where  $\boldsymbol{\kappa} = [\kappa^I, \kappa^H]$ ,  $\mathbf{w} = [\mathbf{w}^I, \mathbf{w}^H]$  with  $\mathbf{w}^I$ ,  $\mathbf{w}^H$  denoting the input-to-hidden and hidden-to output weights, respectively. The gradient of the validation error then yields,

$$\frac{\partial S_{\mathcal{V}_j}}{\partial \kappa^I}(\hat{\mathbf{w}}_j) = -2(\hat{\mathbf{w}}_j^I)^\top \cdot \mathbf{g}_j^I, \quad \frac{\partial S_{\mathcal{V}_j}}{\partial \kappa^H}(\hat{\mathbf{w}}_j) = -2(\hat{\mathbf{w}}_j^H)^\top \cdot \mathbf{g}_j^H \quad (12)$$

where  $\mathbf{g}_j$  is the vector

$$\mathbf{g}_j = [\mathbf{g}_j^I, \mathbf{g}_j^H] = \mathbf{J}_j^{-1}(\hat{\mathbf{w}}_j) \cdot \frac{\partial S_{\mathcal{V}_j}}{\partial \mathbf{w}}(\hat{\mathbf{w}}_j). \quad (13)$$

In summary, the algorithm for adapting regularization parameters consists of the following 8 steps:

<sup>12</sup> We have recently extended this algorithm incorporating second order information via the Conjugate Gradient technique [10].

1. Choose the split ratio; hence, the number of folds,  $K$ .
2. Initialize  $\boldsymbol{\kappa}$  and the weights of the network<sup>13</sup>.
3. Train the  $K$  networks with fixed  $\boldsymbol{\kappa}$  on  $\mathcal{T}_j$  to achieve  $\hat{\boldsymbol{w}}_j(\boldsymbol{\kappa})$ ,  $j = 1, 2, \dots, K$ . Calculate the validation errors  $S_{\mathcal{V}_j}$  and the cross-validation estimate  $\hat{\Gamma}$ .
4. Calculate the gradients  $\partial S_{\mathcal{V}_j} / \partial \boldsymbol{\kappa}$  and  $\partial \hat{\Gamma} / \partial \boldsymbol{\kappa}$  cf. Eq. (9) and (10). Initialize the step-size  $\eta$ .
5. Update  $\boldsymbol{\kappa}$  using Eq. (7).
6. Retrain the  $K$  networks from the previous weight estimates and recalculate the cross-validation error  $\hat{\Gamma}$ .
7. If no decrease in cross-validation error then perform a bisection of  $\eta$  and go to step 5; otherwise, continue.
8. Repeat steps 4–7 until the relative change in cross-validation error is below a small percentage or, e.g., the 2-norm of the gradient  $\partial \hat{\Gamma} / \partial \boldsymbol{\kappa}$  is below a small number.

Compared to standard neural network training the above algorithm does generally not lead to severe computational overhead. First of all, the standard approach of tuning regularization parameters by, more or less systematic search, requires a lot of training sessions. The additional terms to be computed in the adaptive algorithm are: 1) the derivative of the regularization functions w.r.t. the weights,  $\partial \boldsymbol{r} / \partial \boldsymbol{w}$ , 2) the gradient of the validation errors,  $\partial S_{\mathcal{V}_j} / \partial \boldsymbol{w}$ , and 3) the inverse Hessians,  $\boldsymbol{J}_j^{-1}$ . The first term is often a simple function of the weights<sup>14</sup> and computationally inexpensive. In the case of feed-forward neural networks, the second term is computed by one pass of the validation examples through a standard back-propagation algorithm. The third term is computationally more expensive. However, if the network is trained using a second order scheme, which requires computation of the inverse Hessian<sup>15</sup>, there is no computational overhead.

The adaptive algorithm requires of the order of  $K \cdot itr_{\boldsymbol{\kappa}} \cdot itr_{\eta}$  weight retrainings. Here  $itr_{\boldsymbol{\kappa}}$  is the number of iterations in the gradient descent scheme for  $\boldsymbol{\kappa}$  and  $itr_{\eta}$  is the average number of bisections of  $\eta$  in step 7 of the algorithm. In the experiments carried out the number of retrainings is approx. 100–300 times  $K$ . Recall, since we keep on retraining from the current weight estimate, the number of training epochs is generally small.

The number of weight retrainings is somewhat higher than that involved when optimizing the network by using a pruning technique like validation set based Optimal Brain Damage (vOBD) [24], [26]. vOBD based on  $K$ -fold cross-validation requires of the order of  $K \cdot m$  retrainings, where  $m = \dim(\boldsymbol{w})$ . The adaptive regularization algorithm is easily integrated with the pruning algorithm as demonstrated in the experimental section.

<sup>13</sup> In Sec. 4.1 a practical initialization procedure for  $\boldsymbol{\kappa}$  is described.

<sup>14</sup> For weight decay, it is  $2\boldsymbol{w}$ .

<sup>15</sup> Often the computations are reduced by using a Hessians approximation, e.g., the Gauss-Newton approximation. Many studies have reported significant training speed-up by using second order methods, see e.g., [21], [34].

## 4 Numerical Experiments

### 4.1 Potentials and Limitations in the Approach

The purpose of the section is to demonstrate the potential and limitations of the suggested adaptive regularization framework. We consider the simple linear data generating *system*, viz. estimating the mean of a Gaussian variable,

$$y(k) = w^\circ + \varepsilon(k) \quad (14)$$

where  $w^\circ$  is the true mean and the noise  $\varepsilon(k) \sim \mathcal{N}(0, \sigma_\varepsilon^2)$ .

We employ 2-fold cross-validation, i.e.,  $\mathcal{D} = \mathcal{T}_1 \cup \mathcal{T}_2$ , where  $\mathcal{T}_j$ ,  $j = 1, 2$  denote the two training sets in the validation procedure containing approximately half the examples<sup>16</sup>. The linear *model*  $y(k) = w + e(k)$  is trained using the mean square cost function augmented by simple weight decay, as shown by

$$C_{\mathcal{T}_j}(w) = \frac{1}{N_{t_j}} \sum_{k=1}^{N_{t_j}} (y(k) - w)^2 + \kappa \cdot w^2 \quad (15)$$

where  $k$  runs over examples of the data set in question. The estimated weights are  $\hat{w}_j = \bar{y}_j / (1 + \kappa)$  where  $\bar{y}_j = N_{t_j}^{-1} \sum_{k=1}^{N_{t_j}} y(k)$  are the estimated mean. For this simple case, the minimization of the cross-validation error given by,

$$\hat{\Gamma}(\kappa) = \frac{1}{2} \sum_{j=1}^2 S_{\mathcal{V}_j}(\hat{w}_j(\kappa)), \quad S_{\mathcal{V}_j}(\hat{w}_j(\kappa)) = \frac{1}{N_{v_j}} \sum_{k=1}^{N_{v_j}} (y(k) - \hat{w}_j)^2, \quad (16)$$

can be done exactly. The optimal  $\kappa$  is given by

$$\kappa_{\text{opt}} = \frac{\bar{y}_1^2 + \bar{y}_2^2}{2\bar{y}_1\bar{y}_2} - 1. \quad (17)$$

Assuming  $N$  to be even, the ensemble average of the estimated weights<sup>17</sup>,  $\hat{w}_j(\kappa_{\text{opt}})$ , leads to the final estimate

$$\hat{w}_{\text{reg}} = \frac{1}{2} (\hat{w}_1(\kappa_{\text{opt}}) + \hat{w}_2(\kappa_{\text{opt}})) = \frac{\bar{y}_1\bar{y}_2(\bar{y}_1 + \bar{y}_2)}{\bar{y}_1^2 + \bar{y}_2^2}. \quad (18)$$

Notice two properties: First, the estimate is self-consistent as  $\lim_{N \rightarrow \infty} \hat{w}_{\text{reg}} = \lim_{N \rightarrow \infty} \hat{w}_{\mathcal{D}} = w^\circ$  where  $\hat{w}_{\mathcal{D}} = N^{-1} \sum_{k=1}^N y(k) = (\bar{y}_1 + \bar{y}_2)/2$  is the unregularized estimate trained on all data. Secondly, it is easy to verify that  $\bar{y}_j \sim \mathcal{N}(w^\circ, 2\sigma_\varepsilon^2/N)$ . That is, if the *normalized true weight*  $\theta \equiv w^\circ/A$  where  $A = \sqrt{2/N} \cdot \sigma_\varepsilon$  is large then  $\bar{y}_j \approx w^\circ$  which means,  $\hat{w}_{\text{reg}} \approx \hat{w}_{\mathcal{D}}$ .

<sup>16</sup> That is,  $N_{t_1} = \lfloor N/2 \rfloor$  and  $N_{t_2} = N - N_{t_1}$ . Note that these training sets are also the two validation sets,  $\mathcal{V}_1 = \mathcal{T}_2$ , and vice versa.

<sup>17</sup> The ensemble average corresponds to retraining on all data using  $\kappa_{\text{opt}}$ . The weighting of the two estimates is only valid for  $N$  even (see Sec. 2 for the general case).



The objective is now to test whether using  $\hat{w}_{\text{reg}}$  results in lower generalization error than employing the unregularized estimate  $\hat{w}_{\mathcal{D}}$ . The generalization error associated with using the weight  $w$  is given by

$$G(w) = \sigma_\varepsilon^2 + (w - w^\circ)^2. \quad (19)$$

Further define the generalization error improvement,

$$Z = G(\hat{w}_{\mathcal{D}}) - G(\hat{w}_{\text{reg}}) = (\hat{w}_{\mathcal{D}} - w^\circ)^2 - (\hat{w}_{\text{reg}} - w^\circ)^2. \quad (20)$$

Note that  $Z$  merely is a function of the random variables  $\bar{y}_1, \bar{y}_2$  and the true weight  $w^\circ$ , i.e., it suffices to get samples of  $\bar{y}_1, \bar{y}_2$  when evaluating properties of  $Z$ . Define the normalized variables

$$\tilde{y}_j = \frac{\bar{y}_j}{A} \sim \mathcal{N}\left(\frac{w^\circ}{\sigma_\varepsilon} \cdot \sqrt{\frac{N}{2}}, 1\right) = \mathcal{N}(\theta, 1). \quad (21)$$

It is easily shown that the normalized generalization error improvement  $Z/A^2$  is a function of  $\tilde{y}_1, \tilde{y}_2$  and  $\theta$ ; hence, the distribution of  $Z/A^2$  is parameterized solely by  $\theta$ .

As a quality measure we consider the *probability of improvement* in generalization error given by  $\text{Prob}\{Z > 0\}$ . Note that  $\text{Prob}\{Z > 0\} = 1/2$  corresponds to equal preference of the two estimates. The probability of improvement depends only on the normalized weight  $\theta$  since  $\text{Prob}\{Z > 0\} = \text{Prob}\{Z/A^2 > 0\}$ .

Moreover, we consider the *relative generalization error improvement*, defined as

$$\text{RGI} = 100\% \cdot \frac{Z}{G(\hat{w}_{\mathcal{D}})}. \quad (22)$$

In particular, we focus on the probability that the relative improvement in generalization is bigger than<sup>18</sup>  $x$ , i.e.,  $\text{Prob}(\text{RGI} > x)$ . Optimally  $\text{Prob}(\text{RGI} > x)$  should be close to 1 for  $x \leq 0\%$  and slowly decaying towards zero for  $0\% < x \leq 100\%$ . Using the notation  $\tilde{w}_{\text{reg}} = \hat{w}_{\text{reg}}/A$ ,  $\tilde{w}_{\mathcal{D}} = \hat{w}_{\mathcal{D}}/A$ , RGI can be written as

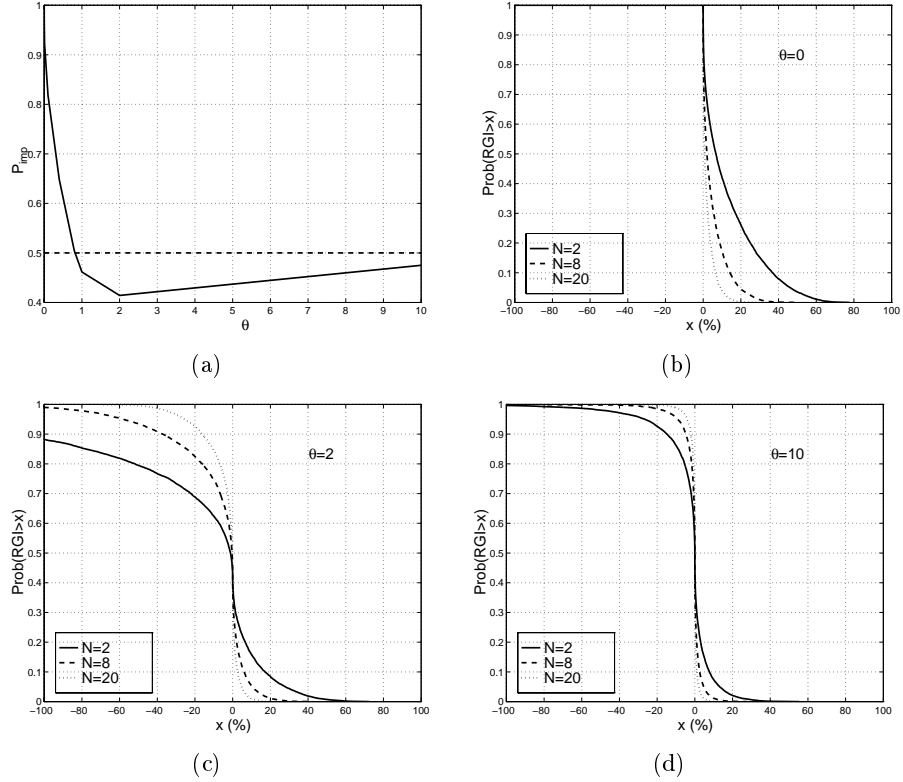
$$\text{RGI} = 100\% \cdot \frac{(\tilde{w}_{\mathcal{D}} - \theta)^2 - (\tilde{w}_{\text{reg}} - \theta)^2}{N/2 + (\tilde{w}_{\mathcal{D}} - \theta)^2}. \quad (23)$$

Thus, the distribution of RGI is parameterized by  $\theta$  and  $N$ .

The quality measures are computed by generating  $Q$  independent realizations of  $\tilde{y}_1, \tilde{y}_2$ , i.e.,  $\{\tilde{y}_1^{(i)}, \tilde{y}_2^{(i)}\}_{i=1}^Q$ . E.g., the probability of improvement is estimated by  $P_{\text{imp}} = Q^{-1} \sum_{i=1}^Q \mu(Z^{(i)})$  where  $\mu(x) = 1$  for  $x > 0$ , and zero otherwise.

The numerical results of comparing  $\hat{w}_{\text{reg}}$  to the unregularized estimate  $\hat{w}_{\mathcal{D}}$  is summarized in Fig. 1.

<sup>18</sup> Note that,  $\text{Prob}(\text{RGI} > 0) = \text{Prob}(Z > 0)$ .



**Fig. 1.** Result of comparing the optimally regularized estimate  $\hat{w}_{\text{reg}}$  of the mean of a Gaussian variable to the unregularized estimate  $\hat{w}_{\mathcal{D}}$ . The results are based on  $Q = 10^5$  independent realizations. The probability of improvement  $P_{\text{imp}}$ , shown in panel (a), is one for when the normalized true weight  $\theta = \sqrt{N/2} \cdot w^o / \sigma_\varepsilon = 0$ , and above 0.5 for  $\theta \lesssim 0.8$ . That is, when the prior information of the weight decay regularizer is correct (true weight close to zero), when  $N$  is small or when  $\sigma_\varepsilon$  is large. As  $\theta$  becomes large  $P_{\text{imp}}$  tends to 0.5 due to the fact that  $\tilde{w} \approx \hat{w}_{\mathcal{D}}$ . Panel (b)–(d) display  $\text{Prob}(\text{RGI} > x)$  for  $\theta \in \{0, 2, 10\}$ . The ideal probability curve is 1 for  $x < 0$  and a slow decay towards zero for  $x > 0$ . The largest improvement is attained for small  $\theta$  and small  $N$ . Panel (c) and (d) indicate that small  $N$  gives the largest probability for  $x > 0$ ; however, also the smallest probability for negative  $x$ . That is, a higher chance of getting a good improvement also increases the change of deterioration. Notice, even though  $P_{\text{imp}} < 0.5$  for  $\theta = 2, 10$  there is still a reasonable probability of getting a significant improvement.

## 4.2 Classification

We test the performance of the adaptive regularization algorithm on a vowel classification problem. The data are based on the Peterson and Barney database [35]. The classes are vowel sounds characterized by the first four formant frequen-

cies. 76 persons (33 male, 28 female and 15 children) have pronounced  $c = 10$  different vowels (IY IH EH AE AH AA AO UH UW ER) two times. This results in a data base of totally 1520 examples. The database is the verified database described in [40] where all data<sup>19</sup> are used, including examples where utterance failed of unanimous identification in the listening test (26 listeners). All examples were included to make the task more difficult.

The regularization was adapted using a hold-out validation error estimator, thus the examples were split into a data set,  $\mathcal{D}$ , consisting of  $N = 760$  examples (16 male, 14 female and 8 children) and an independent test set of the remaining 760 examples. The regularization was adapted by splitting the data set  $\mathcal{D}$  equally into a validation set of  $N_v = 380$  examples and a training set of  $N_t = 380$  examples (8 male, 7 female and 4 children in each set).

We used a feed-forward 2-layer neural network with hyperbolic tangent neurons in the hidden layer and modified SoftMax normalized outputs,  $\hat{y}_i$ , see e.g., [5], [17], [3]. Thus, the outputs estimates the posterior class probabilities  $p(\mathcal{C}_i|\mathbf{x})$ , where  $\mathcal{C}_i$  denotes the  $i$ 'th class,  $i = 1, 2, \dots, c$ . Bayes rule (see e.g., [5]) is used to assign  $\mathcal{C}_i$  to input  $\mathbf{x}$  if  $i = \arg \max_j p(\mathcal{C}_j|\mathbf{x})$ . Suppose that the network weights are given by  $\mathbf{w} = [\mathbf{w}^I, \mathbf{w}_{\text{bias}}^I, \mathbf{w}^H, \mathbf{w}_{\text{bias}}^H]$  where  $\mathbf{w}^I, \mathbf{w}^H$  are input-to-hidden and hidden-to-output weights, respectively, and the bias weights are assembled in  $\mathbf{w}_{\text{bias}}^I$  and  $\mathbf{w}_{\text{bias}}^H$ . Suppose that the targets  $y_i(k) = 1$  if  $\mathbf{x}(k) \in \mathcal{C}_i$ , and zero otherwise. The network is optimized using a log-likelihood loss function augmented by a weight decay regularizer using 4 regularization parameters,

$$C(\mathbf{w}) = \frac{1}{N_t} \sum_{k=1}^{N_t} \sum_{i=1}^c y_i(k) \log(\hat{y}_i(k, \mathbf{w})) \\ + \kappa^I \cdot |\mathbf{w}^I|^2 + \kappa_{\text{bias}}^I \cdot |\mathbf{w}_{\text{bias}}^I|^2 + \kappa^H \cdot |\mathbf{w}^H|^2 + \kappa_{\text{bias}}^H \cdot |\mathbf{w}_{\text{bias}}^H|^2. \quad (24)$$

We further define unnormalized weight decays as  $\alpha \equiv \kappa \cdot N_t$ . This regularizer is motivated by the fact that the bias, input and hidden layer weights play a different role, e.g., the input, hidden and bias signals normally have different scale (see also [5, Ch. 9.2]).

The simulation set-up is:

- Network: 4 inputs, 5 hidden neurons, 9 outputs<sup>20</sup>.
- Weights are initialized uniformly over  $[-0.5, 0.5]$ , regularization parameters are initialized at zero. One step in a gradient descent training algorithm (see e.g., [28]) is performed and the weight decays are re-initialized at  $\lambda_{\text{max}}/10^2$ , where  $\lambda_{\text{max}}$  is the max. eigenvalue of the Hessian matrix of the cost function. This initialization scheme is motivated by the following observations:
  - Weight decays should be so small that they do not reduce the approximation capabilities of the network significantly.

<sup>19</sup> The database can be retrieved from `ftp://eivind.imm.dtu.dk/dist/data/vowel/PetersonBarney.tar.Z`

<sup>20</sup> We only need 9 outputs since the posterior class probability of the 10th class is given by  $1 - \sum_{j=1}^9 p(\mathcal{C}_j|\mathbf{x})$ .

- They should be so large that the algorithm is prevented from being trapped in a local optimum and numerical instabilities are eliminated.
- Training is now done using a Gauss-Newton algorithm (see e.g., [28]). The Hessian is inverted using the Moore-Penrose pseudo inverse ensuring that the eigenvalue spread<sup>21</sup> is less than  $10^8$ .
- The regularization step-size  $\eta$  is initialized at 1.
- When the adaptive regularization scheme has terminated 3% of the weights are pruned using a validation set based version of the Optimal Brain Damage (vOBD) recipe [24], [26].
- Alternation between pruning and adaptive regularization continues until the validation error has reached a minimum.
- Finally, remaining weights are retrained on all data using the optimized weight decay parameters.

**Table 1.** Probability of misclassification ( $pmc$ ) and log-likelihood cost function (without reg. term, see Eq. (24)) for the classification example. The neural network averages and standard deviations are computed from 10 runs. In the case of small fixed regularization, weight decays were set at initial values equal to  $\lambda_{\max}/10^6$  where  $\lambda_{\max}$  is the largest eigenvalue of the Hessian matrix of the cost function. Optimal regularization refers to the case of optimizing 4 weight decay parameters. Pruning refers to validation set based OBD. KNN refers to  $k$ -nearest-neighbor classification.

**Probability of Misclassification ( $pmc$ )**

	NN small fixed reg.	NN opt. reg.+prun.	KNN ( $k = 9$ )
<b>Training Set</b>	$0.075 \pm 0.026$	$0.107 \pm 0.008$	0.150
<b>Validation Set</b>	$0.143 \pm 0.014$	$0.115 \pm 0.004$	0.158
<b>Test Set</b>	$0.146 \pm 0.010$	$0.124 \pm 0.006$	0.199
<b>Test Set (train. on all data)</b>	$0.126 \pm 0.010$	$0.119 \pm 0.004$	0.153

**Log-likelihood Cost Function**

	NN small fixed reg.	NN opt. reg.+prun.
<b>Training Set</b>	$0.2002 \pm 0.0600$	$0.2881 \pm 0.0134$
<b>Validation Set</b>	$0.7016 \pm 0.2330$	$0.3810 \pm 0.0131$
<b>Test Set</b>	$0.6687 \pm 0.2030$	$0.3773 \pm 0.0143$
<b>Test Set (train. on all data)</b>	$0.4426 \pm 0.0328$	$0.3518 \pm 0.0096$

Table 1 reports the average and standard deviations of the probability of

<sup>21</sup> Eigenvalue spread should not be larger than the square root of the machine precision [6].

misclassification ( $pmc$ ) and log-likelihood cost function over 10 runs for pruned networks using the optimal regularization parameters. Note that retraining on the full data set decreases the test  $pmc$  slightly on the average. In fact, improvement was noticed in 9 out of 10 runs. The table further shows the gain of the combined adaptive regularization/pruning algorithm relative to using a small fixed weight decay. However, recall, cf. Sec. 4.1, that the actual gain is *very* dependent on the noise level, data set size, etc. The objective is not to demonstrate high gain for a specific problem, rather to demonstrate that algorithm runs fairly robust in a classification set-up. For comparison we used a  $k$ -nearest-neighbor (KNN) classification (see e.g., [5]) and found that  $k = 9$  neighbors was optimal by minimizing  $pmc$  on the validation set. The neural network performed significantly better. Contrasting the obtained results to other work is difficult. In [36] results on the Peterson-Barney vowel problem are reported, but their data are not exactly the same; only the first 2 formant frequencies were used. Furthermore, different test sets have been used for the different methods presented. The best result reported [27] is obtained by using KNN and reach  $pmc = 0.186$  which is significantly higher than our results.

Fig. 2 shows the evolution of the adaptive regularization as well as the pruning algorithm.

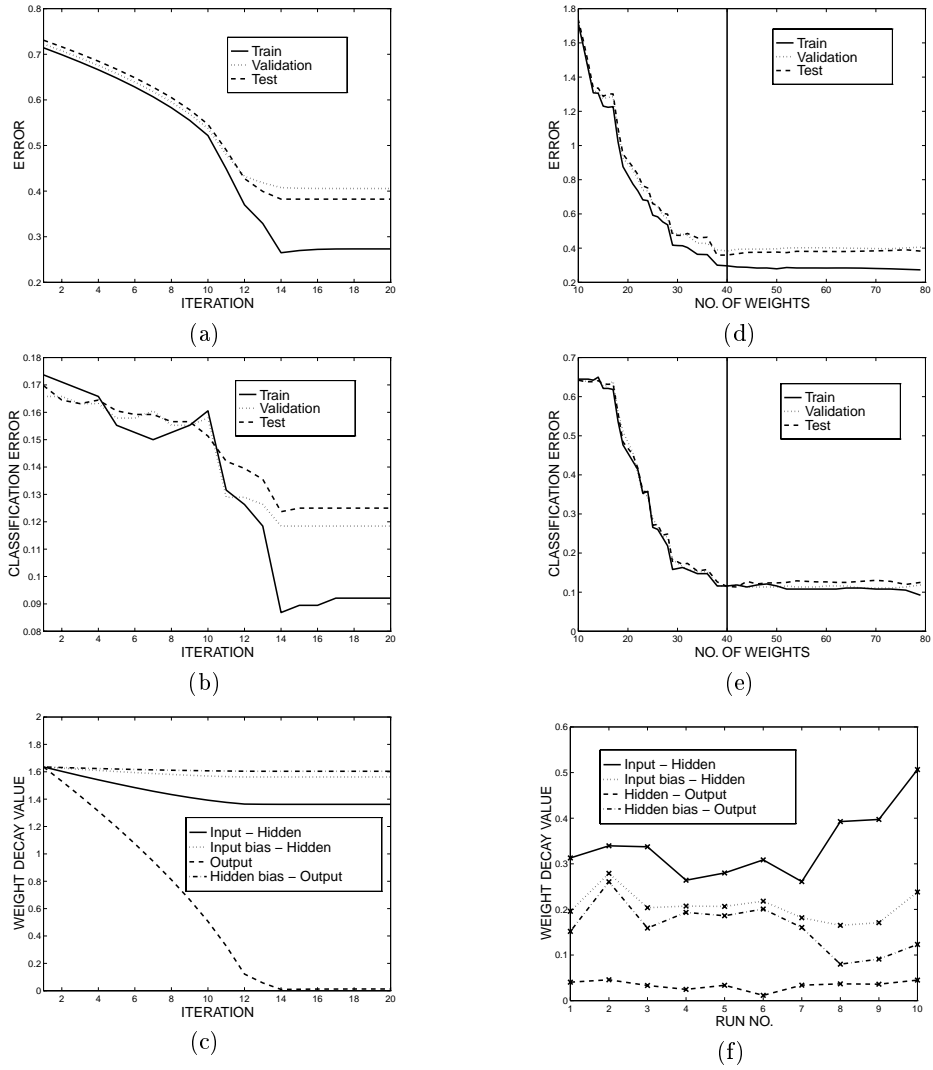
### 4.3 Time Series Prediction

We tested the performance of the adaptive regularization schemes on the Mackey-Glass chaotic time series prediction problem, see e.g., [21], [39]. The goal is to predict the series 100 steps ahead based on previous observations. The feed-forward net configuration is an input lag-space  $\mathbf{x}(k) = [x(k), x(k-6), x(k-12), x(k-18)]$  of 4 inputs, 25 hidden hyperbolic tangent neurons, and a single linear output unit  $\hat{y}(k)$  which predicts  $y(k) = x(k+100)$ . The cost function is the squared error,  $N_t^{-1} \sum_{k=1}^{N_t} (y(k) - \hat{y}(k, \mathbf{w}))^2$ , augmented by a weight decay regularizer using 4 different weight decays as described in Section 4.2.

The simulation set-up is:

- The data set,  $\mathcal{D}$ , has  $N = 500$  examples and an independent test has 8500 examples.
- The regularization parameters are optimized using a hold-out validation error with an even split<sup>22</sup> of the data set into training and validation sets each having 250 examples.
- Weight decays are initialized at zero and one Gauss-Newton iteration is performed, then weight decays were re-initialized at  $\lambda_{\max}/10^6$ , where  $\lambda_{\max}$  is the max. eigenvalue of the Hessian matrix of the cost function.
- The network is trained using a Gauss-Newton training scheme. The Hessian is inverted using the Moore-Penrose pseudo inverse ensuring that the eigenvalue spread is less than  $10^8$ .
- The regularization step-size  $\eta$  is initialized at  $10^{-2}$ .

<sup>22</sup> The sensitivity to different splits are considered in [24].



**Fig. 2.** Classification example. Panels (a), (b) and (c) show the evolution of the adaptive regularization algorithm in a typical run (fully connected network). The weight decays are optimized aiming at minimizing the validation error in panel (a). Note that also the test error decreases. This tendency is also evident in panel (b) displaying *pmc* even though a small increase noticed. In panel (c) the convergence unnormalized weight decays,  $\alpha = \kappa \cdot N_t$ , are depicted. Panels (d) and (e) show the evolution of errors and *pmc* during the pruning session. The optimal network is chosen as the one with minimal validation error, as indicated by the vertical line. There is only a marginal effect of pruning in this run. Finally, in panel (f), the variation of the optimal (end of pruning)  $\alpha$ 's in different runs is demonstrated. A clear similarity over runs is noticed.

- Alternation between adapting the 4 weight decays and validation set based pruning [24].
- The pruned network is retrained on all data using the optimized weight decay parameters.

Table 2 reports the average and standard deviations of the normalized squared error (i.e., the squared error normalized with the estimated variance of  $x(k)$ , denoted  $\hat{\sigma}_x^2$ ) over 10 runs for optimal regularization parameters. Retraining on the full data set decreases the test error somewhat on the average. Improvement was noticed in 10 out of 10 runs. We tested 3 different cases: small fixed regularization, small fixed regularization assisted by pruning and combined adaptive regularization/pruning. It turns that pruning alone does not improve performance; however, supplementing by adaptive regularization gives a test error reduction.

We furthermore tried a flexible regularization scheme, viz. individual weight decay where  $R(\mathbf{w}, \boldsymbol{\kappa}) = \sum_{i=1}^m \kappa_i w_i^2$  and  $\kappa_i \geq 0$  are imposed. In the present case it turned out that the flexible regularizer was not able to outperform the joint adaptive regularization/pruning scheme; possibly due to training and validation set sizes.

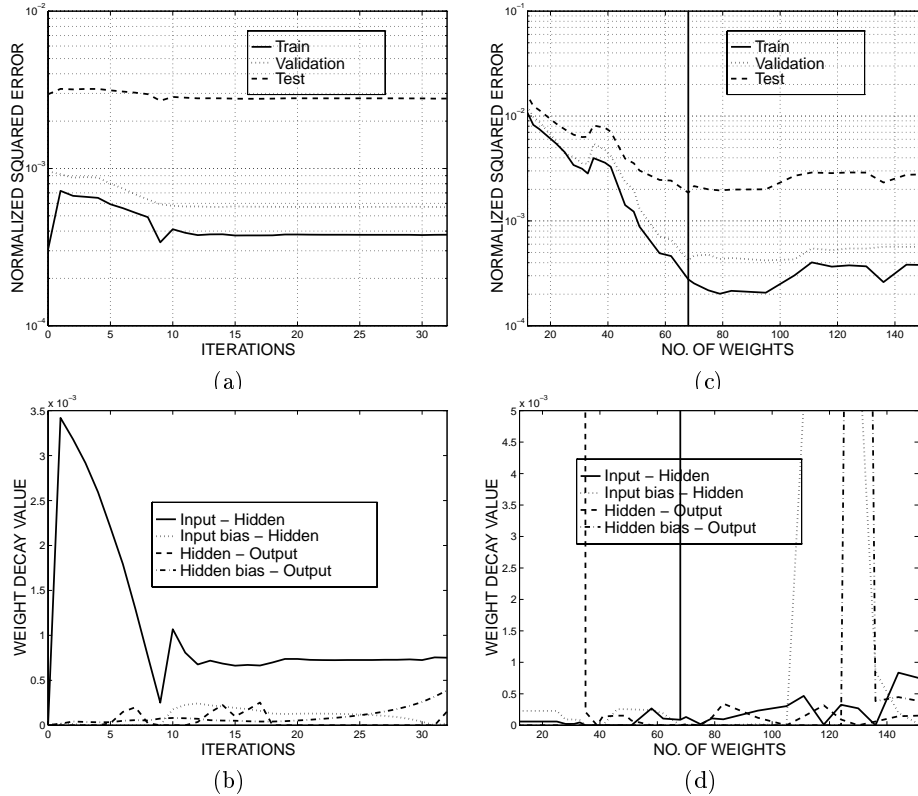
**Table 2.** Normalized squared error performance for the time series prediction examples. All figures are in units of  $10^{-3}\hat{\sigma}_x^2$  and averages and standard deviations are computed from 10 runs. In the case of small fixed regularization, weight decays were set at initial values equal to  $\lambda_{\max}/10^6$  where  $\lambda_{\max}$  is the largest eigenvalue of the Hessian matrix of the cost function. Optimal regularization refers to the case of optimizing 4 weight decay parameters. Pruning refers to validation set based OBD.

	NN small fixed reg.	NN small fixed reg.+prun.	NN opt. reg.+prun.
<b>Training Set</b>	0.17 ± 0.07	0.12 ± 0.04	0.10 ± 0.07
<b>Validation Set</b>	0.53 ± 0.26	0.36 ± 0.07	0.28 ± 0.14
<b>Test Set</b>	1.91 ± 0.68	1.58 ± 0.21	1.29 ± 0.46
<b>Test Set (train. on all data)</b>	1.33 ± 0.43	1.34 ± 0.26	1.17 ± 0.48

Fig. 3 demonstrates adaptive regularization and pruning in a typical case using 4 weight decays.

## 5 Conclusions

In this paper it was suggested to adapt regularization parameters by minimizing the cross-validation error or a simple hold-out validation error. We derived a simple gradient descent scheme for optimizing regularization parameters which has a small programming overhead and an acceptable computational overhead compared to standard training. Numerical examples with a toy linear model showed



**Fig. 3.** Time series prediction example. Panels (a) and (b) show a typical evolution of errors and unnormalized weight decay values  $\alpha$  when running the adaptive regularization algorithm using 4 weight decays. The normalized validation error drops approx. a factor of 2 when adapting weight decays. It turns out that some regularization of the input-to-hidden and output bias weights are needed whereas the other weights essentially requires no regularization<sup>23</sup>. In panel (c) and (d) it is demonstrated that pruning reduces the test error slightly. The optimal network is chosen as the one with minimal validation error, as indicated by the vertical line.

limitations and advantages of the adaptive regularization approach. Moreover, numerical experiments on classification and time series prediction problems successfully demonstrated the functionality of the algorithm. Adaptation of regularization parameters resulted in lower generalization error; however, it should be emphasized that the actual yield is very dependent on the problem and the choice of the regularizers functional form.

<sup>23</sup> Recall that if a weight decay  $\kappa$  is below  $\lambda_{\max}/10^8$  it does not influence the Moore-Penrose pseudo inversion of the Hessian.



## Acknowledgments

This research was supported by the Danish Natural Science and Technical Research Councils through the Computational Neural Network Center. JL furthermore acknowledge the Radio Parts Foundation for financial support. Mads Hintz-Madsen and Morten With Pedersen are acknowledged for stimulating discussions.

## Appendix

Assume that the regularization term is linear in the regularization parameters, i.e.,

$$R(\mathbf{w}, \boldsymbol{\kappa}) = \boldsymbol{\kappa}^\top \mathbf{r}(\mathbf{w}) = \sum_{i=1}^q \kappa_i r_i(\mathbf{w}) \quad (25)$$

The gradient of the cross-validation error Eq. (4) is

$$\frac{\partial \hat{I}}{\partial \boldsymbol{\kappa}}(\boldsymbol{\kappa}) = \frac{1}{K} \sum_{j=1}^K \frac{\partial S_{\mathcal{V}_j}}{\partial \boldsymbol{\kappa}}(\hat{\mathbf{w}}_j(\boldsymbol{\kappa})) \quad (26)$$

Using the chain rule the gradient vector of the validation error,  $S_{\mathcal{V}_j}$ , can be written as

$$\frac{\partial S_{\mathcal{V}_j}}{\partial \boldsymbol{\kappa}}(\hat{\mathbf{w}}_j(\boldsymbol{\kappa})) = \frac{\partial \mathbf{w}^\top}{\partial \boldsymbol{\kappa}}(\hat{\mathbf{w}}_j(\boldsymbol{\kappa})) \cdot \frac{\partial S_{\mathcal{V}_j}}{\partial \mathbf{w}}(\hat{\mathbf{w}}_j(\boldsymbol{\kappa})) \quad (27)$$

where  $\partial \mathbf{w}^\top / \partial \boldsymbol{\kappa}$  is the  $q \times m$  derivative matrix of the estimated weights w.r.t. the regularization parameters and  $m = \dim(\mathbf{w})$ . In order to find this derivative matrix, consider the gradient of the cost function w.r.t. to the weights as a function of  $\boldsymbol{\kappa}$  and use the following expansion around the current estimate  $\boldsymbol{\kappa}_{(n)}$ ,

$$\frac{\partial C_{\mathcal{T}_j}}{\partial \mathbf{w}}(\boldsymbol{\kappa}) = \frac{\partial C_{\mathcal{T}_j}}{\partial \mathbf{w}}(\boldsymbol{\kappa}_{(n)}) + \frac{\partial^2 C_{\mathcal{T}_j}}{\partial \mathbf{w} \partial \boldsymbol{\kappa}^\top}(\boldsymbol{\kappa}_{(n)}) \cdot (\boldsymbol{\kappa} - \boldsymbol{\kappa}_{(n)}) + o(|\boldsymbol{\kappa} - \boldsymbol{\kappa}_{(n)}|). \quad (28)$$

Requiring  $\hat{\mathbf{w}}(\boldsymbol{\kappa}_{(n+1)})$  in the next iteration to be an optimal weight vector, i.e.,  $\partial C_{\mathcal{T}_j} / \partial \mathbf{w}(\boldsymbol{\kappa}_{(n+1)}) = \mathbf{0}$  implies

$$\frac{\partial^2 C_{\mathcal{T}_j}}{\partial \mathbf{w} \partial \boldsymbol{\kappa}^\top}(\hat{\mathbf{w}}(\boldsymbol{\kappa}_{(n)})) = \mathbf{0}. \quad (29)$$

Recall that  $\partial C_{\mathcal{T}_j} / \partial \mathbf{w}(\boldsymbol{\kappa}_{(n)}) = \mathbf{0}$  by assumption. Eq. (29) can be used for determining  $\partial \mathbf{w}^\top / \partial \boldsymbol{\kappa}$ . Recognizing that the cost function  $C_{\mathcal{T}_j}(\hat{\mathbf{w}}(\boldsymbol{\kappa})) = S_{\mathcal{T}_j}(\hat{\mathbf{w}}(\boldsymbol{\kappa})) + R(\hat{\mathbf{w}}(\boldsymbol{\kappa}), \boldsymbol{\kappa})$  depends *implicitly* (through  $\hat{\mathbf{w}}(\boldsymbol{\kappa})$ ) and *explicitly* on  $\boldsymbol{\kappa}$  it is possible, by using Eq. (25), to derive the following relation<sup>24</sup>:

$$\frac{\partial \mathbf{w}^\top}{\partial \boldsymbol{\kappa}}(\hat{\mathbf{w}}_j) = -\frac{\partial \mathbf{r}}{\partial \mathbf{w}^\top}(\hat{\mathbf{w}}_j) \cdot \mathbf{J}_j^{-1}(\hat{\mathbf{w}}_j) \quad (30)$$

<sup>24</sup> For convenience, here  $\hat{\mathbf{w}}$ 's explicit  $\boldsymbol{\kappa}$ -dependence is omitted.

where  $\mathbf{J}_j = \partial^2 C_{\tau_j} / \partial \mathbf{w} \partial \mathbf{w}^\top$  is the Hessian of the cost function which e.g., might be evaluated using the Gauss-Newton approximation [28]. Finally, substituting Eq. (30) into (27) gives

$$\frac{\partial S_{v_j}}{\partial \kappa}(\hat{\mathbf{w}}_j) = -\frac{\partial \mathbf{r}}{\partial \mathbf{w}^\top}(\hat{\mathbf{w}}_j) \cdot \mathbf{J}_j^{-1}(\hat{\mathbf{w}}_j) \cdot \frac{\partial S_{v_j}}{\partial \mathbf{w}}(\hat{\mathbf{w}}_j) \quad (31)$$

$\partial S_{v_j} / \partial \mathbf{w}$  is found by ordinary back-propagation on the validation set while  $\partial \mathbf{r} / \partial \mathbf{w}^\top$  is calculated from the specific assumptions on the regularizer.

## References

1. H. Akaike: Fitting Autoregressive Models for Prediction. *Annals of the Institute of Statistical Mathematics* **21** (1969) 243–247
2. S. Amari, N. Murata, K.R. Müller, M. Finke & H. Yang: Asymptotic Statistical Theory of Overtraining and Cross-Validation. Technical report METR 95-06 (1995). Accepted for *IEEE Transactions on Neural Networks*. Available via `ftp://archive.cis.ohio-state.edu/pub/neuroprose/amari.overtraining.ps.Z`
3. L. Nonboe Andersen, J. Larsen, L.K. Hansen & M. Hintz-madsen: Adaptive Regularization of Neural Classifiers. In J. Principe *et al.* (eds.), *Proceedings of the IEEE Workshop on Neural Networks for Signal Processing VII*, Piscataway, New Jersey: IEEE, (1997) 24–33
4. C.M. Bishop: Curvature-Driven Smoothing: A Learning Algorithm for Feedforward Neural Networks. *IEEE Transactions on Neural Networks* **4**(4) (1993) 882–884
5. C.M. Bishop: *Neural Networks for Pattern Recognition*. Oxford, UK: Oxford University Press (1995)
6. J.E. Dennis & R.B. Schnabel: *Numerical Methods for Unconstrained Optimization and Non-linear Equations*. Englewood Cliffs, NJ: Prentice-Hall, (1983)
7. H. Drucker & Y. Le Cun: Improving Generalization Performance in Character Recognition. In B.H. Juang *et al.* (eds.), *Neural Networks for Signal Processing: Proceedings of the 1991 IEEE-SP Workshop*, Piscataway, New Jersey: IEEE (1991) 198–207
8. S. Geisser: The Predictive Sample Reuse Method with Applications. *Journal of the American Statistical Association* **50** (1975) 320–328
9. S. Geman, E. Bienenstock & R. Doursat: Neural Networks and the Bias/Variance Dilemma. *Neural Computation* **4** (1992) 1–58
10. C. Goutte & J. Larsen: Adaptive Regularization of Neural Networks using Conjugate Gradient, in *Proceedings of ICASSP'98, Seattle USA* **2** (1998) 1201–1204
11. C. Goutte: Note on Free Lunches and Cross-Validation. *Neural Computation* **9**(6) (1997) 1211–1215
12. C. Goutte: Regularization with a Pruning Prior. To appear in *Neural Networks* (1997)
13. L.K. Hansen and C.E. Rasmussen: Pruning from Adaptive Regularization. *Neural Computation* **6** (1994) 1223–1232
14. L.K. Hansen, C.E. Rasmussen, C. Svarer & J. Larsen: Adaptive Regularization. In J. Vlontzos, J.-N. Hwang & E. Wilson (eds.), *Proceedings of the IEEE Workshop on Neural Networks for Signal Processing IV*, Piscataway, New Jersey: IEEE (1994) 78–87

15. L.K. Hansen & J. Larsen: Linear Unlearning for Cross-Validation. *Advances in Computational Mathematics* **5** (1996) 269–280
16. J. Hertz, A. Krogh & R.G. Palmer: *Introduction to the Theory of Neural Computation*. Redwood City, California: Addison-Wesley Publishing Company (1991)
17. M. Hintz-Madsen, M. With Pedersen, L.K. Hansen, & J. Larsen: Design and Evaluation of Neural Classifiers. In S. Usui, Y. Tohkura, S. Katagiri & E. Wilson (eds.), *Proceedings of the IEEE Workshop on Neural Networks for Signal Processing VI*, Piscataway, New Jersey: IEEE, (1996) 223–232
18. K. Hornik: Approximation Capabilities of Multilayer Feedforward Networks. *Neural Networks* **4** (1991) 251–257
19. M. Kearns: A Bound on the Error of Cross Validation Using the Approximation and Estimation Rates, with Consequences for the Training-Test Split. *Neural Computation* **9**(5) (1997) 1143–1161
20. J. Larsen: A Generalization Error Estimate for Nonlinear Systems. In S.Y. Kung *et al.* (eds.), *Neural Networks for Signal Processing 2: Proceedings of the 1992 IEEE-SP Workshop*, Piscataway, New Jersey: IEEE (1992) 29–38
21. J. Larsen: *Design of Neural Network Filters*, Ph.D. Thesis, Electronics Institute, Technical University of Denmark (1993). Available via `ftp://eivind.imm.dtu.dk/dist/PhD_thesis/jlarsen.thesis.ps.Z`
22. J. Larsen & L.K. Hansen: Generalization Performance of Regularized Neural Network Models. In J. Vlontzos *et al.* (eds.), *Proceedings of the IEEE Workshop on Neural Networks for Signal Processing IV*, Piscataway, New Jersey: IEEE (1994) 42–51
23. J. Larsen & L.K. Hansen: Empirical Generalization Assessment of Neural Network Models. In F. Girosi *et al.* (eds.), *Proceedings of the IEEE Workshop on Neural Networks for Signal Processing V*, Piscataway, New Jersey: IEEE (1995) 30–39
24. J. Larsen, L.K. Hansen, C. Svarer & M. Ohlsson: Design and Regularization of Neural Networks: The Optimal Use of a Validation Set. In S. Usui, Y. Tohkura, S. Katagiri & E. Wilson (eds.), *Proceedings of the IEEE Workshop on Neural Networks for Signal Processing VI*, Piscataway, New Jersey: IEEE, (1996) 62–71
25. J. Larsen *et al.* : Optimal Data Set Split Ratio for Empirical Generalization Error Estimates. In preparation.
26. Y. Le Cun, J.S. Denker & S.A. Solla: Optimal Brain Damage. In D.S. Touretzky (ed.), *Advances in Neural Information Processing Systems 2*, *Proceedings of the 1989 Conference*, San Mateo, California: Morgan Kaufmann Publishers (1990) 598–605
27. D. Lowe: Adaptive Radial Basis Function Nonlinearities and the Problem of Generalisation. *Proc. IEE Conf. on Artificial Neural Networks*, (1989) 171– 175
28. L. Ljung: *System Identification: Theory for the User*. Englewood Cliffs, New Jersey: Prentice-Hall (1987)
29. D.J.C. MacKay: A Practical Bayesian Framework for Backprop Networks. *Neural Computation* **4**(3) (1992) 448–472
30. J. Moody: Prediction Risk and Architecture Selection for Neural Networks. In V. Cherkassky *et al.* (eds.), *From Statistics to Neural Networks: Theory and Pattern Recognition Applications*, Berlin, Germany: Springer-Verlag Series F **136** (1994)
31. J. Moody, T. Rögnavaldsson: Smoothing Regularizers for Projective Basis Function Networks. In *Advances in Neural Information Processing Systems 9*, *Proceedings of the 1996 Conference*, Cambridge, Massachusetts: MIT Press (1997)
32. N. Murata, S. Yoshizawa & S. Amari: Network Information Criterion — Determining the Number of Hidden Units for an Artificial Neural Network Model. *IEEE*

- Transactions on Neural Networks **5**(6) (1994) 865–872
33. S. Nowlan & G. Hinton: Simplifying Neural Networks by Soft Weight Sharing. *Neural Computation* **4**(4) (1992) 473–493
  34. M. With Pedersen: Training Recurrent Networks. In Proceedings of the IEEE Workshop on Neural Networks for Signal Processing VII, Piscataway, New Jersey: IEEE, (1997)
  35. G.E. Peterson & H.L. Barney: Control Methods Used in a Study of the Vowels. *JASA* **24** (1952) 175–184
  36. R.S. Shadafan & M. Niranjan: A Dynamic Neural Network Architecture by Sequential Partitioning of the Input Space. *Neural Computation* **6**(6) (1994) 1202–1222
  37. J. Sjöberg: Non-Linear System Identification with Neural Networks, Ph.D. Thesis no. 381, Department of Electrical Engineering, Linköping University, Sweden, (1995)
  38. M. Stone: Cross-validators Choice and Assessment of Statistical Predictors. *Journal of the Royal Statistical Society B* **36**(2) (1974) 111–147
  39. C. Svarer, L.K. Hansen, J. Larsen & C. E. Rasmussen: Designer Networks for Time Series Processing. In C.A. Kamm *et al.* (eds.), Proceedings of the IEEE Workshop on Neural Networks for Signal Processing 3, Piscataway, New Jersey: IEEE (1993) 78–87
  40. R.L. Watrous: Current Status of PetersonBarney Vowel Formant Data. *JASA* **89** (1991) 2459– 2460
  41. A.S. Weigend, B.A. Huberman & D.E. Rumelhart: Predicting the Future: A Connectionist Approach. *International Journal of Neural Systems* **1**(3) (1990) 193–209
  42. P.M. Williams: Bayesian Regularization and Pruning using a Laplace Prior. *Neural Computation* **7**(1) (1995) 117–143
  43. D.H. Wolpert & W.G. Macready: The Mathematics of Search. Technical Report SFI-TR-95-02-010, Santa Fe Institute (1995)
  44. L. Wu & J. Moody: A Smoothing Regularizer for Feedforward and Recurrent Neural Networks. *Neural Computation* **8**(3) 1996
  45. H. Zhu & R. Rohwer: No Free Lunch for Cross Validation. *Neural Computation* **8**(7) (1996) 1421–1426