

ON COMPARISON OF ADAPTIVE REGULARIZATION METHODS

Sigurdur Sigurdsson, Jan Larsen and Lars Kai Hansen
Department of Mathematical Modeling, Building 321
Technical University of Denmark, DK-2800 Lyngby, Denmark
Phone: +45 4525 3920,3923,3889
Email: siggi,jl,lkhansen@imm.dtu.dk, Web: eivind.imm.dtu.dk

INTRODUCTION

Modeling with flexible models, such as neural networks, requires careful control of the model complexity and generalization ability of the resulting model which finds expression in the ubiquitous bias-variance dilemma [4].

Regularization is a tool for optimizing the model structure reducing variance at the expense of introducing extra bias. The overall objective of adaptive regularization is to tune the amount of regularization ensuring minimal generalization error. Regularization is an supplement to direct model selection techniques like step-wise selection and one would prefer a hybrid scheme; however, a very flexible regularization may substitute the need for selection procedures.

This paper investigates recently suggested adaptive regularization schemes. Some methods focus directly on minimizing an estimate of the generalization error (either algebraic or empirical) [1], [3], [5], [6], [7], [12], [13], whereas others starts from different criteria, e.g., the Bayesian evidence [2, Ch. 10], [7], [15], [16]. The evidence expresses basically the probability of the model, which is conceptually different from generalization error; however, asymptotically for large training data sets they will converge¹ [15].

The papers is organized as follows: first the basic model definition, training and generalization is presented. Next, different adaptive regularization schemes are reviewed and extended. Finally, the experimental section presents a comparative study concerning linear models and feed-forward neural networks models for regression/time-series problems.

TRAINING AND GENERALIZATION

Suppose that our model, \mathcal{M} (e.g., neural network), is described by the function $\mathbf{f}(\mathbf{x}; \mathbf{w})$ where \mathbf{x} is the input vector and \mathbf{w} is the vector of parameters

¹Up to a scaling factor and an additive constant.

(or weights) with dimensionality m . The objective is to use the model for approximating the true conditional input-output distribution $p(\mathbf{y}|\mathbf{x})$, or some moments thereof. For regression and signal processing problems we normally model the conditional expectation $E\{\mathbf{y}|\mathbf{x}\}$. Define the training set $\mathcal{T} = \{\mathbf{x}(k); \mathbf{y}(k)\}_{k=1}^{N_{\mathcal{T}}}$ of $N_{\mathcal{T}}$ input-output examples sampled from the unknown but fixed joint input-output probability density $p(\mathbf{x}, \mathbf{y})$. The model is trained by minimizing a cost function, $C_{\mathcal{T}}(\mathbf{w})$, which is usually the sum of a loss function (or training error), $S_{\mathcal{T}}(\mathbf{w})$, and a regularization term $R(\mathbf{w}, \boldsymbol{\kappa})$ parameterized by a set of regularization parameters $\boldsymbol{\kappa}$,

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

where $\ell(\cdot)$ measures the cost of estimating the output $\mathbf{y}(k)$ with the model prediction $\hat{\mathbf{y}}(k; \mathbf{w}) = \mathbf{f}(\mathbf{x}(k); \mathbf{w})$, e.g., log-likelihood loss or the simple squared error loss function $\ell = \|\mathbf{y} - \hat{\mathbf{y}}\|^2$. Often we will consider linear regularization, i.e., $R(\mathbf{w}, \boldsymbol{\kappa}) = \boldsymbol{\kappa}^{\top} \mathbf{r}(\mathbf{w}) = \sum_{i=1}^q \kappa_i r_i(\mathbf{w})$ where $r_i(\mathbf{w})$ are associated regularization functions. Many suggested regularizers are linear; this includes the popular weight decay regularization and regularizers imposing smooth functions such as the Tikhonov regularizer [2]. Training provides the estimated weight vector $\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} C_{\mathcal{T}}(\mathbf{w})$. The *Generalization error* is defined as the expected loss on a future independent sample (\mathbf{x}, \mathbf{y}) ,

$$G(\hat{\mathbf{w}}) = E_{\mathbf{x}, \mathbf{y}}\{\ell(\mathbf{y}, \hat{\mathbf{y}}(\hat{\mathbf{w}}))\} = \int \ell(\mathbf{y}, \hat{\mathbf{y}}) p(\mathbf{x}, \mathbf{y}) d\mathbf{x}d\mathbf{y}, \quad (2)$$

and the *average generalization error* Γ is defined by averaging $G(\hat{\mathbf{w}})$ over all possible training sets:² $\Gamma = E_{\mathcal{T}}\{G(\hat{\mathbf{w}})\} = \int G(\hat{\mathbf{w}}) p(\mathcal{T}) d\mathcal{T}$.

ADAPTIVE REGULARIZATION

Validation Error Approach

Adapting regularization so as to minimize an empirical estimate of the generalization error, viz. the K -fold cross-validation [18], leads to an adaptive regularization scheme originally suggested in [12], which was further improved in [1], [3], [5], [13]. Suppose that all available data $\mathcal{D} = \{\mathbf{x}(k); \mathbf{y}(k)\}_{k=1}^N$ of N input-output examples, split 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 K -fold cross-validation estimate is then given by the average validation error estimates,

$$\hat{\Gamma}_{\text{cv}} = \frac{1}{K} \sum_{j=1}^K S_{\mathcal{V}_j}(\hat{\mathbf{w}}_j), \quad S_{\mathcal{V}_j}(\hat{\mathbf{w}}_j) = \frac{1}{N_{\mathcal{V}_j}} \sum_{k \in \mathcal{V}_j} \ell(\mathbf{y}(k), \hat{\mathbf{y}}(k; \hat{\mathbf{w}}_j)) \quad (3)$$

²For more details on empirical generalization error, generalization error distribution and average generalization error, see e.g., [11].

where $\hat{\mathbf{w}}_j$ are the weights estimated from training set \mathcal{T}_j . N_{v_j} is number of validation examples. $\hat{\Gamma}_{cv}$ is an estimate of the average generalization error over all possible training sets of size N_{t_j} , see [13].

The optimal regularization can be found by using gradient descent³,

$$\boldsymbol{\kappa}_{(n+1)} = \boldsymbol{\kappa}_{(n)} - \eta \frac{\partial \hat{\Gamma}_{cv}}{\partial \boldsymbol{\kappa}}(\hat{\mathbf{w}}(\boldsymbol{\kappa}_{(n)})) \quad (4)$$

where $\eta > 0$ is a step-size (learning rate) and $\boldsymbol{\kappa}_{(n)}$ is the estimate of the regularization parameters in iteration n . After convergence⁴ it is recommended to retrain on all available data using the optimized regularization parameters.

In case of linear regularization, the gradient of the cross-validation error can be written as [12], [13]

$$\frac{\partial \hat{\Gamma}_{cv}}{\partial \boldsymbol{\kappa}}(\boldsymbol{\kappa}) = \frac{1}{K} \sum_{j=1}^K \frac{\partial S_{v_j}}{\partial \boldsymbol{\kappa}}(\hat{\mathbf{w}}_j), \quad \frac{\partial S_{v_j}}{\partial \boldsymbol{\kappa}}(\hat{\mathbf{w}}_j) = -\frac{\partial \mathbf{r}}{\partial \mathbf{w}^\top}(\hat{\mathbf{w}}_j) \cdot \mathbf{J}_j^{-1}(\hat{\mathbf{w}}_j) \cdot \frac{\partial S_{v_j}}{\partial \mathbf{w}}(\hat{\mathbf{w}}_j). \quad (5)$$

where $\mathbf{J}_j(\mathbf{w}) = \partial^2 C_{\mathcal{T}_j}(\mathbf{w}) / \partial \mathbf{w} \partial \mathbf{w}^\top$ 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 of a neural network:

$$R(\mathbf{w}, \boldsymbol{\kappa}) = \kappa^I \cdot |\mathbf{w}^I|^2 + \kappa^H \cdot |\mathbf{w}^H|^2 \quad (6)$$

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,

$$\frac{\partial S_{v_j}}{\partial \kappa^I}(\hat{\mathbf{w}}_j) = -2(\hat{\mathbf{w}}_j^I)^\top \cdot \mathbf{g}_j^I, \quad \frac{\partial S_{v_j}}{\partial \kappa^H}(\hat{\mathbf{w}}_j) = -2(\hat{\mathbf{w}}_j^H)^\top \cdot \mathbf{g}_j^H \quad (7)$$

where \mathbf{g}_j is the vector $\mathbf{g}_j = [\mathbf{g}_j^I, \mathbf{g}_j^H] = \mathbf{J}_j^{-1}(\hat{\mathbf{w}}_j) \cdot \partial S_{v_j}(\hat{\mathbf{w}}_j) / \partial \mathbf{w}$.

Algebraic Generalization Error Approach

The literature suggests many algebraic estimators of the generalization error, including: FPER [10], GEN [8], GPE [14] and NIC [17]. The various estimators differ mainly in assumptions regarding model bias and dependence among data examples. In particular, they are all $o(1/N_{\mathcal{T}})$ estimators where $N_{\mathcal{T}}$ is the number of training examples. In many practical modeling scenarios the large training set assumption may be violated, however, the adaptive regularization based on this algebraic estimate might still be useful, as demonstrated in the experimental section. The major advantage of algebraic estimators is that all available data can be used to train the model,

³Optimization can be improved by using second order information [5], [3].

⁴E.g., small norm of gradient or small change in validation error.

i.e., $\mathcal{T} = \mathcal{D}$. This is not the case when using the empirical validation error approach discussed above.

In [7] 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 generalization error, and [6] proposed an adaptive regularization scheme for neural networks based on an algebraic estimate. In the following we present an extended version of this scheme where regularization parameters are adapted by an iterative gradient descent scheme aiming at minimizing the GEN/NIC [8], [9], [17], estimate of the generalization error. We use GEN/NIC in this work as an representative for the family of algebraic estimators. GEN/NIC has the advantage that model biased is not assumed negligible. The presented procedure can in principle be invoked for any of the mentioned estimators.

The $o(1/N_{\mathcal{T}})$ GEN/NIC estimate of the average generalization error is

$$\Gamma_{\text{GEN}} = E_{\mathcal{T}}\{S_{\mathcal{T}}(\hat{\mathbf{w}})\} + \frac{m_{\text{eff}}}{N_{\mathcal{T}}} - \frac{A}{N_{\mathcal{T}}} \cdot \frac{\partial R}{\partial \mathbf{w}^{\top}}(\mathbf{w}^*) \mathbf{J}^{-1}(\mathbf{w}^*) \frac{\partial R}{\partial \mathbf{w}}(\mathbf{w}^*) \quad (8)$$

where \mathbf{w}^* are the optimal model weights, i.e., $\mathbf{w}^* = \arg \min_{\mathbf{w}} G(\mathbf{w})$. m_{eff} is the effective number of parameters⁵ (weights) in the model [8], [9]

$$\begin{aligned} m_{\text{eff}} &= \text{tr} \left[\mathbf{J}^{-1}(\mathbf{w}^*) \left(\mathbf{K}(0) + \sum_{n=1}^{\bar{M}} \frac{N_{\mathcal{T}} - n}{N_{\mathcal{T}}} (\mathbf{K}(n) + \mathbf{K}^{\top}(n)) \right) \right] \\ &= \text{tr} [\mathbf{J}^{-1}(\mathbf{w}^*) \mathbf{L}] \end{aligned} \quad (9)$$

where $\bar{M} = \min(M, N_{\mathcal{T}} - 1)$, M is the time dependence length (for i.i.d. examples $M = 0$), $A = \bar{M} + 1 - \bar{M}(\bar{M} + 1)/2N_{\mathcal{T}}$, and $\mathbf{K}(n) = E\{\partial \ell(k)/\partial \mathbf{w} \cdot \partial \ell(k+n)/\partial \mathbf{w}^{\top}\}$ with $\ell(k) \equiv \ell(\mathbf{y}(k), \hat{\mathbf{y}}(k; \mathbf{w}^*))$, as $E\{\cdot\}$ denotes expectation w.r.t. joint input-output distribution. $\mathbf{J}(\mathbf{w})$ is the Hessian matrix of the expected cost function $E_{\mathcal{T}}\{C_{\mathcal{T}}(\mathbf{w})\}$, i.e., $\mathbf{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w}) + \partial^2 R/\partial \mathbf{w} \partial \mathbf{w}^{\top}$. If data are independent $\mathbf{K}(n) \equiv \mathbf{0}$ for $n > 0$, and if the cost is the log-likelihood loss then $\mathbf{K}(0)$ becomes the Hessian matrix of the unregularized cost, i.e., $\mathbf{K}(0) = \mathbf{H}(\mathbf{w}^*) = \partial^2 G(\mathbf{w}^*)/\partial \mathbf{w} \partial \mathbf{w}^{\top}$.

For practical implementation the quantities in Eq. (8) are estimated from data, as shown by,

$$\hat{\Gamma}_{\text{GEN}} = S_{\mathcal{T}}(\hat{\mathbf{w}}) + \frac{\hat{m}_{\text{eff}}}{N_{\mathcal{T}}} - \frac{A}{N_{\mathcal{T}}} \cdot \frac{\partial R}{\partial \mathbf{w}^{\top}}(\hat{\mathbf{w}}) \mathbf{J}_{\mathcal{T}}^{-1}(\hat{\mathbf{w}}) \frac{\partial R}{\partial \mathbf{w}}(\hat{\mathbf{w}}) \quad (10)$$

where \hat{m}_{eff} is calculated via Eq. (9) by substituting $\mathbf{J}_{\mathcal{T}}^{-1}(\hat{\mathbf{w}})$ for $\mathbf{J}^{-1}(\mathbf{w}^*)$,

$$\mathbf{J}_{\mathcal{T}}(\mathbf{w}) = \mathbf{H}_{\mathcal{T}}(\mathbf{w}) + \partial^2 R(\mathbf{w}, \kappa)/\partial \mathbf{w} \partial \mathbf{w}^{\top}, \quad \mathbf{H}_{\mathcal{T}}(\mathbf{w}) = \frac{\partial^2 S_{\mathcal{T}}}{\partial \mathbf{w} \partial \mathbf{w}^{\top}}(\mathbf{w}). \quad (11)$$

Further,

$$\mathbf{K}_{\mathcal{T}}(n) = \frac{1}{N_{\mathcal{T}}} \sum_{k=1}^{N_{\mathcal{T}}-n} \frac{\partial \ell(\mathbf{y}(k), \hat{\mathbf{y}}(k; \hat{\mathbf{w}}))}{\partial \mathbf{w}} \cdot \frac{\partial \ell(\mathbf{y}(k+n), \hat{\mathbf{y}}(k+n; \hat{\mathbf{w}}))}{\partial \mathbf{w}^{\top}} \quad (12)$$

⁵ For some cost functions, e.g., mean square error, m_{eff} is scaled by noise variance.

is substituted for $\mathbf{K}(n)$. To proceed as in [6], a simple gradient descent optimization as in Eq. (4) can be used⁶. The gradient of $\hat{\Gamma}_{\text{GEN}}$, noting that all quantities are evaluated at $\hat{\mathbf{w}} = \hat{\mathbf{w}}(\kappa)$, is according to Eq. (10) and (5) in case of linear regularization, given by

$$\frac{\partial \hat{\Gamma}_{\text{GEN}}}{\partial \kappa} = \frac{\partial S_{\mathcal{T}}}{\partial \kappa} + \frac{1}{N_{\mathcal{T}}} \frac{\partial \hat{m}_{\text{eff}}}{\partial \kappa} - \frac{A}{N_{\mathcal{T}}} \cdot \frac{\partial}{\partial \kappa} \left(\kappa^{\top} \frac{\partial \mathbf{r}}{\partial \mathbf{w}^{\top}} \mathbf{J}_{\mathcal{T}}^{-1} \frac{\partial \mathbf{r}^{\top}}{\partial \mathbf{w}} \kappa \right). \quad (13)$$

Eq. (13) can be written as

$$\begin{aligned} \frac{\partial \hat{\Gamma}_{\text{GEN}}}{\partial \kappa_i} &= \left(1 - \frac{2A}{N_{\mathcal{T}}} \right) \left\{ \frac{\partial \mathbf{r}}{\partial \mathbf{w}^{\top}} \mathbf{J}_{\mathcal{T}}^{-1} \frac{\partial \mathbf{r}^{\top}}{\partial \mathbf{w}} \kappa \right\}_i - \frac{1}{N_{\mathcal{T}}} \text{tr} \left[\mathbf{L}_{\mathcal{T}} \mathbf{J}_{\mathcal{T}}^{-1} \frac{\partial^2 r_i}{\partial \mathbf{w} \partial \mathbf{w}^{\top}} \mathbf{J}_{\mathcal{T}}^{-1} \right] \\ &\quad + \frac{A}{N_{\mathcal{T}}} \kappa^{\top} \frac{\partial \mathbf{r}}{\partial \mathbf{w}^{\top}} \mathbf{J}_{\mathcal{T}}^{-1} \frac{\partial^2 r_i}{\partial \mathbf{w} \partial \mathbf{w}^{\top}} \mathbf{J}_{\mathcal{T}}^{-1} \frac{\partial \mathbf{r}^{\top}}{\partial \mathbf{w}} \kappa. \end{aligned} \quad (14)$$

Evidence Approximation Approach

The Bayesian evidence approach adapts regularization parameters so as to minimize the evidence [2, Ch. 10], [15], [16]. The evidence is the probability of data⁷ given the model, $p(\mathcal{T}|\mathcal{M}) = \int p(\mathcal{T}|\mathbf{w}, \mathcal{M}) \cdot p(\mathbf{w}|\mathcal{M}) d\mathbf{w}$, where $p(\mathcal{T}|\mathbf{w}, \mathcal{M})$ is the likelihood and $p(\mathbf{w}|\mathcal{M})$ is the prior. In terms of the cost function components in Eq. (1) the likelihood and prior are expressed by:

$$p(\mathcal{T}|\mathbf{w}, \mathcal{M}) = Z_S^{-1} \exp(-\beta N_{\mathcal{T}} S_{\mathcal{T}}(\mathbf{w})), \quad p(\mathbf{w}|\mathcal{M}) = Z_R^{-1} \exp(-R(\mathbf{w}, \kappa)) \quad (15)$$

where β plays the role of the precision (inverse noise variance), and Z_S, Z_R are normalization constants. The evidence approximation framework consists in expanding the evidence to second order around the maximum a posteriori solution $\hat{\mathbf{w}}$. According to [15] the negative log-evidence is

$$\begin{aligned} -\log p(\mathcal{T}|\mathcal{M}) &\approx \beta N_{\mathcal{T}} S_{\mathcal{T}}(\hat{\mathbf{w}}) + R(\hat{\mathbf{w}}, \kappa) + \log Z_S + \log Z_R \\ &\quad + \frac{\log |\mathbf{J}_{\mathcal{T}}(\hat{\mathbf{w}})|}{2} + \frac{m}{2} (\log \beta + \log N_{\mathcal{T}} - \log 2\pi). \end{aligned} \quad (16)$$

If the likelihood and the weight prior are assumed to be Gaussian distributed⁸, which corresponds to using mean square loss and weight decay regularization as in Eq. (6), then the negative log-evidence is approximated by

$$\begin{aligned} -\log p(\mathcal{T}|\mathcal{M}) &\approx \beta N_{\mathcal{T}} S_{\mathcal{T}}(\hat{\mathbf{w}}) + \alpha^I |\hat{\mathbf{w}}^I|^2 + \alpha^H |\hat{\mathbf{w}}^H|^2 - \frac{(N_{\mathcal{T}} - m) \log \beta}{2} \\ &\quad - \frac{m^I \log \alpha^I + m^H \log \alpha^H}{2} + \frac{\log |\mathbf{J}_{\mathcal{T}}(\hat{\mathbf{w}})|}{2} + \frac{N_{\mathcal{T}} \log \pi + m(\log N_{\mathcal{T}} - \log 2)}{2} \end{aligned} \quad (17)$$

⁶[6] proceeds by finding the gradient of Eq. (8) and then use plug in estimates of unknown quantities. Here we proceed from the computable estimate Eq. (10). The difference between these approaches turn out to be minor.

⁷Also for this approach no validation data is required, i.e., $\mathcal{T} = \mathcal{D}$.

⁸If these assumptions are not fulfilled, the evidence framework becomes much more complicated and closed form solution can generally not be obtained. In such cases Monte Carlo techniques are required.

where $\alpha_i = \kappa_i/(\beta N_{\mathcal{T}})$ are the normalized weight decays, m^I, m^H are the number of hidden-to-input and hidden-to-output weights, respectively.

Minimizing the negative log-evidence by solving the equation, the derivative of Eq. (18) w.r.t. $\beta, \alpha^I, \alpha^H$ equal to zero⁹ results in optimal updated choices for $\beta, \alpha^I, \alpha^H$ based on current values ($\hat{\mathbf{w}} = \hat{\mathbf{w}}(\beta_{(n)}, \boldsymbol{\alpha}_{(n)})$),

$$\beta_{(n+1)} = \frac{N - \hat{m}_{\text{eff}}}{N_{\mathcal{T}} S_{\mathcal{T}}(\hat{\mathbf{w}})}, \quad \alpha_{(n+1)}^I = \frac{\hat{m}_{\text{eff}}^I}{2|\hat{\mathbf{w}}^I|^2}, \quad \alpha_{(n+1)}^H = \frac{\hat{m}_{\text{eff}}^H}{2|\hat{\mathbf{w}}^H|^2} \quad (18)$$

where $\hat{m}_{\text{eff}} = \hat{m}_{\text{eff}}^I + \hat{m}_{\text{eff}}^H$ and $\hat{m}_{\text{eff}}^I = \sum_{i \in \mathcal{I}} \lambda_i / (\lambda_i + \kappa^I)$, $\hat{m}_{\text{eff}}^H = \sum_{i \in \mathcal{H}} \lambda_i / (\lambda_i + \kappa^H)$ with \mathcal{I}, \mathcal{H} defining the indices for input-to-hidden and hidden-to-output weights, respectively, and λ_i is the i 'th eigenvalue of $\mathbf{H}_{\mathcal{T}}(\hat{\mathbf{w}})$. Thus, the evidence based scheme consists in alternating between weight optimization and update of weight decay and precision parameters, like the generalization based schemes.

EXPERIMENTS

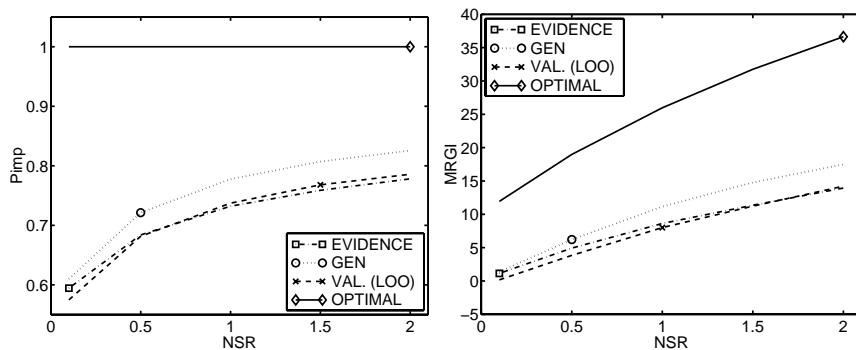


Figure 1: Results for Hessian with low eigenvalue spread. The GEN method is the most effective method for determining the regularization parameter. At a low NSR the effectiveness of all the methods is similar. The evidence method and leave-one-out (LOO) validation based method have similar performance. Optimal κ is found by exhaustive search.

Consider modeling a simple linear system $y(n) = \mathbf{x}^\top(n) \mathbf{w}^\circ + \epsilon(n)$. The input $\mathbf{x}(n) = [x_1(n), \dots, x_m(n)]^\top$ is a $m = 10$ dimensional i.i.d. Gaussian distributed vector $\mathbf{x}(n) \sim \mathcal{N}(\mathbf{0}, \mathbf{H})$, where \mathbf{H} is the covariance matrix. The true weight vector is $\mathbf{w}^* = [1, 1, 1, 0, 0, 0, 0, 0, 0, 0]^\top$. The noise $\epsilon(n) \sim \mathcal{N}(0, \sigma_\epsilon^2)$ is i.i.d. and independent of $\mathbf{x}(n)$. The noise variance is determined by $\sigma_\epsilon^2 = \text{NSR} \cdot (\mathbf{w}^*)^\top \mathbf{H} \mathbf{w}^\circ$, where NSR is the noise-to-signal ratio of the output. The weights are estimated using mean square error augmented by a simple weight decay, i.e., $\hat{\mathbf{w}} = \mathbf{J}_{\mathcal{T}}^{-1} \mathbf{X}^\top \mathbf{y} / N_{\mathcal{T}}$, where $\mathbf{y} = [y(1), \dots, y(N_{\mathcal{T}})]^\top$ and

⁹Some negligible terms are omitted, see [16].

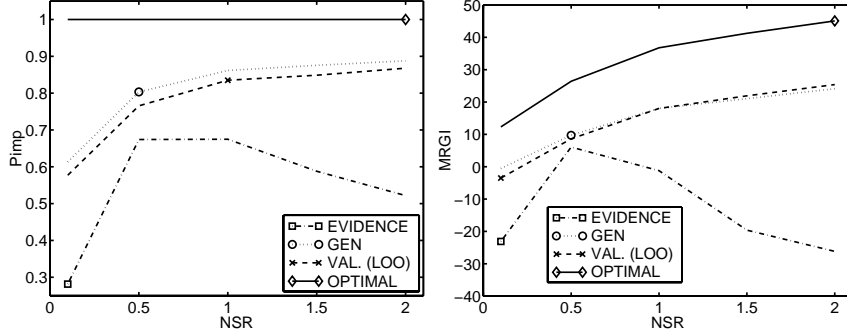


Figure 2: Results for Hessian with high eigenvalue spread. The GEN method has still the highest P_{imp} , but the MRGI is similar to the LOO validation based method. The evidence method has clearly the worst performance. This is caused by extremely low \hat{m}_{eff} , which seems to influence more the evidence method than the GEN. Notice that all methods have negative MRGI at NSR=0.1.

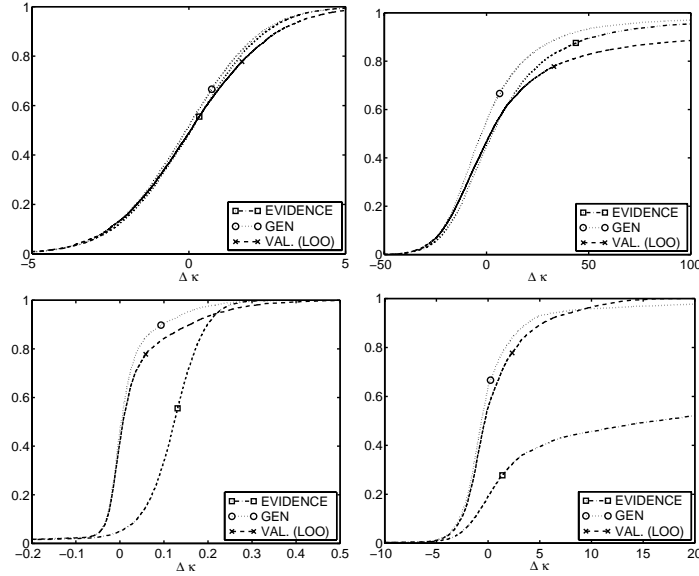


Figure 3: Cumulative probability distribution for the difference between the optimal κ and the κ 's suggested by the different regularization schemes. Left panels NSR = 0.1 and right panels NSR = 2.0. Top and bottom rows are low and high eigenvalue spread for the Hessian matrix, respectively. Top row: When the NSR is low the distributions are similar, while at high NSR the evidence and LOO validation method have a larger tail than the GEN, indicating that they are estimating κ too large. Bottom row: When the NSR is low the evidence method suggests κ 's that are too large, even though it does not have a large tail. When the NSR is large the evidence method has a very large tail, again estimating κ too large. The LOO validation method and the GEN show similar distributions.

$\mathbf{J}_{\mathcal{T}}$ is the Hessian of the regularized cost function given by $\mathbf{J}_{\mathcal{T}} = \mathbf{H}_{\mathcal{T}} + \kappa \mathbf{I}$ where $\mathbf{H}_{\mathcal{T}} = \mathbf{X}\mathbf{X}^{\top}/N_{\mathcal{T}}$ and $\mathbf{X} = [\mathbf{x}(1), \dots, \mathbf{x}(N_{\mathcal{T}})]^{\top}$. The true generalization error of the estimated linear system is easily computed as $G(\hat{\mathbf{w}}) = \sigma_{\epsilon}^2 + (\hat{\mathbf{w}} - \mathbf{w}^{\circ})^{\top} \mathbf{H}(\hat{\mathbf{w}} - \mathbf{w}^{\circ})$.

In order to evaluate the performance of the regularization methods, $Q = 5000$ independent data sets of size $N = 40$ are generated. Two measures are then used to evaluate the performance. The *probability of improvement* measures the fraction of the Q estimated models, using some regularization scheme, which generalize better than using no regularization, and is defined by $P_{\text{imp}} = Q^{-1} \sum_{i=1}^Q \mu(G(\hat{\mathbf{w}}_{\text{unreg}}^{(i)}) - G(\hat{\mathbf{w}}^{(i)}))$, where $\mu(x) = 1$ for $x > 0$, and zero otherwise. $G(\hat{\mathbf{w}}_{\text{unreg}}^{(i)})$ is the generalization error of the model trained on the i th data set using no regularization. The second performance measure is the *mean relative generalization error improvement* defined as $\text{MRGI} = Q^{-1} \sum_{i=1}^Q 100\% \left[G(\hat{\mathbf{w}}_{\text{unreg}}^{(i)}) - G(\hat{\mathbf{w}}^{(i)}) \right] / G(\hat{\mathbf{w}}_{\text{unreg}}^{(i)})$. Two different conditions for the Hessian are considered: small and large eigenvalue spread. Small eigenvalue spread is around 10, while more common large eigenvalue spread around 10^4 is obtained by multiplying a Vandermonde matrix \mathbf{A} to the original input \mathbf{X} using $\tilde{\mathbf{X}} = \mathbf{A}\mathbf{X}$ as the input. The methods are also compared at different NSR's. The performance is demonstrated through Fig. 1–3.

The computational complexity of the adaptive regularization schemes is very different. The leave-one-out (LOO) validation based method has the most computational overhead. Reestimating the weights $N_{\mathcal{T}}$ times is obviously very time consuming. Both the LOO validation based and GEN methods use gradient descent for estimating the regularization parameters. The convergence is very dependent on the step-size η Eq. (4). In particular, when the eigenvalue spread is high a small value has to be used, thus slowing down the convergence. The evidence method is much faster as Eq. (18) are analytical equations for regularization parameter updates.

CONCLUSION

This paper compared generalization error and evidence based schemes for adaptive regularization. We suggested various algorithm extensions and performed numerical experiments with linear models. The generalization error based methods generally performs good, while the evidence method yields comparable performance at low Hessian eigenvalue spread. However, at high eigenvalue spread, which is the common case in neural net applications, the evidence method has very low generalization error improvement.

Acknowledgments. Research supported by the Danish Research Councils through the THOR Center for Neuroinformatics and the Signal and Image Processing for Telemedicine (SITE) program.

REFERENCES

- [1] L.N. Andersen, J. Larsen, L.K. Hansen & M. Hintz-Madsen: "Adaptive Regularization of Neural Classifiers," in J. Principe *et al.* (eds.) **Proceedings of IEEE NNSP VII**, pp. 24-33, 1997.
- [2] C.M. Bishop: **Neural Networks for Pattern Recognition**, Oxford University Press, 1995.
- [3] D. Chen & M. Hagan: "Optimal Use of Regularization and Cross-Validation in Neural Network Modeling," in **Proceedings of IJCNN**, Washington DC, vol. 2, pp. 1275-1280, 1999.
- [4] S. Geman, E. Bienenstock & R. Doursat: "Neural Networks and the Bias/Variance Dilemma," **Neural Computation**, vol. 4, pp. 1-58, 1992.
- [5] C. Goutte & J. Larsen: "Adaptive Regularization of Neural Networks using Conjugate Gradient," in **Proceedings of ICASSP'98**, Seattle, USA, vol. 2, pp. 1201-1204, 1998.
- [6] L.K. Hansen, C.E. Rasmussen, C. Svarer & J. Larsen: "Adaptive Regularization," in J. Vlontzos *et al.* (eds.) **Proceedings of IEEE NNSP IV**, pp. 78-87, 1994.
- [7] L.K. Hansen, and C.E. Rasmussen: "Pruning from Adaptive Regularization," **Neural Computation**, vol. 6, pp. 1223-1232, 1994.
- [8] J. Larsen: A Generalization Error Estimate for Nonlinear Systems. In S.Y. Kung *et al.* (eds.), **Proceedings of IEEE NNSP II**, pp. 29-38, 1992.
- [9] J. Larsen, **Design of Neural Network Filters**, Ph.D. Thesis, Electronics Institute, Techn. Univ. of Denmark, March 1993.
- [10] J. Larsen & L.K. Hansen: Generalization Performance of Regularized Neural Network Models. In J. Vlontzos *et al.* (eds.), **Proceedings of IEEE NNSP IV**, pp. 42-51, 1994.
- [11] J. Larsen & L.K. Hansen: "Empirical Generalization Assessment of Neural Network Models," in F. Girosi *et al.* (eds.) **Proceedings of NNSP V**, Piscataway, New Jersey: IEEE, pp. 30-39, 1995.
- [12] 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 *et al.* (eds.), **Proceedings of IEEE NNSP VI**, pp. 62-71, 1996.
- [13] J. Larsen, C. Svarer, L. Nonboe Andersen & L.K. Hansen: "Adaptive Regularization in Neural Network Modeling," in G.B. Orr, K. Müller (eds.) **Neural Networks: Tricks of the Trade**, Lecture Notes in Computer Science 1524, Germany: Springer-Verlag, Chapter 5, pp. 113-132, 1998.
- [14] J. Moody: "The Effective Number of Parameters: An Analysis of Generalization and Regularization in Nonlinear Learning Systems," in J.E. Moody *et al.* (eds.) **Proceedings of NIPS 4**, Morgan Kaufmann Pub., pp. 847-854, 1992.
- [15] D.J.C. MacKay: "Bayesian Model Comparison and Backprop Nets," in J.E. Moody *et al.* (eds.) **Proceedings of NIPS 4**, Morgan Kaufmann Pub., pp. 839-846, 1992.
- [16] D.J.C. MacKay: "A Practical Bayesian Framework for Backprop Networks," **Neural Computation**, vol. 4, no. 3, pp. 448-472, 1992.
- [17] N. Murata, S. Yoshizawa & S. Amari: "Network Information Criterion — Determining the Number of Hidden Units for an Artificial Neural Network Model," **IEEE Trans. on NN**, vol. 5, no. 6, pp. 865-872, 1994.
- [18] M. Stone: "Cross-validated Choice and Assessment of Statistical Predictors," **Jour. Royal Stat. Soc. B** vol. 36, no. 2, pp. 111-147, 1974.