

Regularization of Neural Network Fitting

Stephen Vardeman
Analytics Iowa LLC
ISU Statistics and IMSE

Penalized training error for a neural net

Suppose that coordinates of input vectors in a training set have been standardized and one wants to regularize the fitting of a neural net. One way of proceeding is to define a penalty function like

$$J(\mathbf{A}) = \sum_{h=0}^H \sum_{i,j} (a_{ij}^h)^2 \quad (1)$$

for \mathbf{A} standing for the entire set of weights in $\mathbf{A}^0, \mathbf{A}^1, \dots, \mathbf{A}^H$ (it is not completely clear whether one wants to include weights on the "bias" terms in the neural net sums in (1)) and seek to optimize the penalized total training loss

$$\sum_{i=1}^N L(\hat{\mathbf{f}}_{\mathbf{A}}(\mathbf{x}_i), y_i) + \lambda J(\mathbf{A}) = N \cdot \overline{\text{err}} + \lambda J(\mathbf{A}) \quad (2)$$

for a $\lambda > 0$.

Gradient descent for penalized fitting

By modifying the basic gradient descent recursion previously discussed to

$$a_{\text{new}} = a_{\text{current}} - \gamma (D(a_{\text{current}}) + 2\lambda a_{\text{current}})$$

one arrives at an algorithm for optimizing the penalized training loss (2). Potentially, an appropriate value for λ might be chosen based on cross-validation.

Bayes analysis

Something that may at first seem quite different comes from a Bayesian point of view. For example, with a univariate regression model for outputs

$$y_i = f(\mathbf{x}_i | \mathbf{A}) + \epsilon_i$$

for the ϵ_i iid $N(0, \sigma^2)$, a likelihood is simply

$$l(\mathbf{A}, \sigma^2) = \prod_{i=1}^N h(y_i | f(\mathbf{x}_i | \mathbf{A}), \sigma^2)$$

for $h(\cdot | \mu, \sigma^2)$ the normal pdf.

If then $g(\mathbf{A}, \sigma^2)$ specifies a prior distribution for \mathbf{A} and σ^2 , a posterior for (\mathbf{A}, σ^2) has density proportional to

$$l(\mathbf{A}, \sigma^2) g(\mathbf{A}, \sigma^2)$$

Priors

For example, one might well assume that *a priori* the a s are iid $N(0, \eta^2)$ (where small η^2 will provide regularization and it is again unclear whether one wants to include the a s corresponding to bias terms in such an assumption or to instead provide more diffuse priors for them, like improper "Uniform($-\infty, \infty$)" or at least large variance normal ones). A standard improper prior for σ^2 is $\ln \sigma \sim \text{Uniform}(-\infty, \infty)$. In any case, whether improper or proper, abuse notation and write $g(\sigma^2)$ for a prior density for σ^2 .

Log-posterior and penalized training error

Then with independent mean 0 variance η^2 priors for all the weights (except possibly the ones for bias terms that might be given Uniform $(-\infty, \infty)$ priors), $\ln (l(\mathbf{A}, \sigma^2) g(\mathbf{A}, \sigma^2))$ is proportional to

$$\begin{aligned} & -NK \ln(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f(\mathbf{x}_i | \mathbf{A}))^2 - \frac{1}{2\eta^2} J(\mathbf{A}) + \ln g(\sigma^2) \\ & = -NK \ln(\sigma) + \ln g(\sigma^2) - \frac{1}{\sigma^2} \left(\sum_{i=1}^N (y_i - f(\mathbf{x}_i | \mathbf{A}))^2 + \frac{\sigma^2}{\eta^2} J(\mathbf{A}) \right) \quad (3) \end{aligned}$$

(flat improper priors for the bias weights correspond to the absence of terms for them in the sums for $J(\mathbf{A})$ in form (1)). This recalls display (2) and suggests that appropriate λ for regularization can be thought of as a variance ratio of "observation variance" and prior variance for the weights.

MCMC and maximum posterior density

It's clear how to define Metropolis-Hastings-within-Gibbs algorithms for sampling from $l(\mathbf{A}, \sigma^2) g(\mathbf{A}, \sigma^2)$. But typically the high dimensionality of the parameter space combined with the symmetry-derived multi-modality of the posterior will prevent one from running an MCMC algorithm long enough to fully detail the posterior. But detailing the posterior may not really be necessary or even desirable. Rather, one might simply run the MCMC algorithm monitoring the values of $l(\mathbf{A}, \sigma^2) g(\mathbf{A}, \sigma^2)$ corresponding to the successive MCMC iterates. An MCMC algorithm will spend much of its time where the corresponding posterior density is large and we can expect that a long MCMC run will identify a nearly modal value for the posterior. Rather than averaging neural nets according to the posterior, one might instead use as a predictor a neural net corresponding to a parameter vector (at least locally) maximizing the posterior.

MCMC and maximum posterior density cont.

One might even take the parameter vector in an MCMC run with the largest $l(\mathbf{A}, \sigma^2) g(\mathbf{A}, \sigma^2)$ value and for a grid of σ^2 values around the empirical maximizer use the back-propagation algorithm modified to fully optimize

$$\sum_{i=1}^N (y_i - f(\mathbf{x}_i | \mathbf{A}))^2 + \frac{\sigma^2}{\eta^2} J(\mathbf{A})$$

over choices of \mathbf{A} . This, in turn, could be used with relationship (3) to perhaps improve somewhat the result of the MCMC "search."