

Regularization in Neural Networks

1. Bayesian Interpretation of Regularization (Extension of Gaussian Priors)

$$\hat{E} = E_D(w) + \lambda E_w(w)$$

where $E_D(w)$ is the error of misclassification on the training set and $E_w(w)$ is the penalty form.

Interpreted in the Bayesian way:

$$P(w|D) \propto P(D|w)P(w)$$

considering:

$$P(D|w) = Z_D^{-1} \exp(-E_D) \text{ and } P(w) = Z_w^{-1} \exp(-\lambda E_w)$$

It follows that the process of minimization is equivalent to

$$\hat{E} = -\log P(D|w)P(w) + \text{const}$$

2. Approaches to Invariance

Let the transformation denoted by $s(x_n, \xi)$, which is parameterized by ξ . $s(x, 0) = x$. The tangent vector at the point x_n is given by:

$$\tau_n = \left. \frac{\partial s(x_n, \xi)}{\partial \xi} \right|_{\xi=0} \quad (5.125)$$

2.1 Regularization through Tangent Propagation

If the input vector x is transformed with parameter ξ , in general the network output will change. The derivative of network output with respect to ξ is given by:

$$\left. \frac{\partial y_k}{\partial \xi} \right|_{\xi=0} = \sum \frac{\partial y_k}{\partial x_i} \left. \frac{\partial x_i}{\partial \xi} \right|_{\xi=0} = \sum J_{ki} \tau_i \quad (5.126)$$

By adding this part into the error function to penalized variance to the transform, **the total error function** is modified to be:

$$\hat{E} = E_D(w) + \lambda \Omega \quad (5.127)$$

where

$$\Omega = \frac{1}{2} \sum_n \sum_k \left(\left. \frac{\partial y_{nk}}{\partial \xi} \right|_{\xi=0} \right)^2 = \frac{1}{2} \sum_n \sum_k \left(\sum_{i=1:D} J_{ki} \tau_i \right)^2 \quad (5.128)$$

2.2 Training with transformed data

By considering the average error function under a distribution $p(\xi)$ of the transformation parameter ξ , **the average error function** becomes:

$$\hat{E} = \frac{1}{2} \iiint \{y(s(x, \xi)) - t\}^2 p(t|x)p(x)p(\xi) dx dt d\xi \quad (5.130)$$

By taking the Taylor expansion of y with respect to x and ξ , in the average error function

$$\hat{E} = E_D + \lambda\Omega \quad (5.131)$$

$$\Omega = \int \left[\{y(x) - E[t|x]\} \frac{1}{2} \{(\tau')^T \nabla y(x) + \tau^T \nabla \nabla y(x) \tau\} + (\tau^T \nabla y(x))^2 \right] p(x) dx \quad (5.132)$$

2.3 Convolutional Networks

The convolutional network is composed of (1) receptive layer, (2) convolutional layer and (3) sub-sampling layer.

Weight sharing: in the convolutional layer all the units in a feature map are constrained to share the same weight.

Subsampling: each unit takes input from a small receptive field in the corresponding feature map of the convolutional layer. The average value of those input are calculated to be the input of a sigmoidal function.

2.4 Soft Weight Sharing

Soft weight sharing is designed to replace the hard constraint of equal weight in the convolutional network with a form of regularization.

Instead of considering a single Gaussian prior distribution of the weight, in the soft weight sharing the distribution is taken as a mixture of Gaussians:

$$p(w) = \prod_i p(w_i) \quad (5.136)$$

$$p(w_i) = \sum_{j=1:M} \pi_j N(w_i | \mu_j, \sigma_j^2) \quad (5.138)$$

So the regularization function takes the form:

$$\Omega(w) = - \sum_i \ln \left(\sum_{j=1}^M \pi_j N(w_i | \mu_j, \sigma_j^2) \right) \quad (5.139)$$

The minimization of the total error function

$$\hat{E} = E_D(w) + \lambda\Omega$$

is performed with the help of considering the mixing coefficient $\{\pi_j\}$ as prior probabilities,

so the constraint:

$$\sum_j \pi_j = 1 \quad (5.145)$$

is added.