# Soft-weight sharing

## Sargur Srihari

# Sharing of Weights

- One way of reducing complexity of a neural network with a large number of weights is:
  - Encouraging groups of weights to have similar values
- What we saw with Conv Nets is only applicable when form of the network can be specified in advance
  - Hard weight sharing:  Constraint of equal weights
- Alternative approach is Soft weight sharing:
  - Groups of weights encouraged to have similar values
- In soft-weight sharing, learning process determines
  - Division of weights into groups
  - Mean weight value for each group, and
  - Spread of values

2

# Mixture of Gaussians for Weights

- Simple weight decay regularizer for neural networks:

$$\tilde{E}(\mathrm{w}) = E(\mathrm{w}) + \frac{\lambda}{2}\mathrm{w}^\mathrm{T}\mathrm{w}$$

  where $\lambda$ is the regularization coefficient

  - This regularizer can be interpreted as negative log of a zero-mean Gaussian prior over $\mathrm{w}$
    - Useful for translation invariance

- Instead, encourage weight values to form several groups rather than a single group
  - Group weights into $M$ groups, each with a Gaussian prior distribution over the weights
    - A mixture distribution
    - Centers, variances and mixing coefficients are learnt from data

3

# Learning weight groups

- Probability density of weight groups

$$p(\mathbf{w}) = \prod_i p(w_i) \qquad p(w_i) = \sum_{j=1}^{M} \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2)$$

  - Negative log leads to regularization function

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

  - Total error function $\widetilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda \Omega(\mathbf{w})$

  - Error minimized wrt weights $w_i$ and wrt parameters $\{\pi_j, \mu_j, \sigma_j\}$

- Optimization method

  - If $w_i$ were fixed, we could use EM algorithm, but they are not

    - If $w_k$ were fixed this is a standard mixture of Gaussians problem

  - So use conjugate gradients or quasi Newton optimization

    - Conjugate gradients is faster than SGD but is sensitive

# Derivatives for Minimization of total error function

$$\widetilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda\Omega(\mathbf{w})$$

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

- Define

$$\gamma_j = \frac{\pi_j N\left(w \mid \mu_j, \sigma_j^2\right)}{\sum_k \pi_k N\left(w \mid \mu_k, \sigma_k^2\right)}$$

- Derivatives of total error function wrt weights $w_i$ is given by

$$\frac{\partial \tilde{E}}{\partial w_i} = \frac{\partial E}{\partial w_i} + \lambda\sum_j \gamma_j\left(w_i\right)\frac{(w_i - \mu_j)}{\sigma_j^2}$$

- Derivatives of error wrt the centers of Gaussians $\mu_j$

$$\frac{\partial \tilde{E}}{\partial \mu_j} = \lambda\sum_j \gamma_j\left(w_i\right)\frac{(\mu_i - w_j)}{\sigma_j^2}$$

- Derivatives of error wrt the variances $\sigma_j$

$$\frac{\partial \tilde{E}}{\partial \sigma_j} = \lambda\sum_j \gamma_j\left(w_i\right)\left(\frac{1}{\sigma_j} - \frac{\left(w_i - \mu_j\right)^2}{\sigma_j^2}\right)$$

- Derivatives of error wrt mixing coefficients

$$\frac{\partial \tilde{E}}{\partial \eta_j} = \lambda\sum_i \left\{\pi_j - \gamma_j\left(w_i\right)\right\}$$     where     $$\pi_j = \frac{\exp\left(\eta_j\right)}{\sum_{k=1}^{M} \exp\left(\eta_k\right)}$$

5