

Dropout as variational inference

Apostolos Psaros

1 Forward pass with dropout

Consider a NN with a single hidden layer. The hidden layer weights and biases are \mathbf{m}_1 and \mathbf{b} , respectively. These are vectors of size K , the width of the layer. The output weights are \mathbf{m}_2 , also of size K . For doing a forward pass with dropout we sample a vector $\hat{\epsilon} \sim p(\epsilon)$ of dimension K . The elements of $\hat{\epsilon}$ take value 0 with probability $0 \leq p \leq 1$. Given the output of the hidden layer $\mathbf{h} = \sigma(x\mathbf{m}_1 + \mathbf{b})$, we set a proportion of it to zero, i.e., $\hat{\mathbf{h}} = \mathbf{h} \odot \hat{\epsilon}$. Finally, $f_{\mathcal{H}}(x_i; \theta, \hat{\epsilon}) = \mathbf{m}_2^T \hat{\mathbf{h}}$, where $\theta = \{\mathbf{m}_1, \mathbf{m}_2, \mathbf{b}\}$.

2 Training with dropout

Training involves obtaining θ by minimizing

$$\mathcal{L}_{dropout}(\theta) = \sum_{i=1}^N \frac{1}{2} |y_i - f_{\mathcal{H}}(x_i; \theta, \epsilon)|^2 + Reg(\theta) \quad (1)$$

where ϵ is a random variable as described above and $Reg(\theta)$ is a regularization term. In each iteration, for obtaining the gradient of $\mathcal{L}_{dropout}$ we use a sample $\hat{\mathcal{L}}_{dropout}$ given as

$$\hat{\mathcal{L}}_{dropout}(\theta) = \sum_{i=1}^N \frac{1}{2} |y_i - f_{\mathcal{H}}(x_i; \theta, \hat{\epsilon})|^2 + Reg(\theta) \quad (2)$$

3 Equivalence with variational inference

Note that

$$\begin{aligned} f_{\mathcal{H}}(x_i; \theta, \hat{\epsilon}) &= \mathbf{m}_2^T \hat{\mathbf{h}} \\ &= \mathbf{m}_2^T \text{diag}(\hat{\epsilon}) \mathbf{h} \\ &= \hat{\mathbf{w}}_2^T \mathbf{h} \\ &= f_{\mathcal{H}}(x_i; \hat{\mathbf{w}}) \end{aligned} \quad (3)$$

where $\hat{\mathbf{w}} = \{\mathbf{m}_1, \hat{\mathbf{w}}_2, \mathbf{b}\}$ is a sample of the random variable $\mathbf{w} = \{\mathbf{m}_1, \mathbf{w}_2, \mathbf{b}\}$, with $\mathbf{w}_2 = \mathbf{m}_2 \odot \epsilon$ and $\hat{\mathbf{w}}_2 = \mathbf{m}_2 \odot \hat{\epsilon}$. We can therefore transfer the uncertainty from the feature space to the model weights. Therefore, the training loss in each iteration is

$$\hat{\mathcal{L}}_{dropout}(\theta) = \sum_{i=1}^N \frac{1}{2} |y_i - f_{\mathcal{H}}(x_i; \hat{\mathbf{w}})|^2 + Reg(\theta) \quad (4)$$

which can also be written as

$$\hat{\mathcal{L}}_{dropout}(\boldsymbol{\theta}) = -\beta \sum_{i=1}^N \log p(y_i | \hat{\mathbf{w}}, x_i, \mathcal{H}) + Reg(\boldsymbol{\theta}) \quad (5)$$

If we write $\mathbf{w} = \mathbf{g}(\boldsymbol{\theta}, \boldsymbol{\epsilon}) = \{\mathbf{m}_1, \mathbf{m}_2 \odot \boldsymbol{\epsilon}, \mathbf{b}\}$, Eq. (5) becomes

$$\hat{\mathcal{L}}_{dropout}(\boldsymbol{\theta}) = -\beta \sum_{i=1}^N \log p(y_i | \mathbf{g}(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}), x_i, \mathcal{H}) + Reg(\boldsymbol{\theta}) \quad (6)$$

Recall that the loss in Bayes by backprop is

$$\hat{\mathcal{L}}_{BBB}(\boldsymbol{\theta}) = -\sum_{i=1}^N \log p(y_i | \mathbf{g}(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i), x_i, \mathcal{H}) + KL(q_{\boldsymbol{\theta}}(\mathbf{w}) || p(\mathbf{w} | \mathcal{H})) \quad (7)$$

4 Summary

Training with dropout is equivalent to Bayes by backprop

1. with a difference in scale in the summation term
2. with reparametrization $\mathbf{g}(\boldsymbol{\theta}, \boldsymbol{\epsilon}) = \{\mathbf{m}_1, \mathbf{m}_2 \odot \boldsymbol{\epsilon}, \mathbf{b}\}$
3. with prior $p(\mathbf{w} | \mathcal{H})$ and approximating distribution $q_{\boldsymbol{\theta}}(\mathbf{w})$ such that $KL(q_{\boldsymbol{\theta}}(\mathbf{w}) || p(\mathbf{w} | \mathcal{H})) = Reg(\boldsymbol{\theta})$.

Two more notes:

1. Other stochastic regularization techniques can be recovered with different reparametrizations and $\mathbf{g}(\boldsymbol{\theta}, \boldsymbol{\epsilon})$
2. after training with dropout the NN can be used exactly as a BNN (MC dropout).

Overall, optimizing any NN with dropout is equivalent to *a form* of variational inference.