

General Gradient of Neural Network

tomas.ukkonen@iki.fi, 2017

Backpropagation is a commonly known algorithm for computing the gradient of error function which arises when we know target values and the loss, cost or error function is one dimensional. Generalizing this to general gradient calculation when we seek to find the maximum or minimum value of a neural network (thought often ill-fated because of local optimas produced by an over-fitted neural network) is then important. This is needed, for example, when implementing certain reinforcement learning methods.

Consider a two-layer neural network

$$y(\mathbf{x}) = f(\mathbf{W}^{(2)} g(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$$

The gradients of the final layer are (non-zero terms are at the j :th row):

$$\begin{aligned} \frac{\partial y(\mathbf{x})}{\partial w_{ji}^{(2)}} &= \text{diag}\left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}\right) \begin{pmatrix} 0 \\ g_i \\ 0 \end{pmatrix} \\ \frac{\partial y(\mathbf{x})}{\partial b_j^{(2)}} &= \text{diag}\left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}\right) \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \end{aligned}$$

The derivation chain-rule can be used to calculate the second (and more deep layers' gradients):

$$\begin{aligned} \frac{\partial y(\mathbf{x})}{\partial w_{ji}^{(1)}} &= \text{diag}\left(\frac{\partial f(\mathbf{x})}{\partial (\mathbf{W}^{(2)} \mathbf{g} + \mathbf{b}^{(2)})}\right) \frac{\partial (\mathbf{W}^{(2)} \mathbf{g} + \mathbf{b}^{(2)})}{\partial \mathbf{g}} \frac{\partial \mathbf{g}(\mathbf{x})}{\partial (\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})} \frac{\partial (\mathbf{W}^{(2)} \mathbf{x} + \mathbf{b}^{(2)})}{\partial w_{ji}^{(1)}} \\ \frac{\partial y(\mathbf{x})}{\partial w_{ji}^{(1)}} &= \text{diag}\left(\frac{\partial f(\mathbf{x})}{\partial (\mathbf{W}^{(2)} \mathbf{g} + \mathbf{b}^{(2)})}\right) \mathbf{W}^{(2)} \text{diag}\left(\frac{\partial \mathbf{g}(\mathbf{x})}{\partial (\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})}\right) \begin{pmatrix} 0 \\ x_i \\ 0 \end{pmatrix} \\ \frac{\partial y(\mathbf{x})}{\partial b_j^{(1)}} &= \text{diag}\left(\frac{\partial f(\mathbf{x})}{\partial (\mathbf{W}^{(2)} \mathbf{g} + \mathbf{b}^{(2)})}\right) \mathbf{W}^{(2)} \text{diag}\left(\frac{\partial \mathbf{g}(\mathbf{x})}{\partial (\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})}\right) \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \end{aligned}$$

By analysing the chain rule we can derive generic backpropagation formula for the full gradient. Let $\mathbf{v}^{(k)}$ be a k :th layers local field, $\mathbf{v}^{(k)} = \mathbf{W}^{(k)} \mathbf{f}(\mathbf{v}^{(k-1)}) + \mathbf{b}^{(k)}$. Then local gradient matrices $\boldsymbol{\delta}^{(k)}$ are

$$\begin{aligned} \boldsymbol{\delta}^{(L)} &= \text{diag}\left(\frac{\partial \mathbf{f}(\mathbf{v}^{(L)})}{\partial \mathbf{v}^{(L)}}\right) \\ \boldsymbol{\delta}^{(k-1)} &= \boldsymbol{\delta}^{(k)} \mathbf{W}^{(k)} \text{diag}\left(\frac{\partial \mathbf{f}(\mathbf{v}^{(k-1)})}{\partial \mathbf{v}^{(k-1)}}\right) \end{aligned}$$

And network's parameter gradient matrices for each layer are (only j :th element of each row is non-zero):

$$\begin{aligned} \frac{\partial y(\mathbf{x})}{\partial w_{ji}^{(k)}} &= \boldsymbol{\delta}^{(k)} \begin{pmatrix} 0 \\ f(v_i^{(k-1)}) \\ 0 \end{pmatrix} \\ \frac{\partial y(\mathbf{x})}{\partial b_j^{(k)}} &= \boldsymbol{\delta}^{(k)} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \end{aligned}$$

To test that gradient matrix is correctly computed it can be compared with normal squared error calculations (normal backpropagation).

$$\begin{aligned} \varepsilon(\mathbf{x}|\mathbf{w}) &= \frac{1}{2} \|y_i - y(\mathbf{x}|\mathbf{w})\|^2 \\ \frac{\partial \varepsilon(\mathbf{x}|\mathbf{w})}{\partial \mathbf{w}} &= (y(\mathbf{x}|\mathbf{w}) - y_i)^T \frac{\partial y(\mathbf{x}|\mathbf{w})}{\partial \mathbf{w}} \end{aligned}$$

Sometimes also needs gradient with respect to \mathbf{x} and not weights parameters \mathbf{w} . This can be calculated using the chain rule again and is simply (diag() entries are square matrices which diagonal is nonzero):

$$\nabla_{\mathbf{x}} \mathbf{g}(\mathbf{x}|\mathbf{w}) = \text{diag}(\nabla_{\mathbf{v}_L} f(\mathbf{v}_L)) \mathbf{W}_L \dots \text{diag}(\nabla_{\mathbf{v}_2} f(\mathbf{v}_2)) \mathbf{W}_2 \text{diag}(\nabla_{\mathbf{v}_1} f(\mathbf{v}_1)) \mathbf{W}_1$$