

## 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):

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

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

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,  $\mathbf{v}^{(k)} = \mathbf{W}^{(k)} f(\mathbf{v}^{(k-1)}) + \mathbf{b}^{(k)}$ . Then local gradient matrices  $\boldsymbol{\delta}^{(k)}$  are

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

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. For simplicity, let's consider two-layer case initially.

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

The gradient is:

$$\begin{aligned} \frac{\partial \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}} &= \frac{\partial \mathbf{f}(\mathbf{v}^{(2)})}{\partial (\mathbf{W}^{(2)} \mathbf{h} + \mathbf{b}^{(2)})} \frac{\partial (\mathbf{W}^{(2)} \mathbf{h} + \mathbf{b}^{(2)})}{\partial \mathbf{h}} \frac{\partial \mathbf{h}(\mathbf{v}^{(1)})}{\partial (\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})} \frac{\partial (\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})}{\partial \mathbf{x}} \\ \frac{\partial \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}} &= \frac{\partial \mathbf{f}(\mathbf{v}^{(2)})}{\partial \mathbf{v}^{(2)}} \mathbf{W}^{(2)} \frac{\partial \mathbf{h}(\mathbf{v}^{(1)})}{\partial \mathbf{v}^{(1)}} \mathbf{W}^{(1)} \end{aligned}$$

This results into following formula (diag() entries are square matrices which diagonal is nonzero):

$$\nabla_{\mathbf{x}} \mathbf{g}(\mathbf{x}|\mathbf{w}) = \text{diag}(\nabla_{\mathbf{v}^{(L)}} \mathbf{f}(\mathbf{v}^{(L)})) \mathbf{W}^{(L)} \dots \text{diag}(\nabla_{\mathbf{v}^{(2)}} \mathbf{f}(\mathbf{v}^{(2)})) \mathbf{W}^{(2)} \text{diag}(\nabla_{\mathbf{v}^{(1)}} \mathbf{f}(\mathbf{v}^{(1)})) \mathbf{W}^{(1)}$$

### Recurrent Neural Networks and Backpropagation Through Time (BPTT)

The basic learning algorithm for recurrent neural networks (RNN) is **backpropagation through time** (*BPTT*). This is done by unfolding neural net in time and computing the gradients. The recurrent neural network is

$$\mathbf{z}(n+1) = \begin{pmatrix} \mathbf{y}^{(n+1)} \\ \mathbf{r}^{(n+1)} \end{pmatrix} = \mathbf{f}(\mathbf{x}(n), \mathbf{r}(n))$$

The error function to minimize is:

$$E(N) = \frac{1}{2} \sum_{n=0}^N \|\mathbf{d}(n+1) - \mathbf{\Gamma}_y \mathbf{f}(\mathbf{x}(n+1), \mathbf{\Gamma}_r \mathbf{z}(n)|\mathbf{w})\|^2$$

In which  $\mathbf{\Gamma}$  matrices are used to select  $\mathbf{y}(n)$  and  $\mathbf{r}(n)$  vectors from generic output vector and the initial input to feedforward neural network is zero  $\mathbf{z}(0) = \mathbf{0}$ .

It is possible to calculate gradient of  $\mathbf{f}$  using the chain rule

$$\frac{\partial E(N)}{\partial \mathbf{w}} = \sum_{n=0}^N (\mathbf{\Gamma}_y \mathbf{f}(\mathbf{x}(n+1), \mathbf{\Gamma}_r \mathbf{z}(n)|\mathbf{w}) - \mathbf{d}(n+1))^T \mathbf{\Gamma}_y \nabla_{\mathbf{w}} \mathbf{f}(\mathbf{x}(n+1), \mathbf{\Gamma}_r \mathbf{z}(n)|\mathbf{w})$$

To calculate the gradient  $\nabla_{\mathbf{w}} \mathbf{f}$  one must remember that  $\mathbf{z}(n)$  now also depends on  $\mathbf{w}$  resulting into eq:

$$\nabla_{\mathbf{w}} \mathbf{f}(\mathbf{x}(n+1), \mathbf{\Gamma}_r \mathbf{z}(n)) = \frac{\partial \mathbf{f}}{\partial \mathbf{w}} + \frac{\partial \mathbf{f}}{\partial \mathbf{r}} \mathbf{\Gamma}_r \frac{\partial \mathbf{z}(n)}{\partial \mathbf{w}}$$

To further compute gradients we get a generic update rule

$$\frac{\partial \mathbf{z}(n)}{\partial \mathbf{w}} = \frac{\partial \mathbf{f}}{\partial \mathbf{w}} + \frac{\partial \mathbf{f}}{\partial \mathbf{r}} \mathbf{\Gamma}_r \frac{\partial \mathbf{z}(n-1)}{\partial \mathbf{w}}$$

The computation of gradients can be therefore bootstrapped by setting  $\frac{\partial \mathbf{z}(1)}{\partial \mathbf{w}} = \frac{\partial \mathbf{f}(\mathbf{0}, \mathbf{0})}{\partial \mathbf{w}}$  and iteratively updating  $\mathbf{z}$  gradient while computing current error for the timestep.