

# THE INCREDIBLE SHRINKING NEURAL NETWORK: PRUNING TO OPERATE IN CONSTRAINED MEMORY ENVIRONMENTS

**Nikolas Wolfe, Aditya Sharma, Lukas Drude & Bhiksha Raj**

School of Computer Science

Carnegie Mellon University

Pittsburgh, PA 15213, USA

{nwolfe, adityasharma, bhiksha}@cmu.edu

## ABSTRACT

We propose and evaluate a method for pruning neural networks to operate in constrained memory environments such as mobile or embedded devices. We evaluate a simple pruning technique using first-order derivative approximations of the gradient of each neuron in an optimally trained network, and turning off those neurons which contribute least to the output of the network. We then show the limitations of this type of approximation by comparing against the ground truth value for the change in error resulting from the removal of a given neuron. We attempt to improve on this using a second-order derivative approximation. We also explore the correlation between neurons in a trained network and attempt to improve our choice of candidate neurons for removal to account for faults that can occur from the removal of a single neuron at a time. We argue that this method of pruning allows for the optimal tradeoff in network size versus accuracy in order to operate within the memory constraints of a particular device or application environment.

## 1 SECOND DERIVATIVE GRADIENT TERMS

Given a neuron  $n$  with output  $a_i$ , and outgoing weights  $[w_{i,1}, w_{i,2}, \dots, w_{i,j}]$ , the input  $x$  to each of its  $j$  forward-connected neurons is given by:

$$x_j = \sum_i (w_{i,j} \cdot a_i) \quad (1)$$

For simplicity's sake, we will drop the index variable  $i$  and examine only the connection between the output  $a$  of neuron  $n$  and each of its  $j$  connections to the next forward layer. So the contribution  $c_j$  from neuron  $n$  to the input of each forward-connected neuron is given by:

$$c_j = w_j \cdot a \quad (2)$$

Where  $w_j$  is the weight connecting the output  $a$  from  $n$  to the input of the  $j$ th neuron. We will denote the error function of an optimally trained network as  $E$ . The second-derivative of  $E$  with respect to the output of neuron  $n$  is given by:

$$\frac{d^2 E}{da^2} = \frac{d}{da} \frac{dE}{da} = \frac{d}{da} \sum_j \left( \frac{dE}{dc_j} \cdot \frac{dc_j}{da} \right) = \sum_j \left( \frac{d}{da} \frac{dE}{dc_j} \cdot \frac{dc_j}{da} \right) \quad (3)$$

Which states that the 2nd derivative of  $E$  with respect to  $a$  is the sum of the 2nd derivative terms of all outgoing connections.

### 1.1 SECOND DERIVATIVE BACK-PROPAGATION

Name and network definitions:

$$E = \frac{1}{2} \sum_i (o_i^{(0)} - t_i)^2 \quad o_i^{(m)} = \sigma(x_i^{(m)}) \quad x_i^{(m)} = \sum_j w_{ji}^{(m)} o_j^{(m+1)} \quad c_{ji}^{(m)} = w_{ji}^{(m)} o_j^{(m+1)} \quad (4)$$

Superscripts represent the index of the layer of the network in question, with 0 representing the output layer.  $E$  is the squared-error network cost function.  $o_i^{(m)}$  is the  $i$ th output in layer  $m$  generated by the activation function  $\sigma$ , which in this paper is the standard logistic sigmoid.  $x_i^{(m)}$  is the weighted sum of inputs to the  $i$ th neuron in the  $m$ th layer, and  $c_{ji}^{(m)}$  is the contribution of the  $j$ th neuron in the  $m + 1$  layer to the input of the  $i$ th neuron in the  $m$ th layer.

## 2 FIRST AND SECOND DERIVATIVES

The first and second derivatives of the cost function with respect to the output:

$$\frac{\partial E}{\partial o_i^{(0)}} = o_i^{(0)} - t_i \quad (5)$$

$$\frac{\partial^2 E}{\partial (o_i^{(0)})^2} = 1 \quad (6)$$

The first and second derivatives of the sigmoid function in forms depending only on the output:

$$\sigma' = \sigma (1 - \sigma) \quad (7)$$

$$\sigma'' = \sigma' (1 - 2\sigma) \quad (8)$$

And for future convenience:

$$\frac{\partial o_i^{(m)}}{\partial x_i^{(m)}} = \frac{\partial}{\partial x_i^{(m)}} \left( o_i^{(m)} = \sigma(x_i^{(m)}) \right) = \left( o_i^{(m)} \right) \left( 1 - o_i^{(m)} \right) \quad (9)$$

$$\frac{\partial^2 o_i^{(m)}}{\partial x_i^{(m)2}} \quad (10)$$

Derivative of the error with respect to the  $i$ th neuron's input  $x_i^{(0)}$  in the output layer:

$$\frac{\partial E}{\partial x_i^{(0)}} = \frac{\partial E}{\partial o_i^{(0)}} \frac{\partial o_i^{(0)}}{\partial x_i^{(0)}} \quad (11)$$

$$= \underbrace{\left( o_i^{(0)} - t_i \right)}_{\text{from (5)}} \underbrace{\left( \sigma \left( x_i^{(0)} \right) \left( 1 - \sigma \left( x_i^{(0)} \right) \right) \right)}_{\text{from (7)}} \quad (12)$$

$$= \left( o_i^{(0)} - t_i \right) \left( o_i^{(0)} \left( 1 - o_i^{(0)} \right) \right) \quad (13)$$

$$= \left( o_i^{(0)} - t_i \right) (\sigma') \quad (14)$$

Second derivative of the error with respect to the  $i$ th neuron's input  $x_i^{(0)}$  in the output layer:

$$\frac{\partial^2 E}{\partial (x_i^{(0)})^2} = \frac{\partial}{\partial x_i^{(0)}} \left( \frac{\partial E}{\partial o_i^{(0)}} \frac{\partial o_i^{(0)}}{\partial x_i^{(0)}} \right) \quad (15)$$

$$= \frac{\partial^2 E}{\partial x_i^{(0)} \partial o_i^{(0)}} \frac{\partial o_i^{(0)}}{\partial x_i^{(0)}} + \frac{\partial E}{\partial o_i^{(0)}} \frac{\partial^2 o_i^{(0)}}{\partial x_i^{(0)2}} \quad (16)$$

$$= \frac{\partial^2 E}{\partial x_i^{(0)} \partial o_i^{(0)}} \underbrace{\left( o_i^{(0)} \left( 1 - o_i^{(0)} \right) \right)}_{\text{from (7)}} + \underbrace{\left( o_i^{(0)} - t_i \right)}_{\text{from (5)}} \underbrace{\left( o_i^{(0)} \left( 1 - o_i^{(0)} \right) \right)}_{\text{from 8}} \left( 1 - 2o_i^{(0)} \right) \quad (17)$$

$$\left( \frac{\partial^2 E}{\partial x_i^{(0)} \partial o_i^{(0)}} \right) = \frac{\partial}{\partial x_i^{(0)}} \frac{\partial E}{\partial o_i^{(0)}} = \frac{\partial}{\partial x_i^{(0)}} \underbrace{\left( o_i^{(0)} - t_i \right)}_{\text{from (5)}} = \frac{\partial o_i^{(0)}}{\partial x_i^{(0)}} = \underbrace{\left( o_i^{(0)} \left( 1 - o_i^{(0)} \right) \right)}_{\text{from (7)}} \quad (18)$$

$$\frac{\partial^2 E}{\partial (x_i^{(0)})^2} = \left( o_i^{(0)} \left( 1 - o_i^{(0)} \right) \right)^2 + \left( o_i^{(0)} - t_i \right) \left( o_i^{(0)} \left( 1 - o_i^{(0)} \right) \right) \left( 1 - 2o_i^{(0)} \right) \quad (19)$$

$$= (\sigma')^2 + \left( o_i^{(0)} - t_i \right) (\sigma'') \quad (20)$$

First derivative of the error with respect to a single input contribution  $c_{ji}^{(0)}$  from neuron  $j$  to neuron  $i$  with weight  $w_{ji}^{(0)}$  in the output layer:

$$\frac{\partial E}{\partial c_{ji}^{(0)}} = \frac{\partial E}{\partial o_i^{(0)}} \frac{\partial o_i^{(0)}}{\partial x_i^{(0)}} \frac{\partial x_i^{(0)}}{\partial c_{ji}^{(0)}} \quad (21)$$

$$= \underbrace{\left(o_i^{(0)} - t_i\right)}_{\text{from (5)}} \underbrace{\left(o_i^{(0)} \left(1 - o_i^{(0)}\right)\right)}_{\text{from (7)}} \frac{\partial x_i^{(0)}}{\partial c_{ji}^{(0)}} \quad (22)$$

$$\left(\frac{\partial x_i^{(m)}}{\partial c_{ji}^{(m)}}\right) = \frac{\partial}{\partial c_{ji}^{(m)}} \left(x_i^{(m)} = \sum_j w_{ji}^{(m)} o_j^{(m+1)}\right) = \frac{\partial}{\partial c_{ji}^{(m)}} \left(c_{ji}^{(m)} + k\right) = 1 \quad (23)$$

$$\frac{\partial E}{\partial c_{ji}^{(0)}} = \left(o_i^{(0)} - t_i\right) \left(o_i^{(0)} \left(1 - o_i^{(0)}\right)\right) \quad (24)$$

Second derivative of the error with respect to a single input contribution  $c_{ji}^{(0)}$ :

$$\frac{\partial^2 E}{\partial c_{ji}^{(0)^2}} = \frac{\partial^2 E}{\partial c_{ji}^{(0)} \partial o_i^{(0)}} \frac{\partial o_i^{(0)}}{\partial x_i^{(0)}} \underbrace{\frac{\partial x_i^{(0)}}{\partial c_{ji}^{(0)}}}_{(24)} + \underbrace{\frac{\partial E}{\partial o_i^{(0)}}}_{(5)} \frac{\partial^2 o_i^{(0)}}{\partial c_{ji}^{(0)} \partial x_i^{(0)}} \underbrace{\frac{\partial x_i^{(0)}}{\partial c_{ji}^{(0)}}}_{(24)} + \underbrace{\frac{\partial E}{\partial o_i^{(0)}} \frac{\partial o_i^{(0)}}{\partial x_i^{(0)}} \frac{\partial^2 x_i^{(0)}}{\partial c_{ji}^{(0)^2}}}_{\frac{\partial x_i^{(0)}}{\partial c_{ji}^{(0)}}=1, \text{ so } \frac{\partial^2 x_i^{(0)}}{\partial c_{ji}^{(0)^2}}=0} \quad (25)$$

$$= \left(o_i^{(0)} - t_i\right) \quad (26)$$

## 2.1 RETRIEVAL OF STYLE FILES

The style files for ICLR and other conference information are available on the World Wide Web at

<http://www.iclr.cc/>

The file `iclr2016_conference.pdf` contains these instructions and illustrates the various formatting requirements your ICLR paper must satisfy. Submissions must be made using  $\text{\LaTeX}$  and the style files `iclr2016_conference.sty` and `iclr2016_conference.bst` (to be used with  $\text{\LaTeX}2\epsilon$ ). The file `iclr2016_conference.tex` may be used as a “shell” for writing your paper. All you have to do is replace the author, title, abstract, and text of the paper with your own.

The formatting instructions contained in these style files are summarized in sections ??, ??, and ?? below.