

Backpropagation on the exact gradient as solution to 1st order PDE.

March 3, 2019

1 Normal Backpropagation:

Through out this document, we follow the convention that indices that appear on one side of the equation and not on the other are summed over. To make it convinient later on, we split the usual neural network function into two steps (the linear transformation step and the non-linear activation step):

$$\begin{aligned} O_{by}^K &= W_{yx}^K I_{bx}^K + B_y^K \\ \text{For : } K &\geq 1 \\ \text{We have : } I_{bx}^K &\rightarrow O_{bx}^{K-1} \\ O_{by}^{K+1} &= \sigma [O_{by}^K] \\ \sigma &\Rightarrow \text{Activation Function.} \end{aligned}$$

To make this more friendly for matrix multiplication, the linear part can be written as:

$$O_{by}^K = I_{bx}^K (W^K)^T_{xy} + B_y^K$$

To learn features by backpropagation (denoting the last layer by the index L):

$$\begin{aligned} \epsilon &\equiv (O_{by}^L - \mathcal{Y}_{by})^2 \\ \delta\epsilon &= 2 (O_{by}^L - \mathcal{Y}_{by}) (\delta O_{by}^L) \end{aligned}$$

Let us denote:

$$\Delta_{by}^L \equiv 2 (O_{by}^L - \mathcal{Y}_{by})$$

Hence, we get:

$$\delta\epsilon = \Delta_{by}^L (\delta O_{by}^L)$$

Now, consider the first (linear) cases:

$$\begin{aligned} \Delta_{by}^K \delta O_{by}^K &= (\delta W_{yx}^K) \Delta_{by}^K I_{bx}^K + \Delta_{by}^K W_{yx}^K (\delta I_{bx}^K) + \Delta_{by}^K (\delta B_y^K) \\ \Delta_{by}^K \delta O_{by}^K &= (\delta W_{yx}^K) \Delta_{by}^K O_{bx}^{K-1} + \Delta_{by}^K W_{yx}^K (\delta O_{bx}^{K-1}) + \Delta_{by}^K (\delta B_y^K) \end{aligned}$$

This gives the backpropagation recursion relation:

$$\Delta_{bx}^{K-1} = \Delta_{by}^K W_{yx}^K$$

$$\text{Hence : } \Delta_{by}^K (\delta O_{by}^K) \rightarrow \Delta_{bx}^{K-1} (\delta O_{bx}^{K-1})$$

And also the gradient of W_{yx}^K and B_y^K which will be used for the gradient descent:

$$\delta W_{yx}^K = \Delta_{by}^K O_{bx}^{K-1} = (\Delta^K)_{yb}^T O_{bx}^{K-1}$$

$$\delta B_y^K = \Delta_{by}^K$$

Using the sum convention mentioned at the beginning, we can interpret the second equation as:

$$\delta B_y^K = \sum_b \Delta_{by}^K$$

Similar relations can be derived for the activation layer:

$$\Delta_{by}^K (\delta O_{by}^K) = (\sigma' [O_{by}^{K-1}] \Delta_{by}^K) (\delta O_{by}^{K-1})$$

$$\text{Hence : } \Delta_{by}^{K-1} = \Delta_{by}^K \sigma' [O_{by}^{K-1}]$$

We now have all the recursion relations for derivatives (backpropagation) required to perform gradient descent.

2 Backpropagation on the gradient:

We now show that the above procedure can also be used to learn on the gradient of the neural network instead of the plain output of the network.

$$\begin{aligned}
O_{by}^K &= W_{yx}^K I_{bx}^K + B_y^K \\
\left(\frac{dO_{by}^K}{dI_{bm}^0} \right) &= W_{yx}^K \left(\frac{dI_{bx}^K}{dI_{bm}^0} \right) \\
\left(\frac{dO_{by}^K}{dI_{bm}^0} \right) &= W_{yx}^K \left(\frac{dO_{bx}^{K-1}}{dI_{bm}^0} \right) \\
\mathfrak{D}_{bym}^K &\equiv \left(\frac{dO_{by}^K}{dI_{bm}^0} \right) \\
\mathfrak{D}_{bym}^K &= W_{yx}^K \mathfrak{D}_{bxm}^{K-1}
\end{aligned}$$

$$\begin{aligned}
\text{if: } \mathfrak{D}_{bxm}^{K-1} &= \delta_{xm} \quad (\text{first layer}) \\
\mathfrak{D}_{bym}^K &= W_{yx}^K \delta_{xm} \\
\mathfrak{D}_{bym}^K &= W_{ym}^K
\end{aligned}$$

$$\begin{aligned}
O_{bi}^K &= \sigma(I_{bi}^K) \\
\frac{dO_{bi}^K}{dI_{bm}^0} &= \sigma'(I_{bi}^K) \frac{dI_{bi}^K}{dI_{bm}^0} \\
\left(\frac{dO_{bi}^K}{dI_{bm}^0} \right) &= \sigma'(I_{bi}^K) \left(\frac{dO_{bi}^{K-1}}{dI_{bm}^0} \right) \\
\mathfrak{D}_{bim}^K &= \sigma'(O_{bi}^{K-1}) \mathfrak{D}_{bim}^{K-1} \\
(\mathfrak{D}_{bim}^K - \mathcal{A}_{bim})^2 &= \epsilon \\
2(\mathfrak{D}_{bim}^K - \mathcal{A}_{bim}) \partial \mathfrak{D}_{bim}^K &= \partial \epsilon
\end{aligned}$$

So:

$$(\mathfrak{D}_{bim}^K - \mathcal{A}_{bim}) \equiv \Delta_{bim}^K$$

$$\begin{aligned}
\mathfrak{D}_{bym}^K &= W_{yx}^K \mathfrak{D}_{bxm}^{K-1} \\
\Delta_{bym}^K \partial \mathfrak{D}_{bym}^K &= W_{yx}^K \Delta_{bym}^K \partial \mathfrak{D}_{bxm}^{K-1} \\
\Rightarrow \Delta_{bxm}^{K-1} &= W_{yx}^K \Delta_{bym}^K \\
\Rightarrow (\Delta_b^{K-1})_{xm} &= (W^K)_{xy}^T (\Delta_b^K)_{ym} \\
(\Delta_b^{K-1})_{xm} &= \left[(W^K)^T \Delta_b^K \right]_{xm}
\end{aligned}$$

For evaluating W :

$$\begin{aligned}
\Delta_{bym}^K \partial \mathfrak{D}_{bym}^K &= (\Delta_{bym}^K \partial) [W_{yx}^K \mathfrak{D}_{bxm}^{K-1}] \\
&= \Delta_{bym}^K \mathfrak{D}_{bxm}^{K-1} (\partial W_{yx}^K) \\
\delta W_{yx}^K &= \Delta_{bym}^K \mathfrak{D}_{bxm}^{K-1} \\
\delta W_{yx}^K &= (\Delta_b^K)_{ym} (\mathfrak{D}_b^{K-1})_{xm} \\
\delta W_{yx}^K &= (\Delta_b^K)_{ym} (\mathfrak{D}_b^{K-1})_{mx}^T \\
\delta W_{yx}^K &= \left[(\Delta_b^K) (\mathfrak{D}_b^{K-1})^T \right]_{yx} \\
\text{if : } (\mathfrak{D}_b^{K-1})_{xm} &= \delta_{xm} \\
\delta W_{yx}^K &= (\Delta_b^K)_{ym} (\delta_{xm}) \\
\delta W_{yx}^K &= \sum_b (\Delta_b^K)_{yx} \\
\delta W_{yx}^K &= (\Delta_b^K)_{ym} (\mathfrak{D}_b^{K-1})_{mx}^T
\end{aligned}$$

For the activation layers:

$$\begin{aligned}
\mathfrak{D}_{bim}^K &= \sigma' (O_{bi}^{K-1}) \mathfrak{D}_{bim}^{K-1} \\
\partial \mathfrak{D}_{bim}^K &= [\partial \sigma' (O_{bi}^{K-1})] \mathfrak{D}_{bim}^{K-1} + \sigma' (O_{bi}^{K-1}) [\partial \mathfrak{D}_{bim}^{K-1}] \\
\partial \mathfrak{D}_{bim}^K &= [\partial O_{bi}^{K-1}] \mathfrak{D}_{bim}^{K-1} \sigma'' (O_{bi}^{K-1}) + \sigma' (O_{bi}^{K-1}) [\partial \mathfrak{D}_{bim}^{K-1}] \\
\Delta_{bim}^K \partial \mathfrak{D}_{bim}^K &= \{ \Delta_{bim}^K \mathfrak{D}_{bim}^{K-1} \sigma'' (O_{bi}^{K-1}) \} [\partial O_{bi}^{K-1}] + \{ \Delta_{bim}^K \sigma' (O_{bi}^{K-1}) \} [\partial \mathfrak{D}_{bim}^{K-1}] \\
\Delta_{bim}^K \partial \mathfrak{D}_{bim}^K &= (\delta_{bi}^{K-1}) [\partial O_{bi}^{K-1}] + (\Delta_{bim}^{K-1}) [\partial \mathfrak{D}_{bim}^{K-1}] \\
(\delta_{bi}^{K-1}) &\equiv \{ \Delta_{bim}^K \mathfrak{D}_{bim}^{K-1} \sigma'' (O_{bi}^{K-1}) \} \\
(\Delta_{bim}^{K-1}) &\equiv \{ \Delta_{bim}^K \sigma' (O_{bi}^{K-1}) \}
\end{aligned}$$

Where the first term represents a normal back propagation while the second term represents a gradient backpropagation.

3 Solving a toy example

Using this, we are now in a position to solve simple partial differential equations. For benchmarking purposes, consider:

$$\frac{dy}{dx} = \frac{1}{x}$$

The solution to this equation is: $y = \log(x) + c$. Solving this using the gradient backpropagation, we get the remarkable graph [Figure 1](#)

An implementation of this method (in C++) is available at: https://github.com/aravindhv10/CPP_Wrappers/tree/master/LatestHeaders. The implementation so far only works on the CPU (uses CPU blas and doesnot run on GPU), this result was obtained with just a few minutes of training (ofcourse, the differential equation is also very simple).

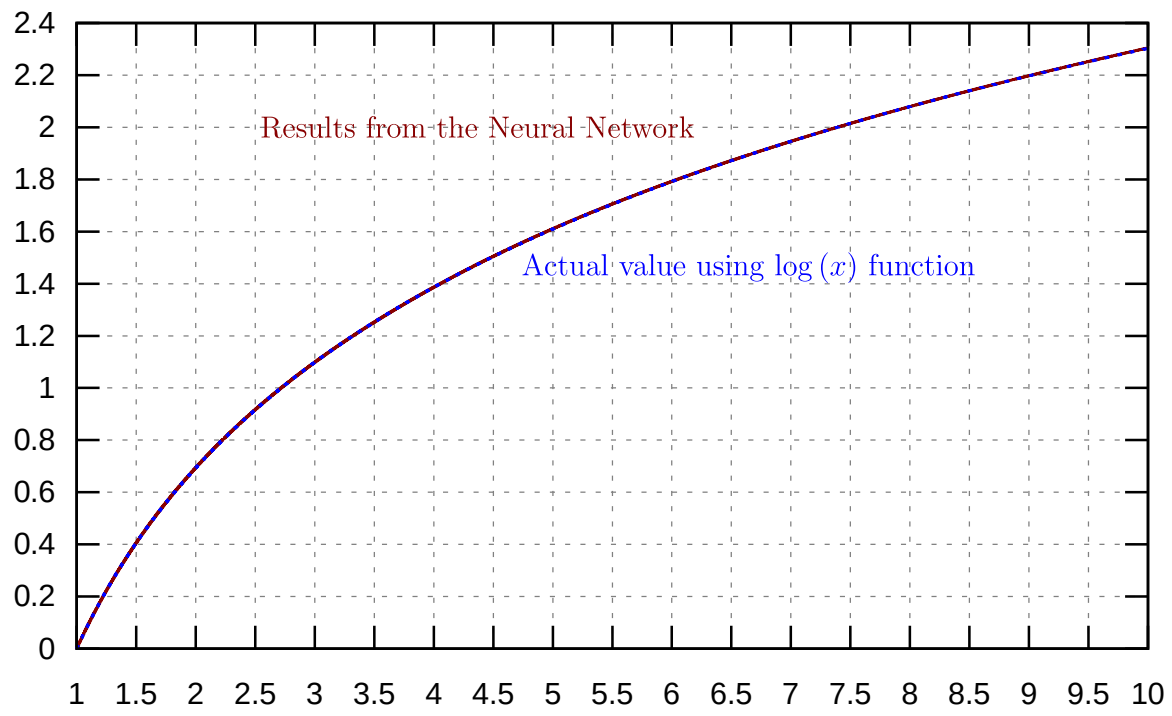


Figure 1: Comparing solution obtained from Neural Network to the exact solution for the simple differential equation mentioned above.