

1

Fully Connected NN Model

1.1 Derivatives with respect to model parameters and input

Consider a fully connected neural network (FNN) model, with $k - 1$ hidden layers, with input Z and output Y of the form

$$\begin{aligned} W_1 &= \mathbf{A}_1 Z + \mathcal{B}_1 \\ Z_1 &= f_1(W_1) \\ &\vdots \\ W_{k-1} &= \mathbf{A}_{k-1} Z_{k-2} + \mathcal{B}_{k-1} \\ Z_{k-1} &= f_{k-1}(W_{k-1}) \\ W_k &= \mathbf{A}_k Z_{k-1} + \mathcal{B}_k \\ Y &= f_k(W_k), \end{aligned} \tag{1.1}$$

where $f_i(W), i = 1, \dots, k$ are scalar functions applied to each elements of the vector W . The parameters of the network are $\theta = \{\mathbf{A}_1, \mathcal{B}_1, \dots, \mathbf{A}_k, \mathcal{B}_k\}$ and the goal is to compute the derivative with respect to each elements of these matrices and vectors and, also, with respect to the input Z . In fact as we will see we can group the derivatives into matrices for each level separately and compute them recursively.

Let us start our derivations. For (stochastic) gradient type algorithms we usually have some form of scalar cost function $\phi(Y)$ that we like to compute its gradient with respect to the network parameters and input. Therefore, starting from the output and going backwards, we can write

$$\nabla_{[\mathbf{A}_k \mathcal{B}_k]} \phi(Y) = (\nabla_Y \phi(Y) \odot f'_k(W_k)) [Z_{k-1}^\top \ 1] \tag{1.2}$$

$$\nabla_{Z_{k-1}} \phi(Y) = \mathbf{A}_k^\top (\nabla_Y \phi(Y) \odot f'_k(W_k)), \tag{1.3}$$

where $f'(W)$ denotes the derivative of the scalar function $f(\cdot)$ applied to each element of the vector W and \odot denotes element-by-element multiplication. Let us now go one step deeper into the network. We then have

$$\nabla_{[\mathbf{A}_{k-1} \mathcal{B}_{k-1}]_{ij}} \phi(Y) = (\nabla_{Z_{k-1}} \phi(Y))^\top \frac{\partial Z_{k-1}}{\partial [\mathbf{A}_{k-1} \mathcal{B}_{k-1}]_{ij}},$$

from which we conclude that

$$\nabla_{[\mathbf{A}_{k-1}\mathbf{B}_{k-1}]} \phi(Y) = (\nabla_{Z_{k-1}} \phi(Y) \odot f'_{k-1}(W_{k-1})) [Z_{k-2}^\top \ 1]. \quad (1.4)$$

Similarly

$$\nabla_{[Z_{k-2}]_i} \phi(Y) = (\nabla_{Z_{k-1}} \phi(Y))^\top \frac{\partial Z_{k-1}}{\partial [Z_{k-2}]_i},$$

therefore

$$\nabla_{Z_{k-2}} \phi(Y) = \mathbf{A}_{k-2}^\top (\nabla_{Z_{k-1}} \phi(Y) \odot f'_{k-1}(W_{k-1})). \quad (1.5)$$

We can now deduce the following recursion formula (adopting the engineering induction). For $i = k, k-1, \dots, 1$ define

$$\begin{aligned} \mathcal{V}_i &= \mathcal{U}_i \odot f'_i(W_i), \text{ initializing with } \mathcal{U}_k = \nabla_Y \phi(Y), \\ \mathcal{U}_{i-1} &= \mathbf{A}_i^\top \mathcal{V}_i, \end{aligned} \quad (1.6)$$

which generates $\{\mathcal{V}_k, \dots, \mathcal{V}_1\}$ and $\{\mathcal{U}_{k-1}, \dots, \mathcal{U}_0\}$. Then

$$\begin{aligned} \nabla_{[\mathbf{A}_i\mathbf{B}_i]} \phi(Y) &= \mathcal{V}_i [Z_{i-1}^\top \ 1] \\ \nabla_Z \phi(Y) &= \mathcal{U}_0. \end{aligned} \quad (1.7)$$

1.2 Update of Network Parameters

If for example our goal is to select the network parameters to *minimize* the expectation of $\phi(Y)$ then we can apply a stochastic gradient descent of the form

$$[\mathbf{A}_k\mathbf{B}_k]_t = [\mathbf{A}_k\mathbf{B}_k]_{t-1} - \mu \nabla_{[\mathbf{A}_k\mathbf{B}_k]_{t-1}} \phi(Y_t).$$

However it has been observed in practice that normalizing the derivatives produces more stable solutions. In particular we can use the ADAM scheme (Kingma, D. P., Ba, J. L. (2015). Adam: a Method for Stochastic Optimization. International Conference on Learning Representations, 1–13) which consists in computing an estimate of the power of each derivative element

$$\mathbf{P}_{[\mathbf{A}_k\mathbf{B}_k]}(t) = (1 - \lambda) \mathbf{P}_{[\mathbf{A}_k\mathbf{B}_k]}(t-1) + \lambda (\nabla_{[\mathbf{A}_k\mathbf{B}_k]_{t-1}} \phi(Y_t))^{(.2)}$$

where for a matrix \mathbf{Q} we denote with $\mathbf{Q}^{(.2)}$ the element-by-element raise to the power 2 and $\lambda \ll 1$. We initialize with

$$\mathbf{P}_{[\mathbf{A}_k\mathbf{B}_k]}(1) = (\nabla_{[\mathbf{A}_k\mathbf{B}_k]_0} \phi(Y_1))^{(.2)}.$$

The power estimate must be updated after we compute the gradient $\nabla_{[\mathbf{A}_k\mathbf{B}_k]_{t-1}} \phi(Y_t)$.

When updating the network parameters, we can normalize each gradient element using the corresponding power estimate as follows

$$[\mathbf{A}_k\mathbf{B}_k]_t = [\mathbf{A}_k\mathbf{B}_k]_{t-1} - \mu \nabla_{[\mathbf{A}_k\mathbf{B}_k]_{t-1}} \phi(Y_t) ./ \sqrt{c + \mathbf{P}_{[\mathbf{A}_k\mathbf{B}_k]}(t)}$$

where $./$ denotes element-by-element division, c is a small number to avoid division by 0, and the square-root ($\sqrt{}$) is again applied element-by element onto the matrix $c + \mathbf{P}_{[\mathbf{A}_k\mathbf{B}_k]}(t)$. This way we normalize each derivative and make the elements of the same order of magnitude which then allows for the use of a single step size μ for all adaptations. Otherwise, if we do not normalize, and use a single μ then we experience completely different convergence rates per parameter element!