

How to vectorize the equations that we have derived so far so that we can exploit parallelism available due to GPU?

Instead of passing one example through the network we want to pass one batch of examples so that we get prediction for everything at same time and we can propagate all the gradients back at the same time. Instead of computing gradient wrt one example we are going to do vectorization for m examples.

We will see what changes in equation we have to make to make this work?

For 1 Example	For m Examples
<p>Error in output layer $\rightarrow \delta^L = (a^L - y^L)$</p> <p>hidden layer $\rightarrow \delta^L = (w^{L+1} \delta^{L+1}) \odot \sigma'(z^L)$</p> <p>updation of weights and biases $\left\{ \begin{array}{l} \frac{\partial L}{\partial b^L} = \delta^L \\ \frac{\partial L}{\partial w^L} = a^{L-1} (\delta^L)^T \end{array} \right.$</p>	<p>$\delta^L = a^L - y^L$</p> <p>$\delta^L = (\delta^{L+1}, w^{L+1^T}) \odot \sigma'(z^L)$</p> <p>$\frac{\partial L}{\partial b^L} = \frac{1}{m} \text{np.sum}(\delta^L, \text{axis}=0)$</p> <p>$\frac{\partial L}{\partial w^L} = a^{L-1^T} \cdot \delta^L$</p>

In 1 example,

Each example is a column vector

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \end{bmatrix}_{(n,1)}$$

$$X = \begin{bmatrix} \text{---} & x^{(1)} & \text{---} \\ \text{---} & x^{(2)} & \text{---} \\ & \vdots & \\ \text{---} & x^{(m)} & \text{---} \end{bmatrix}_{m \times n}$$

Error in output layer:

activation vector
(if you have C units in last layer)

$$a_L = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_C \end{bmatrix}$$

Then you take argmax over these to find out which class the input belongs to

Everything which was a column vector earlier will now become a row vector.

m examples and each example has n features

activation values for 1st example

$$a_L = \begin{bmatrix} a^{(1)} \\ a^{(2)} \\ \vdots \\ a^{(i)} \\ \vdots \\ a^{(m)} \end{bmatrix}_{m \times C}$$

m = no. of example
 C = no. of output classes/
 no. of units in output layer

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{matrix} \leftarrow \text{Cat} \\ \leftarrow \text{Dog} \\ \leftarrow \text{Horse} \end{matrix}$$

$$y = \begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}_{m \times C}$$

We can easily do $a_L - y$ which will give $m \times C$ matrix.

δ^L will be $m \times C$ matrix

Update of biases:

$$b^L = \begin{bmatrix} b_1^L \\ b_2^L \\ b_3^L \end{bmatrix} \quad \delta^L = \begin{bmatrix} \delta b_1 \\ \delta b_2 \\ \delta b_3 \end{bmatrix}$$

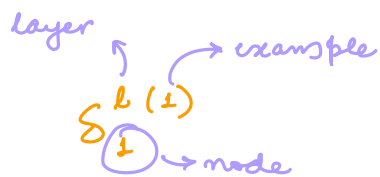
$$b^L = b^L - \eta \delta^L$$

Each nodes has bias b_1, b_2 and b_3

Error w.r.t 1st example

$$\begin{bmatrix} \delta^{L(1)} \\ \delta^{L(2)} \\ \vdots \\ \delta^{L(m)} \end{bmatrix}$$

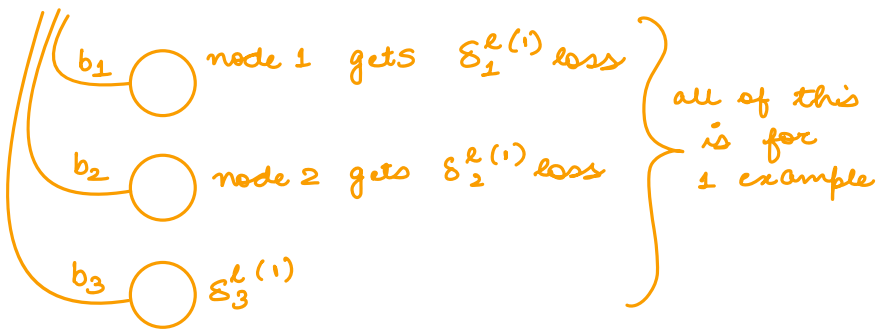
Error w.r.t m^{th} example



when you apply gradient descent and you want to update b_1 :

$$b_1^L = b_1^L - \eta \sum_{i=1}^m \frac{\partial L}{\partial b_1^{L(i)}}$$

we want to keep hidden unit (subscript) same and iterate over all examples



If there are more examples then you will have $\delta_1^{L(2)}, \delta_2^{L(2)}, \dots$

$$\begin{bmatrix} \delta_1^{L(1)} & \delta_2^{L(1)} & \delta_3^{L(1)} & \dots \\ \delta_1^{L(2)} & \delta_2^{L(2)} & \delta_3^{L(2)} & \dots \\ \delta_1^{L(3)} & & & \\ \vdots & & & \\ \delta_1^{L(n)} & & & \end{bmatrix}$$

when you are going to update b_1 then do sum along this column.

this sum will update b_1 bias.

$$\frac{\partial L}{\partial b^L} = \frac{1}{m} \text{np.sum}(\delta^L, \text{axis}=0)$$

In matrix notation we will use this formula



we are finding average gradient along bias

update of weights

w^l will have a shape

of h_{l-1}, h_l

hidden unit
in previous layer

hidden unit
in current layer

$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix}_{(h_{l-1}, h_l)} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}_{(h_{l-1}, 1)} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}_{(1, h_l)}$$

$$\frac{\partial L}{\partial w^l} = a^{l-1} \boxed{(s^l)^T}$$

$$\frac{\partial L}{\partial w^l} = \boxed{a^{l-1}{}^T} s^l$$

$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix}_{(h^{l-1}, h^l)} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}_{(h^{l-1}, m)} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}_{(m, h^l)}$$