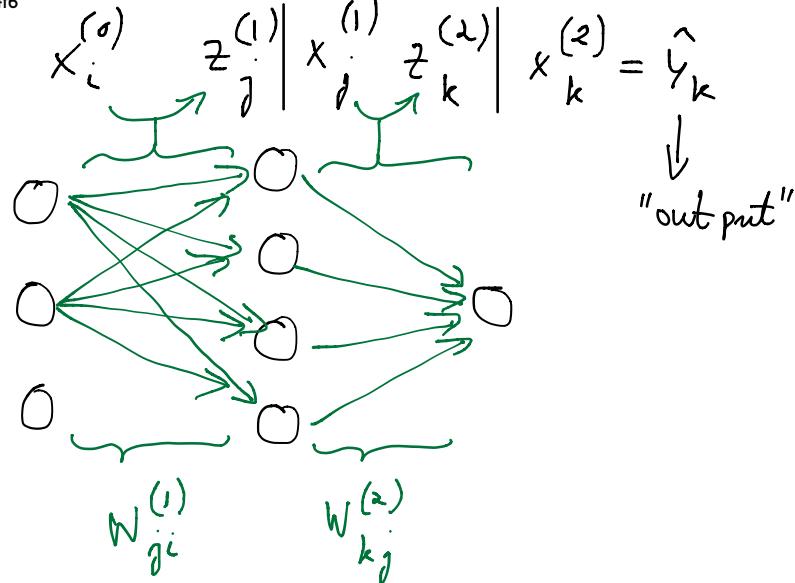


## 2-layer neural net with backpropagation

<https://towardsdatascience.com/how-to-build-your-own-neural-network-from-scratch-in-python-68998a08e4f6>



$m$  ↓

X1	X2	X3	Y
0	0	1	0
0	1	1	1
1	0	1	1
1	1	1	0

↑

$(x_{m,i}, y_{m,k})$  goal:  $\hat{y}_{m,k} = y_{m,k}$   
 for each example  $m$ : for all  $k$  and  
 $x_i^{(0)} = x_{m,i}$  "input" all examples  $m$

$$x_i^{(0)} \rightarrow z_j^{(1)} = W_{ji}^{(1)} x_i^{(0)} + b_j^{(1)}$$

$$\rightarrow x_j^{(1)} = \sigma(z_j^{(1)})$$

$$\rightarrow z_k^{(2)} = W_{kj}^{(2)} x_j^{(1)} + b_k^{(2)}$$

$$\rightarrow \hat{y}_k = x_k^{(2)} = \sigma(z_k^{(2)})$$

define loss function  $\Delta \mathcal{L}$ :  
 "mean squared error"

$$\Delta \mathcal{L} = \sum_i (y_i - \hat{y}_i)^2$$

$$\hat{y} = x^{(2)} = \sigma(W^{(2)} \cdot \sigma(W^{(1)} \cdot x^{(0)} + b^{(1)}) + b^{(2)})$$

back propagate through last layer  
 - compute gradient of  $\Delta L$  w.r.t.  
 $w^{(2)}$ :

$$\Rightarrow \frac{\partial \Delta L}{\partial w_{mn}^{(2)}} = \frac{\partial \Delta L}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial z_k^{(2)}} \cdot \frac{\partial z_k^{(2)}}{\partial w_{mn}^{(2)}}$$

$$= 2 \cdot (\hat{y}_k - y_k) \cdot \sigma' \Big|_{z_k^{(2)}} \cdot \delta_{mk} \delta_{nj} \cdot x_j^{(1)}$$

$$= 2 \cdot (\hat{y}_m - y_m) \cdot \sigma' \Big|_{z_m^{(2)}} \cdot x_n^{(1)}$$

$$= 2 \cdot (\hat{y}_m - y_m) \cdot \sigma(z_m^{(2)}) \cdot (1 - \sigma(z_m^{(2)}))$$

$$\quad \quad \quad \cdot x_n^{(1)}$$

now propagate further back  
 one layer:

$$\frac{\partial \Delta L}{\partial w_{mn}^{(1)}} = \frac{\partial \Delta L}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial z_k^{(2)}} \cdot \frac{\partial z_k^{(2)}}{\partial x_j^{(1)}} \cdot \frac{\partial x_j^{(1)}}{\partial z_j^{(1)}} \cdot \frac{\partial z_j^{(1)}}{\partial w_{mn}^{(1)}}$$

$$= 2 \cdot (\hat{y}_k - y_k) \cdot \sigma' \Big|_{z_k^{(2)}} \cdot W_{kj}^{(2)} \cdot \sigma' \Big|_{z_j^{(1)}}$$

$$\quad \quad \quad \cdot \delta_{mj} \delta_{ni} \cdot x_i^{(0)}$$

$$= 2 \cdot W_{km}^{(2)} \cdot \sigma' \Big|_{z_k^{(2)}} \cdot (\hat{y}_k - y_k) \cdot \sigma' \Big|_{z_m^{(1)}} \cdot x_n^{(0)}$$

finally, for the biases we get:

$$\frac{\partial \mathcal{L}}{\partial b_m^{(2)}} = 2 \cdot (\hat{y}_m - y_m) \cdot \sigma'|_{z_m^{(2)}}$$

$$\frac{\partial \mathcal{L}}{\partial b_n^{(1)}} = 2 w_{km}^{(2)} \cdot \sigma'|_{z_k^{(2)}} \cdot (\hat{y}_k - y_k) \cdot \sigma'|_{z_m^{(1)}}$$

gradient descent update rule:

$$w_{mn}^{(i)} \rightarrow w_{mn}^{(i)} - \xi \cdot \frac{\partial \mathcal{L}}{\partial w_{mn}^{(i)}}$$

$$b_m^{(i)} \rightarrow b_m^{(i)} - \xi \cdot \frac{\partial \mathcal{L}}{\partial b_m^{(i)}}$$

$\xi > 0$ : learning rate  
(a hyperparameter  
of the network)

We need to choose an activation function  $\sigma$ . Clever choice allows to compute derivative  $\sigma'$  exactly solely in terms of already computed function values  $\sigma$ . Here, we choose as an example the sigmoid:

$$\left\{ \begin{array}{l} \sigma'|_{z_k^{(2)}} = \sigma(z_k^{(2)}) \cdot (1 - \sigma(z_k^{(2)})) \\ \quad = \hat{y}_k \cdot (1 - \hat{y}_k) \\ \sigma'|_{z_m^{(1)}} = x_m^{(1)} \cdot (1 - x_m^{(1)}) \end{array} \right.$$

sigmoid function  $\sigma(x)$  as activation function & its derivative:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\Rightarrow \sigma'(x) = -\frac{1}{(1+e^{-x})^2} \cdot (-e^{-x})$$

$$= \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1+e^{-x}-1}{(1+e^{-x})^2}$$

$$= \frac{1}{1+e^{-x}} - \frac{1}{(1+e^{-x})^2}$$

$$= \sigma(x) - \sigma^2(x)$$

$$= \sigma(x)(1-\sigma(x))$$

$$\Rightarrow \left. \frac{d\sigma}{dx} \right|_x = \sigma|_x \cdot (1-\sigma)|_x$$

Final comment:

By including biases into the simple Boolean-function learning  $NN$  and generalizing back propagation - based gradient descent to stochastic gradient descent using randomized mini batching, this  $NN$  can learn handwritten digits 0-9 from the MNIST image dataset at  $\simeq 95\%$  recognition accuracy, if the  $NN$  has 1 hidden layer with 30 neurons densely connected.