- Activation Function: Sigmoid Function:

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

↳ Considering the input of an activation function $z = \sum_j w_j x_j + b$

$$\sigma(z) = \frac{1}{1 + e^{-(\sum_j w_j x_j + b)}}$$

↳ Derivative of it:

$$\frac{d\sigma}{dz} = \frac{d}{dz}\left(\frac{1}{1+e^{-z}}\right) = \frac{0\cdot(1+e^{-z}) - 1\cdot(-e^{-z})}{(1+e^{-z})^2} = \frac{e^{-z}}{(1+e^{-z})^2} = \left(\frac{1}{1+e^{-z}}\right)$$

$$\Rightarrow \boxed{\frac{\partial\sigma}{\partial z} = \sigma(1-\sigma)}$$

- Input:

↳ Hand Written Digit image of size $28 \times 28 = 784$ pixels

↳ The input pixels are greyscale, with $0.0 \to$ white and $1.0 \to$

↳ Thus, we have an input layer with $784$ "neurons".

↳ 50.000 training inputs/samples

- **Output / Desired output:** $\overset{\frown}{y}$

  $\hookrightarrow$ Desired output: column array One hot-encoded where a 1 means the desired Digit. In other words: $y[i] = 1 \Rightarrow$ the image is a hand written number $i$. // Ex.: $(0,0,0,1,0,0,0,0,0,0)^T \Rightarrow$ Image with a 3 //

  $\hookrightarrow$ Thus, the output layer has 10 neurons.

  $\hookrightarrow$ the predicted output is given by the activation function of the last layer.

- **Loss / Cost Function:**

$$C = \frac{1}{2m} \sum_{x=0}^{m} || y(x) - \hat{y}(x) ||^2$$

- **Gradient Descent:**

  $\hookrightarrow$ Let's suppose $C$ as a function of many variables $v_1 \dots v_m$:

  $$C(v_1 \dots v_m) \Rightarrow \Delta C \approx \sum_{i=1}^{n} \frac{\partial C}{\partial v_i} \Delta v_i \text{, for a tiny } \Delta v_i .$$

  $$\Rightarrow \Delta C \approx \nabla C \cdot \Delta v$$

  $\hookrightarrow$ Our goal is to choose wisely $\Delta v$ in order to make $\Delta C$ negative ( decrease the loss)

2

- **Output / Desired output:** $\overset{\frown}{y}$

  $\hookrightarrow$ Desired output: column array One hot-encoded where a 1 means the desired Digit. In other words: $y[i] = 1 \Rightarrow$ the image is a hand written number $i$. // Ex.: $(0,0,0,1,0,0,0,0,0,0)^T \Rightarrow$ Image with a 3 //

  $\hookrightarrow$ Thus, the output layer has 10 neurons.

  $\hookrightarrow$ the predicted output is given by the activation function of the last layer.

- **Loss / Cost Function:**

$$C = \frac{1}{2m} \sum_{x=0}^{m} || y(x) - \hat{y}(x) ||^2$$

- **Gradient Descent:**

  $\hookrightarrow$ Let's suppose $C$ as a function of many variables $v_1 \dots v_m$:

  $$C(v_1 \dots v_m) \Rightarrow \Delta C \approx \sum_{i=1}^{n} \frac{\partial C}{\partial v_i} \Delta v_i \text{, for a tiny } \Delta v_i .$$

  $$\Rightarrow \Delta C \approx \nabla C \cdot \Delta v$$

  $\hookrightarrow$ Our goal is to choose wisely $\Delta v$ in order to make $\Delta C$ negative ( decrease the loss)

2

↳ A way of doing that is choosing: $\Delta v = -\eta \nabla C$, → learning rate

because, in that way, $\Delta C = -\eta \|\nabla C\|^2$ //

↳ After a small $\Delta v$ chosen that way we will have:

$$v' = v - \eta \nabla C$$

↳ Finally, since C is a function of all weights and biases we want to modify and adjust to increase accuracy; the rules we will apply is:

$$w_k' = w_k - \eta \frac{\partial C}{\partial w_k}$$

$$b_\ell' = b_\ell - \eta \frac{\partial C}{\partial b_\ell}$$

↳ To speed up gradient descent algorithm we will implement the stochastic gradient descent, assuming that:

$m < < n$    $$\frac{\sum_{j=1}^{m} \nabla C_{S_j}}{m} \approx \frac{\sum_{j=1}^{n} \nabla C_{S_j}}{n} = \nabla C$$

⎱ — Each loop shuffle the samples array and take

loops $\left(\frac{\text{samples}}{\text{mini-batch-size}}\right)$ batches.

— Use each mini-batch to update the weights and biases. Repeat this process some amount of times. For instance: In one loop take 5.000 batches of 10 samples each. Repeat this process 30 times.

— It is important to say that after implementing a neural network, the values choose to $\eta$, m and the repetitions might require an adjustment to reach a satisfactory result.

- Back Propagation Algorithm :

1- Input a mini-batch

2- For each sample s perform the following steps :

- Feed forward : for each layer compute :
$$\begin{cases} z^l = w^l \cdot a^{l-1} + b \\ a^l = \sigma(z^l) \end{cases}$$

- Output error : Compute the error vector regarding the last layer :
$$\delta^L = \nabla_a C \odot \sigma'(z^L)$$

- Back propagate the error : For each layer starting from L-1 use the eq 2 to compute the error vector of this layer :
$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

- Compute the derivatives regarding weights and biases :
$$\frac{\partial C}{\partial b_j^l} = \delta_j^l \quad ; \quad \frac{\partial C}{\partial w_{k,j}^l} = a_k^{l-1} \cdot \delta_j^l$$

3- Gradient descent : After calculate all $\frac{\partial C}{\partial b}$ and $\frac{\partial C}{\partial w}$ update the weights and biases :
$$\begin{cases} (w_{kj}^l)' = w_{kj}^l - \frac{\eta}{m} \sum_{s=0}^{m} \frac{\partial C_s}{\partial w_{kj}^l} \\ (b_m^l)' = b_m^l - \frac{\eta}{m} \sum_{s=0}^{m} \frac{\partial C_s}{\partial b_m^l} \end{cases}$$

4

↳ Finally, after each mini-batch of training inputs:

$$w'_n = w_n - \frac{\eta}{m} \sum_{j=0}^{m} \frac{\partial C_{S_j}}{\partial w_n}$$

$$b'_l = b_l - \frac{\eta}{m} \sum_{j=0}^{m} \frac{\partial C_{S_j}}{\partial b_l}$$

• Back Propagation

- Relation between the activation output of a layer and the activation output of the previous layer:

$$a_j^l = \sigma \left( \underbrace{\sum_k w_{kj}^l a_k^{l-1} + b_j^l}_{z_j^l} \right)$$

- The quadratic cost for a single training sample:

$$C_x = \frac{1}{2} \| y - a^L \|^2 \qquad ↳ \frac{\partial C_x}{\partial a^L} = a^L - y //$$

Error Function:

- Let's suppose that in the neuron $(l, j)$ a small $\Delta z_j^l$ is applied in the input $z_j^l$ of the activation function. Then, the overall cost would change by an amount of: $\Delta C = \frac{\partial C}{\partial z_j^l} \Delta z_j^l$

— Notice that we will try to find a $\Delta z_j^l$ which makes the cost smaller.

— Notice that, if $\frac{\partial C}{\partial z_j^l}$ has a large absolute value, then we must choose a $\Delta z_j^l$ which has the opposite sign to $\frac{\partial C}{\partial z_j^l}$. But, if $\frac{\partial C}{\partial z_j^l}$ is small, then any choice of a small $\Delta z_j^l$ will not be sufficient to impact $\Delta C$. That means that, when $\frac{\partial C}{\partial z_j^l}$ has a big absolute value, we can use it to decrease the cost function. This neuron must be adjusted. Im another hand, when $\frac{\partial C}{\partial z_j^l}$ is small, we can't adjust this neuron to reach a smaller cost. We need to look for other neurons to adjust. This one, then, is considered already adjusted.

— Thus, we can choose $\frac{\partial C}{\partial z_j^l}$ as the neuron error, since we want to decrease it to reach a adjusted neuron.

Error:
$$\boxed{\delta_j^l = \frac{\partial C}{\partial z_j^l}}$$

Fundamental equations:

1) $\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L) \iff \delta^L = \nabla_a C \odot \sigma'(z^L)$

Computing the error in the last layer

2) • computing the relation between the error of one layer and the error of the next layer

$$\delta^\ell = \left((w^{\ell+1})^T \delta^{\ell+1}\right) \odot \sigma'(z^\ell)$$

3) the rate of change of cost function with respect to a bias:

what we want → $\dfrac{\partial C}{\partial b_j^\ell} = \delta_j^\ell$  ( Yes, it is the same error )

why? ⟹ $\sigma_j^\ell = \dfrac{\partial C}{\partial z_j^\ell}$ ⟹ $\underbrace{\dfrac{\partial C}{\partial b_j^\ell}}_{} = \underbrace{\dfrac{\partial C}{\partial z_j^\ell}}_{\delta_j^\ell} \cdot \underbrace{\dfrac{\partial z_j^\ell}{\partial b_j^\ell}}_{1}$ ⟹ $\dfrac{\partial C}{\partial b_j^\ell} = \delta_j^\ell$ //

4) An equation for the rate of change of the cost with respect to any weight:

$$\dfrac{\partial C}{\partial w_{nj}^\ell} = a_n^{\ell-1} \cdot \delta_j^\ell$$

↑
what we actually want