# Network Convention

n_inputs

n_neurons

n_outputs

input layer

output layer

hidden layers

$0^{th}$ layer ...

$L^{th}$ layer

① input layer

$$\vec{a}^1 = f(\vec{z}^1) \quad \text{where} \quad \vec{z}^1 = \omega^1 \vec{x} + \vec{b}^1$$

$$- \quad a_i^1 = f(z_i^1)$$

$$- \quad \vec{a}, \vec{z} \sim \text{column vectors} \quad \Rightarrow \quad \vec{x}, \vec{b} \sim \text{column vectors}$$

$$- \quad \omega^1 \sim \text{n\_neurons} \times \text{n\_inputs}$$

② hidden layers : $\quad \vec{a}^l = f(\vec{z}^l = \omega^l \vec{a}^{l-1} + \vec{b}^{l-1})$

$$\vec{a}^l, \vec{z}^l \sim (\text{n\_neurons} \times 1)$$

$$\omega^l \sim (\text{n\_neurons} \times \text{n\_neurons})$$

③ output layer :

$$\vec{y} = \vec{a}^L = f(\vec{z}^L)$$

$$\vec{a}^L, \vec{z}^L \sim (\text{n\_outputs} \times 1)$$

$$\omega^L \sim (\text{n\_outputs} \times \text{n\_neurons})$$

## Some Matrix Derivative Notes (Convention)

$$\vec{y} \sim n \atop \vec{x} \sim m \Rightarrow \frac{\partial \vec{y}}{\partial \vec{x}} \sim \begin{pmatrix} \partial y_1 / \partial x_1 & --- & \partial y_1 / \partial x_m \\ & \vdots & \\ \partial y_n / \partial x_1 & --- & \partial y_n / \partial x_m \end{pmatrix} \sim n \times m \sim \vec{y}\, \vec{x}^T$$

**General:** ① If $a \sim$ scalar, $X \sim n \times m$, $\dfrac{\partial a}{\partial X} \sim m \times n$

② If $\vec{z} \sim n \times 1$, $\omega \sim n \times m$, $\dfrac{\partial \vec{z}}{\partial \omega} \sim m \times 1$ ✓

$$\frac{\partial z_i}{\partial \omega_{jk}} = \frac{\partial}{\partial \omega_{jk}} \sum_b \omega_{ib} x_b = \sum_b x_b \, \delta_{ij} \delta_{kb} = x_k \, \delta_{ij}$$

## Back propagation

Given a loss function $\mathcal{L}$, determine gradients
with respect to all $\omega^i$ & $b^i$. Then use

$$C = \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}_i \sim \text{sample average}$$

Denote $\mathcal{L} = \mathcal{L}(\vec{y}, \vec{t})$ where $\vec{t}$ is the target data.

### $L^{th}$ layer

$$\frac{\partial \mathcal{L}}{\partial \omega^L} = \frac{\partial \mathcal{L}}{\partial \vec{y}} \frac{\partial \vec{y}}{\partial \vec{z}^L} \frac{\partial \vec{z}^L}{\partial \omega^L}$$

$$\vec{z}^L = \omega^L a^{L-1} + \vec{b}^L$$

$$; \quad \frac{\partial \vec{z}^L}{\partial \omega^L} \sim n\text{-outputs} / (n\text{-outputs} \times n\text{-neurons})$$
$$\sim (n\text{-neurons} \times 1)$$

$$= \frac{\partial \mathcal{L}}{\partial \vec{y}} f'(\vec{z}^L) \, \vec{a}^{L-1}$$
$$\llcorner (n\text{-neurons} \times 1)$$

$$\frac{\partial y_i}{\partial z_j} = \frac{\partial y_i}{\partial z_i} \, \delta_{ij} = \begin{pmatrix} 0 & \cdot & 0 \\ & \ddots & \\ 0 & \cdot & \cdot \end{pmatrix} \longrightarrow \begin{pmatrix} \partial y_1/\partial z_1 \\ \vdots \\ \partial y_n/\partial z_n \end{pmatrix}^T \sim (1, n\_outputs)$$

$$\frac{\partial f}{\partial \vec{y}} \sim (1, n\_outputs)$$

$$\frac{\partial f}{\partial w^L} \sim (w^L)^T \sim (n\_neurons \times n\_output) = (1 \times n\_outputs) \cdot (n\_neurons \times 1)$$

$$= n\_neurons \times n\_outputs$$

$$\Rightarrow \frac{\partial f}{\partial w^L} = \left( \frac{\partial f}{\partial \vec{y}} \circ f'(\vec{z}^L) \right) \cdot \vec{a}^{L-1}$$

metrix    row vector  $\times$  column vector    $\rightarrow$ needs column vector $\times$ row vector

Define $(\delta^L)^T = \dfrac{\partial f}{\partial \vec{y}} \odot f'(\vec{z}^L)$    $\llcorner$ Hademerd

- $(\delta^L) \sim (n\_outputs \times 1)$
- once we flip to row vectors, this will make senses

Note in Python (red):

$$\delta^L \sim (1 \times n\_outputs)$$

$$w_L \sim (n\_output \times n\_neurons)$$

$$\frac{\partial f}{\partial w^L} = \vec{a}^{L-1} (\delta^L)^T \sim (n\_neurons \times n\_output)$$

$$\frac{\partial f}{\partial \vec{b}^L} = (\delta^L)^T$$

$$(n\_neurons \times 1) \times (1 \times n\_outputs)$$
$$= (n\_neurons \times n\_outputs)$$

## Hidden-Layers

$n\_output \times n\_neurons$

$$\frac{\partial f}{\partial w^{L-1}} = \frac{\partial f}{\partial \vec{y}} \, \frac{\partial \vec{y}}{\partial \vec{z}^L} \, \frac{\partial \vec{z}^L}{\partial \vec{a}^{L-1}} \, \frac{\partial \vec{a}^{L-1}}{\partial \vec{z}^{L-1}} \, \frac{\partial \vec{z}^{L-1}}{\partial w^{L-1}}$$

$(\delta^L)^T \vec{a}^{L-1}$ (crossed out)

$$= (\delta^L)^T (w^L) \odot f'(\vec{z}^{L-1}) (\vec{a}^{L-2})$$

$$\begin{aligned}
\overset{\text{n\_neurons}}{\underset{\text{n\_neurons}}{\times}} = (1 \times \text{n\_outputs})\,\underline{(\text{n\_outputs} \times \text{n\_neurons})}\,(1 \times \text{n\_neurons})\,(\text{n\_neurons} \times 1)
\end{aligned}$$

$$= (1 \times \text{n\_neurons}) \odot (1 \times \text{n\_neurons}) \times (\text{n\_neurons} \times 1)$$

$$= (\text{n\_neurons} \times 1) \times \big[(1 \times \text{n\_neurons}) \odot (1 \times \text{n\_neurons})\big]$$

$$= (\text{n\_neurons}) \times (\text{n\_neurons}) \checkmark$$

$$\frac{\partial f}{\partial \omega^{L-1}} = \vec{a}^{\,L-2}\left( (\delta^L)^T \omega^L \odot f'(z^{L-1}) \right)$$

$$= \vec{a}^{\,L-2} (\delta^{L-1})^T$$

$$\frac{\partial f}{\partial \omega^{L-1}} = (\delta^{L-1})^T$$

$$(\delta^\ell)^T = (\delta^{\ell+1})^T \omega^{\ell+1} \odot f'(z^\ell)$$

$$\Downarrow$$

$$\frac{\partial f}{\partial \omega^\ell} = \vec{a}^{\,\ell-1} (\delta^\ell)^T$$

$$\frac{\partial f}{\partial b^\ell} = (\delta^\ell)^T$$

Note

$$\frac{\partial f}{\partial b^\ell} \sim (1 \times \text{n\_neurons}) \sim (\delta^\ell)^T$$

# Using Coding Conventions

Typical vectors are row vectors

$\implies$ transpose everything! (And keep orders!)

this ensures dimensions for products still work

**Note**: This is not globally consistent: $\delta^L \sim$ row vector, $b^L \sim$ row vector

but $\frac{\partial \mathcal{L}}{\partial b^L} \sim$ column

$$\frac{\partial f}{\partial w^L} = (a^{L-1})^T \delta^L \quad \& \quad \frac{\partial f}{\partial b^L} = \delta^L$$

But, it works for the **computations**

$$\frac{\partial f}{\partial w^\ell} = (a^{\ell-1})^T \delta^\ell \quad \& \quad \frac{\partial f}{\partial b^\ell} = \delta^\ell$$

$$\delta^L = \frac{\partial f}{\partial \vec{y}^L} \circ f'(\bar{z}^L)$$

$$\delta^\ell = f'(\bar{z}^\ell) \circ \left( \delta^{\ell+1} (w^{\ell+1})^T \right)$$

**Check:**

$$z^1 = \vec{x}_{1 \times n\_inputs} \, w^1_{n\_inputs \times n\_neurons} + \bar{b}^1_{1 \times n\_inputs}$$

$L-1$: $\delta^L \sim (1 \times n\_outputs)$

$w^L \sim (n\_neurons \times n\_outputs)$

$$\delta^{L-1} = f'(\bar{z}^{L-1}) \circ \left( \delta^L (w^L)^T \right)$$

$$(1 \times n\_neurons) \circ \left[ (1 \times n\_outputs) \times (n\_outputs \times n\_neurons) \right]$$

$$= (1 \times n\_neurons) \checkmark$$

$$(\delta^\ell)^T = (\delta^{\ell+1})^T \omega^{\ell+1} \odot f'(z^\ell)$$

$$\downarrow$$

$$\delta^\ell = \left[(\delta^{\ell+1})^T \omega^{\ell+1} \odot f'(z^\ell)\right]^T =$$

## Adding in Python data

- For training, we would like to feed forward all the training data simultaneously ~ ==vectorization.==

$$\bar{x} \longrightarrow X \sim (n\_samples, n\_inputs)$$

features

$$\underset{n\_samples \times n\_neurons}{\overset{\rightarrow}{z}^1} = \underset{n\_samples \times n\_inputs}{\overset{\leftrightarrow}{X}} \times \underset{n\_inputs \times n\_neurons}{\omega} + \underset{1 \times n\_neurons}{\overset{\rightarrow}{b}}$$

Python adds $\bar{b}$
to each row of $X\omega$

# Using Softmax Output Activation

If we change the output activation to a softmax,

$$y_i = a_i^L = S(z_i^L) = \frac{e^{z_i^L}}{\sum_k e^{z_k^L}}$$

we only affect $\delta_L$. Furthermore, let's take our NN to do max likelihood on the log likelihood. See "Softmax ..." note for details

Note: $\mathcal{L} = - \log$ likelihood

$$\frac{\partial \mathcal{L}}{\partial w^L} = - \frac{\partial \mathcal{L}}{\partial \vec{y}} \underbrace{\frac{\partial \vec{y}}{\partial \vec{z}}}_{\substack{\text{now a non-diagonal} \\ \text{matrix}}} \frac{\partial \vec{z}}{\partial w^L}$$

$$= - \sum_{i=1}^{N} \left( \vec{a}^{L-1}(x^{(i)}) \right)^T \left( \vec{y}^{(i)} - S(\vec{z}_s^L ; x^{(i)}) \right)$$

$$\frac{\partial C}{\partial w^L} = - \left( \vec{a}^{L-1} \right)^T \left( \vec{y} - S(\vec{z}_s^L) \right)$$

$$= \left( \vec{a}^{L-1} \right)^T \left( S(\vec{z}_s^L) - \vec{y} \right)$$

$$= \left( \vec{a}^{L-1} \right)^T \delta^L \quad \text{where} \quad \delta^L = S(\vec{z}_s^L) - \vec{y}$$

$$\vec{a}^{L-1} \sim (n\text{-samples} \times n\text{-neurons})$$ hidden

$$\vec{y} \sim (n\text{-samples} \times n\text{-classes})$$

$$S(\vec{z}_s^L) \sim (n\text{-samples}, n\text{-classes})$$

$$\frac{\partial C}{\partial w^L} \sim (n\text{-neurons} \times n\text{-samples}) \times (n\text{-samples}, n\text{-classes})$$

$$= (n\text{-neurons} \times n\text{-classes}) \checkmark$$

$( \text{n\_neuron} \times \text{n\_output} )$

$\omega^L \sim ( \text{n\_output} , \text{n\_neuron} )$

$$\frac{\partial L}{\partial \omega^L} = \frac{\partial L}{\partial \vec{y}^L} \frac{\partial \vec{y}^L}{\partial \vec{z}} \frac{\partial \vec{z}}{\partial \omega^L} \longrightarrow \frac{\text{n\_output}}{\text{n\_output} \times \text{n\_neuron}} \sim ( 1 \times \text{n\_neuron} )$$

$\nearrow$ n-row

$\uparrow$ n\_output $\times$ n\_neuron

n-row

$$\partial_{\omega_{ij}} \left( z_{\ell}^L = \sum_k \omega_{\ell k} a_k^{L-1} \right) = \sum_k a_k^{L-1} \delta^i_{\ell} \delta^j_{\cdot k}$$

$$= \boxed{a_j^{L-1}} \delta^{\overline{il}}$$

$$\sum_{i,j} \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial z_j} \frac{\partial z_j}{\partial \omega_{ab}} = \sum_{i,j} \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial z_j} a_b^{L-1} \delta^{aj}$$

$$\frac{\partial L}{\partial \omega_{ab}} = \sum_i \frac{\partial L}{\partial y_i} \left( \frac{\partial y_i}{\partial z_a} \right) a_b^{L-1}$$

$\uparrow \uparrow$ row colum

$\uparrow$ row

$n \times n$ Matrix

$n$ colum vector

row

column

$$\vec{z}^1 = x\,\omega + b \quad \leftarrow \text{row vector}$$

$(1 \times n\_neurons) \qquad (1 \times n\_input)(n\_input \times n\_neurons)$

$(N \times n\_neurons) = (N \times n\_input)(n\_input \times n\_neurons)$

$(N \times n\_neurons)^T$

$$\frac{\partial f}{\partial \omega^L} = (a^{L-1})^T \delta^L \quad \& \quad \frac{\partial f}{\partial b^L} = \delta^L$$

$$\frac{\partial f}{\partial \omega^\ell} = (a^{\ell-1})^T \delta^\ell \quad \& \quad \frac{\partial f}{\partial b^\ell} = \delta^\ell$$

$(N \times n\_outputs)$

$(n\_outputs \times N)$

$$\delta^L = \frac{\partial f}{\partial \vec{y}^L} \circ f'(\vec{z}^L)$$

$$\delta^\ell = f'(\vec{z}^\ell) \circ \left( \delta^{\ell+1}(\omega^{\ell+1})^T \right)$$

$$(\delta^L)^T = \text{row-vector}$$

$\downarrow \text{transpose}$

$$\delta^L =$$

$$(a^{L-1})(\delta^L)^T$$

$\downarrow$

$$\boxed{(a^{L-1})^T}\,\delta^L \qquad (\delta^L)(a^{L-1})^T$$

$\uparrow$ row vector $\qquad \uparrow$