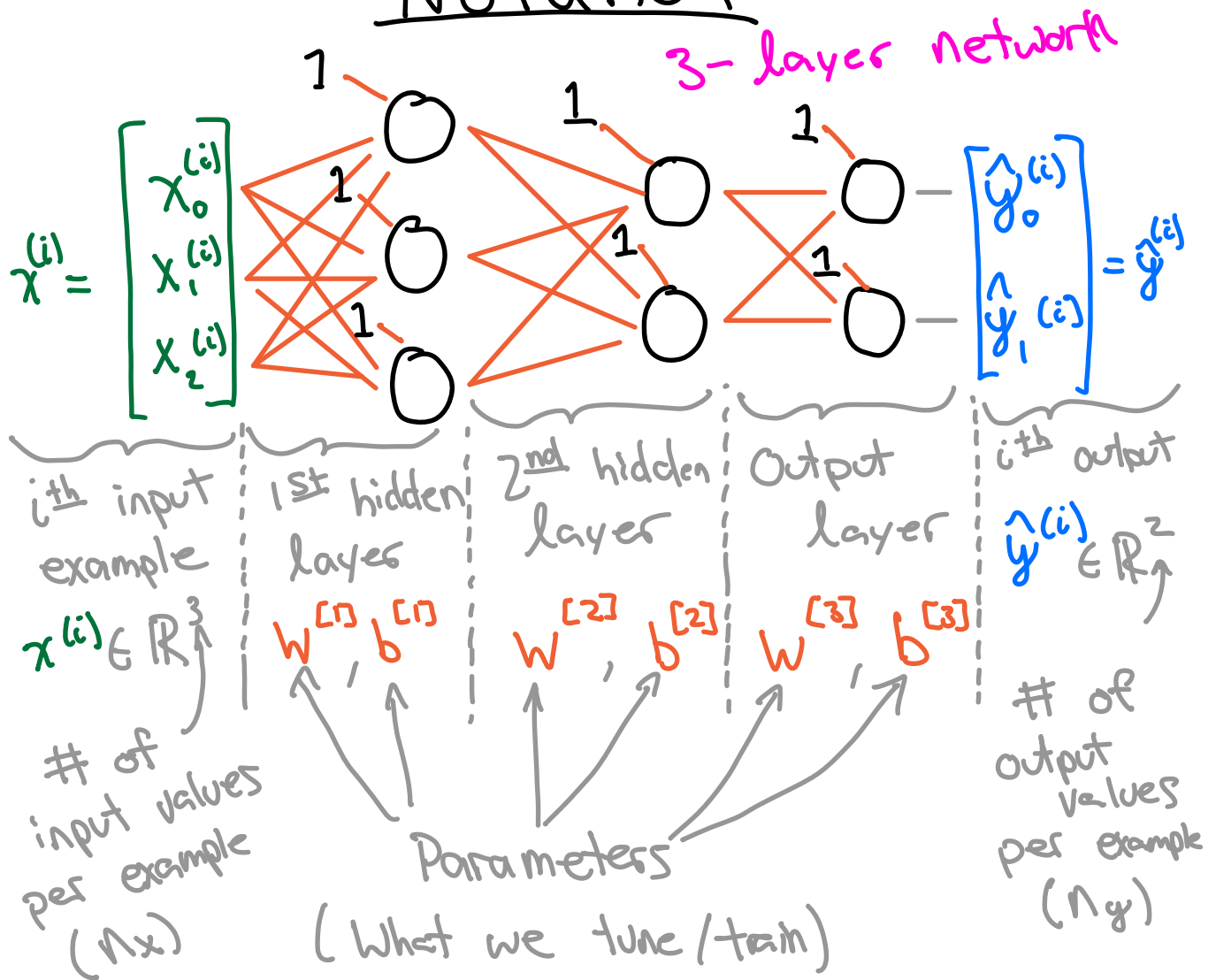


Notation



The input size (n_x) and output size (n_y) depend on the problem.

The number of layers (L) and the number of neurons per layer (n_x) are "Hyperparameters" chosen by you.

Summary

" (i) " parentheses are used to denote a specific input training example/sample

" [1] " brackets are used to denote a specific layer in the network

$x^{(i)}$: a single input example
with a shape $\in \mathbb{R}^{n_x}$,
which is a $(n_x, 1)$ column vector

$\hat{y}^{(i)}$: a single output example prediction
with a shape $\in \mathbb{R}^{n_y}$,
which is a $(n_y, 1)$ column vector

$W^{[l]}$: each hidden and output layer includes a matrix of parameter values (sometimes called the "weight" matrix). This matrix $\in \mathbb{R}^{n_l \times n_{l-1}}$ the shape is (n_l, n_{l-1})

of neurons in l^{th} layer # of neurons in $l-1$ layer

$b^{[L]}$: each hidden and output layer includes a vector of parameter values (sometimes called the "bias" vector). This vector $\in \mathbb{R}^{N_L}$, the shape is $(N_L, 1)$.

Forward Compute

We use $x^{(i)}$ and all $w^{[l]}$, $b^{[l]}$ values to compute $\hat{y}^{(i)}$.

Specifically, we compute the output of every neuron as follows:

$$z^{[l]}(i) = w^{[l]} a^{[l-1]}(i) + b^{[l]}$$

Linear output of all neurons in layer l given example i . The shape $\in \mathbb{R}^{n_l}$.

$$a^{[l]}(i) = g^{[l]}(z^{[l]}(i))$$

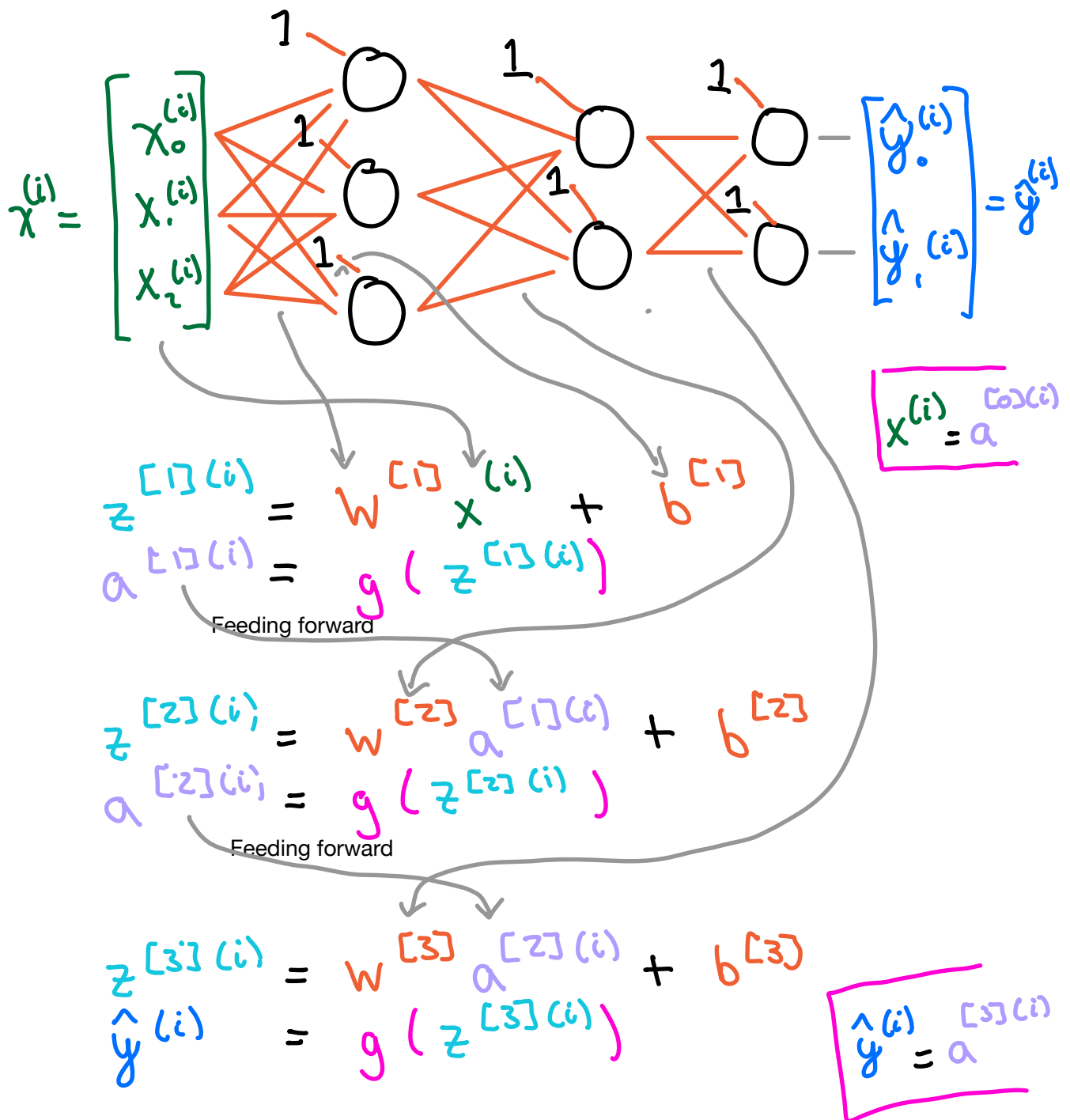
Activation function

"Activation" (output) of all neurons in layer l for the example i . The shape $\in \mathbb{R}^{n_l}$.

$$a^{[0]}(i) = x^{(i)} \quad \text{"Output" of input neurons}$$

$$\hat{y}^{(i)} = a^{[L]}(i) \quad \text{Output of output neurons}$$

We can now compute the output for the example network.



Loss

How do we compute performance?

We want some function that we can optimize. We need to know how to update the parameters $w^{(L)}$, $b^{(L)}$ so that we can drive our prediction to the true value

$$\overset{\text{hat}}{\rightarrow} \hat{y}^{(i)} \approx y^{(i)} \leftarrow \text{no hat}$$

↑ ↑
"prediction" "target" / label

Lets use the "mean-squared-error", aka MSE.

$$\underset{\substack{\uparrow \\ \text{"Loss"}}}{L}(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{2} \underbrace{(\hat{y}^{(i)} - y^{(i)})^2}_{\text{MSE}}$$

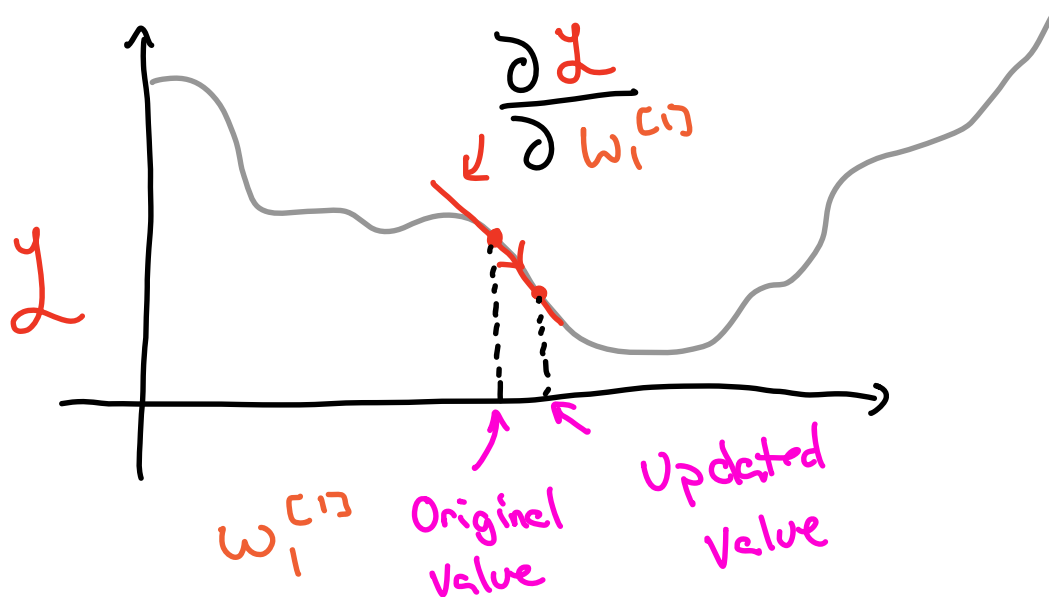
prediction true target

Back propagation

We want to find values for $w^{[L]}$, $b^{[L]}$ to minimize \mathcal{L} . Effectively, we want $\hat{y}^{(i)} \approx y^{(i)}$.

We can do this by starting with random values for $w^{[L]}$, $b^{[L]}$ and then updating them in such a way that we continually reduce \mathcal{L} .

In effect, we want to follow the slope of \mathcal{L} to smaller and smaller values.



So, we need to compute:

$$\frac{\partial \mathcal{L}}{\partial W^{[1]}}, \frac{\partial \mathcal{L}}{\partial b^{[1]}}, \text{ Parameters for layer 1}$$

$$\frac{\partial \mathcal{L}}{\partial W^{[2]}}, \frac{\partial \mathcal{L}}{\partial b^{[2]}}, \text{ Parameters for layer 2}$$

$$\frac{\partial \mathcal{L}}{\partial W^{[3]}}, \frac{\partial \mathcal{L}}{\partial b^{[3]}}, \text{ Parameters for layer 3}$$

Derive equations for these six quantities using the example 3-layer network, the half-MSE loss function, and assuming that the activation function $g(\cdot)$ is the sigmoid function.

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