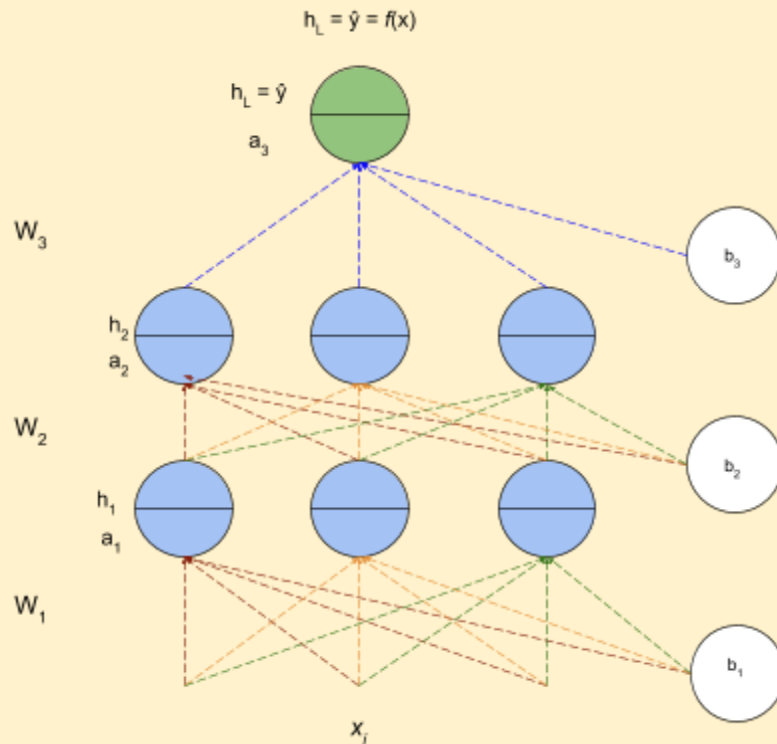


## Backpropagation (Light Math)

### Setting the context

Can we use the same learning algorithm as before?

1. Here is the learning algorithm as discussed in the previous chapter, the no-math version
2. Consider the Neural Network with the following configuration



3. The algorithm
  - a. **Initialise:**  $w_{111}, w_{112}, \dots, w_{313}, b_1, b_2, b_3$  randomly
  - b. **Iterate over data**
    - i. Compute  $\hat{y}$
    - ii. Compute  $L(w,b)$  Cross-entropy loss function
    - iii.  $w_{111} = w_{111} - \eta \Delta w_{111}$
    - iv.  $w_{112} = w_{112} - \eta \Delta w_{112}$
    - ...
    - v.  $w_{313} = w_{313} - \eta \Delta w_{313}$
    - vi.  $b_i = b_i + \eta \Delta b_i$
    - vii. Pytorch/Tensorflow have functions to compute  $\frac{\partial L}{\partial w}$  and  $\frac{\partial L}{\partial b}$
  - c. **Till satisfied**
    - i. Number of epochs is reached ( ie 1000 passes/epochs)
    - ii. Continue till Loss <  $\epsilon$  (some defined value)

# PadhAI: Backpropagation - the light math version

## One Fourth Labs

4. In this section, we will be looking at the light-math version, where we will be computing the derivatives
5. Derivatives for all layers from 1 to L

$$\begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial W_{111}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{11n}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{211}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{21n}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,11}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,1k}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,1k}} & \frac{\partial \mathcal{L}(\theta)}{\partial b_{11}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial b_{L1}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial W_{121}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{12n}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{221}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{22n}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,21}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,2k}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,2k}} & \frac{\partial \mathcal{L}(\theta)}{\partial b_{12}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial b_{L2}} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial W_{1n1}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{1nn}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{2n1}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{2nn}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,n1}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,nk}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,nk}} & \frac{\partial \mathcal{L}(\theta)}{\partial b_{1n}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial b_{Lk}} \end{bmatrix}$$

Layer 1                      Layer 2                      Output Layer nxk weights                      Bias terms

6. Once we know the gradients, we can use them in the Gradient Descent algorithm to compute the weights of the network