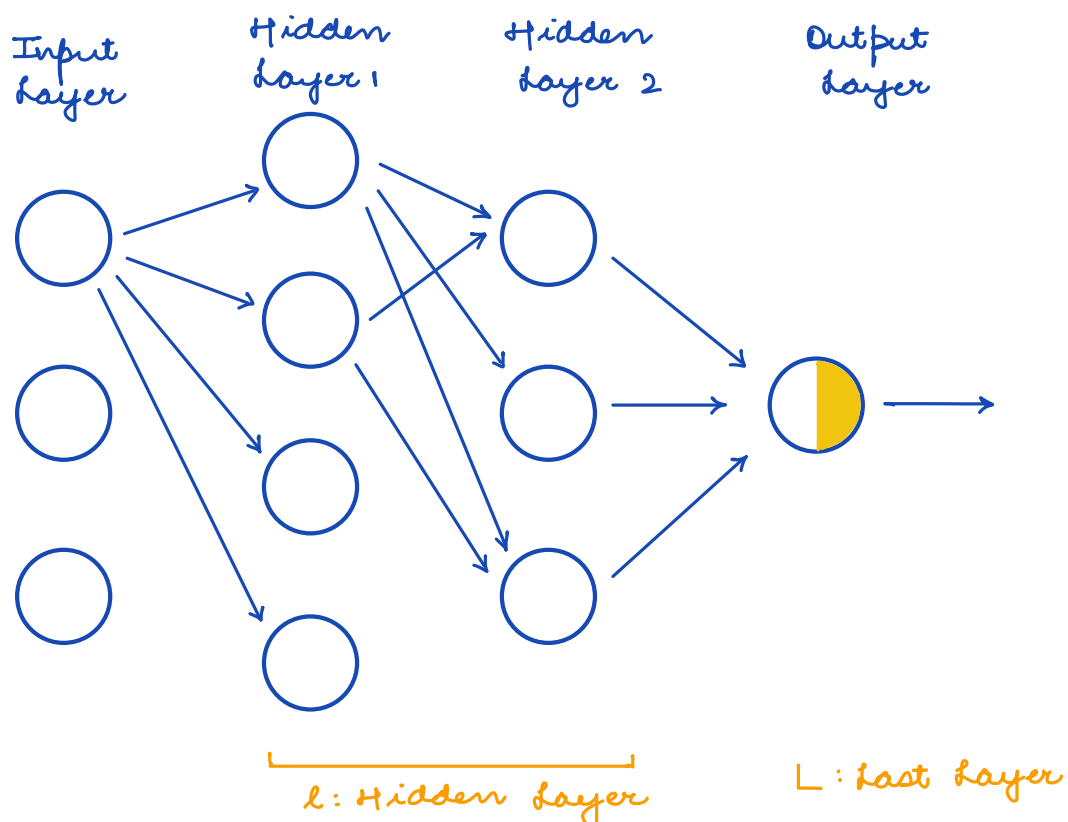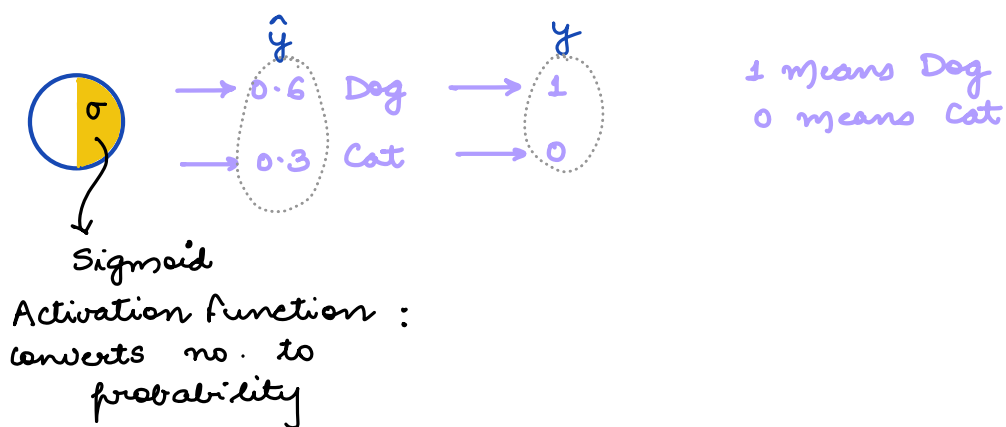Backpropagation is considered as a challenging problem in Machine Learning, reason being it involves knowledge of Linear Algebra and Calculus. Given the formula it is not very difficult to implement but the important thing is how you reached at the formula.

If you have a neural network then you can do a forward pass, get the prediction but prediction are going to be random because you are doing random initialization of weights. We want that our network should be able to learn and should give good results for regression and classification. Libraries like TensorFlow and PyTorch will automatically do Backpropagation for you.



Input Layer    Hidden Layer 1    Hidden Layer 2    Output Layer

$l$: Hidden Layer          $L$: Last Layer

Backpropagation algo is needed because : lets say your algo said that I am 60% sure that it is a dog but actual output you wanted is 1. Your algo said probability of cat is 0.3, you will take sigmoid and it will become 0.

$\hat{y}$ → 0.6 Dog → $y$ 1

→ 0.3 Cat → 0

Sigmoid
Activation function :
converts no. to
probability
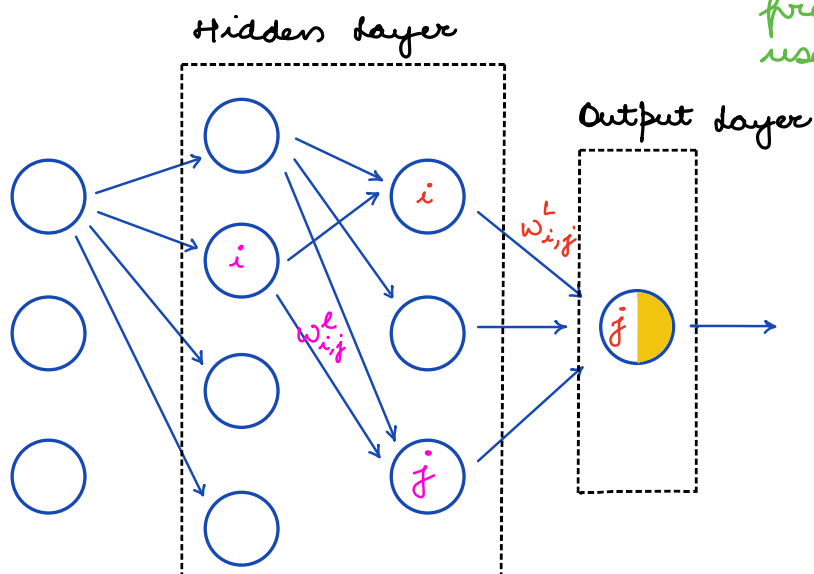
1 means Dog
0 means Cat

we want to figure out the loss. Loss is $\hat{y} - y$. Loss function we generally use in classification is Cross Entropy. For simplification, we are using mean squared error loss.

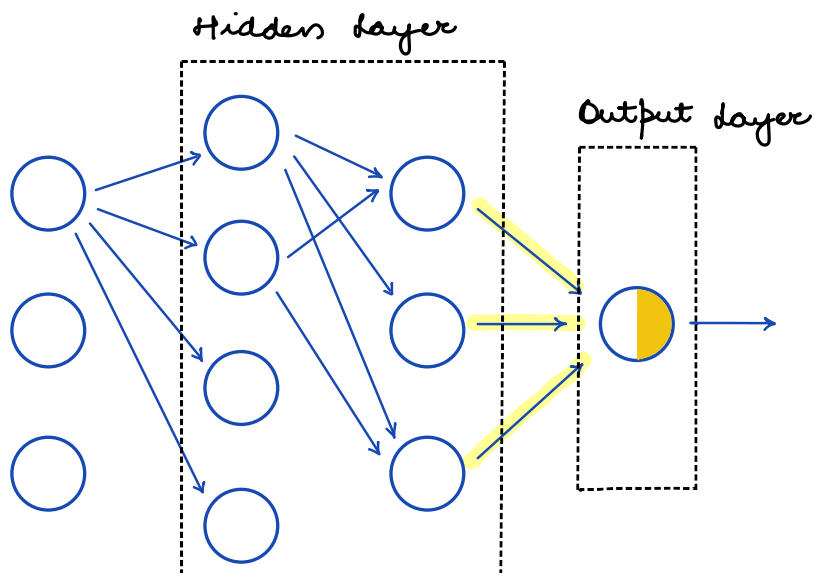$$MSE = L(\hat{y}, y) = \frac{1}{2} \sum_{i=1}^{m} \left( y^{(i)} - \hat{y}^{(i)} \right)^2$$

In general, this type of loss is non convex for a given classification problem, so we generally use cross entropy.

**Hidden Layer**

**Output Layer**

$w_{i,j}^{L}$

$w_{i,j}^{l}$

$i$

$i$

$j$

$j$

I want to analyse how loss function changes on changing the weights, if I make a change is $w_{i,j}^L$ then how does my error, loss and prediction change.

There are 2 type of neurons: one which are is output layer and other neurons are is hidden layer.

Output neurons are directly connected with loss, if you want to change your loss by $\delta L$ then you can find out what change I should make to these weights:



Hidden Layer

Output Layer

You need to compute 2 things:

Case 1: output Layer

$$\frac{\partial L}{\partial w_{i,j}^L}$$

$\downarrow$

$\partial L$ wrt w is last layer

$$\frac{\partial L}{\partial b^L}$$

$\downarrow$

$\partial L$ wrt bias is last layer

Case 2: Hidden Layer

$$\frac{\partial L}{\partial w_{i,j}^l}$$

$\curvearrowright$

$\partial L$ wrt w is hidden layer

$$\frac{\partial L}{\partial b^l}$$

$\downarrow$

$\partial L$ wrt bias is hidden layer

we need to find out how loss changes when weight and bias changes.

## Case 1:

$L(y, a)$

$$z_j = \sum_i^l w_{i,j}^L a_i^{l-1} + b_j^L$$

$i$ iterates over all the neurons in last layer $l$.

we want to find out how $\hat{y}$ changes with $w_{i,j}$

$$\frac{\partial L}{\partial w_{i,j}}$$

$$w_{i,j} \longrightarrow z_j \longrightarrow a_j \longrightarrow L(a, y)$$

$$\frac{\partial L}{\partial w_{i,j}} = \underbrace{\frac{\partial L}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j}}_{} \cdot \frac{\partial z_j}{\partial w_{i,j}} \qquad \text{(Chain Rule)}$$

①   ②   ③

$$\Downarrow$$

$$\delta_j^L = \frac{\partial L}{\partial z_j}$$

**①** $\dfrac{\partial L}{\partial a_j}$

$$L = \dfrac{1}{2} \sum_i \left( y_i^L - a_i^L \right)^2$$

$$\dfrac{\partial L}{\partial a_j} = \dfrac{1}{2} \cdot 2 \cdot \left( y_j - a_j \right)$$

*On differentiating wrt $a_j$, for all other terms it will become 0 and only $(y_j - a_j)$ will be left.*

$$\dfrac{\partial L}{\partial a_j} = -\left( y_j - a_j \right)$$

| $L$ | $a^L$ | $y^L$ |
|---|---|---|
| | 0·3 | 0 |
| $i \quad \rightarrow a_i^L$ | 0·5 | 1 |
| | 0·2 | 0 |

**②** $\dfrac{\partial a_j}{\partial z_j}$

$$a_j = \sigma(z_j)$$

$$\dfrac{\partial a_j}{\partial z_j} = \sigma'(z_j)$$

$$= \left( 1 - \sigma(z_j) \right) \sigma(z_j)$$

**③** $\dfrac{\partial z_j}{\partial w_{i,j}}$

$$z_j^{\ell} = \sum_i w_{i,j} \cdot \underline{a_i^{\ell-1}} + b_j^{\ell}$$
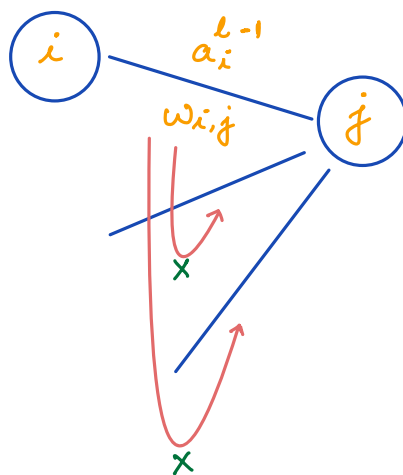
*If you take derivative wrt $w_{i,j}$ then you will get activation. For all other terms it will be 0 and only $a_i$ will be left.*

$$\frac{\partial z_j}{\partial w_{i,j}} = a_i^{l-1}$$

It is the activation associated with previous layer.

$l-1$ layer          $L$ layer



$W_{i,j}$ will depend on $a_i^{l-1}$ activation.

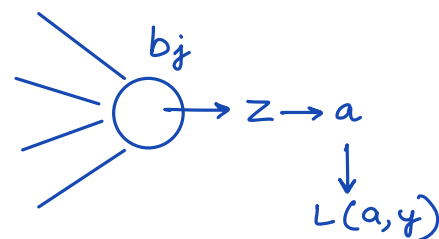$W_{i,j}$ won't be influenced by other activation values.

Combining ①, ② and ③

$$\frac{\partial L}{\partial w_{i,j}} = \underbrace{-(y_j - a_j)\, \sigma'(z_j)}_{\delta_j^L}\; a_i^{l-1}$$

$$\frac{\partial L}{\partial w_{i,j}} = \delta_j^L \cdot a_i^{l-1}$$

Bias

$$\frac{\partial L}{\partial b_j} = \underbrace{\frac{\partial L}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j}}_{\delta_j^L} \cdot \frac{\partial z_j}{\partial b_j}$$



$$z_j^l = \sum_i^l w_{i,j}\, a_i^{l-1} + \boxed{b_j^l}$$

will become 0 (under the sum)    will become 1 (circled $b_j^l$)

$$\frac{\partial z_j}{\partial b_j} = 1$$

$$\frac{\partial L}{\partial b_j} = \delta_j^L \cdot 1$$

for output layer,

$$\frac{\partial L}{\partial w_{i,j}} = \delta_j^L \cdot a_i^{l-1}$$

$$\frac{\partial L}{\partial b_j} = \delta_j^L \cdot 1$$

$$\delta_j^L = -(y_j - a_j)\, \sigma'(z_j)$$

$$\delta_j^L = (a_j - y_j)\, \sigma'(z_j)$$