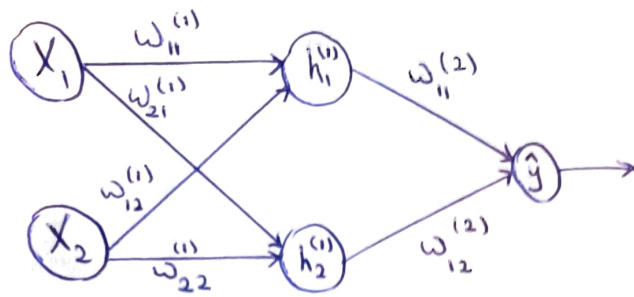


# Backpropagation



Convention:

$w_{ij}^{(l)} \rightarrow$  # of layer  
 $i, j \rightarrow$   $i$ th neuron  $j$ th input

Firstly, let's get the Forward pass equations:

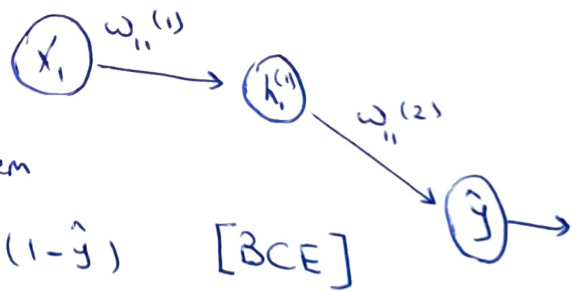
First Hidden Layer  $\begin{cases} z_1^{(1)} = w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + b_1^{(1)} \Rightarrow h_1^{(1)} = \text{ReLU}(z_1^{(1)}) \\ z_2^{(1)} = w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + b_2^{(1)} \Rightarrow h_2^{(1)} = \text{ReLU}(z_2^{(1)}) \end{cases}$

Output Layer  $\begin{cases} z_1^{(2)} = w_{11}^{(2)} h_1^{(1)} + w_{12}^{(2)} h_2^{(1)} + b_1^{(2)} \Rightarrow \hat{y} = \sigma(z_1^{(2)}) \end{cases}$

Secondly, let's try to get eqn of weight update using backpropagation

if we want to see how we want to update  $w_{11}^{(1)}$

Assume  $L = \frac{1}{2} (y - \hat{y})^2$  if we have regression [MSE]



But since that's a classification problem

Assume  $L = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$  [BCE]

$$\frac{\partial L}{\partial w_{11}^{(1)}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial h_1^{(1)}} \cdot \frac{\partial h_1^{(1)}}{\partial z_1^{(1)}} \cdot \frac{\partial z_1^{(1)}}{\partial w_{11}^{(1)}}$$

$$\frac{\partial L}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}, \quad \frac{\partial \hat{y}}{\partial z_1^{(2)}} = \hat{y}(1-\hat{y})$$

$$\frac{\partial z_1^{(2)}}{\partial h_1^{(1)}} = w_{11}^{(2)}$$

$$\frac{\partial h_1^{(1)}}{\partial z_1^{(1)}} = \begin{cases} 1 \\ 0 \end{cases}$$

$$\begin{aligned} z_1^{(1)} &> 0 \\ z_1^{(1)} &\leq 0 \end{aligned}$$

$$\frac{\partial z_1^{(1)}}{\partial w_{11}^{(1)}} = x_1$$

$$\frac{\partial L}{\partial w_{11}^{(1)}} = \begin{cases} \left[ -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right] \hat{y}(1-\hat{y}) w_{11}^{(2)} x_1 \\ 0 \end{cases}$$

$$z_1^{(1)} > 0$$

$$z_1^{(1)} \leq 0$$

$$\frac{\partial L}{\partial w_{11}^{(1)}} = \begin{cases} (\hat{y} - y) w_{11}^{(2)} x_1 \\ 0 \end{cases}$$

$$z_1^{(1)} > 0$$

$$z_1^{(1)} \leq 0$$

Think about  $\frac{\partial L}{\partial b_1^{(1)}}$

$$\frac{\partial L}{\partial b_1^{(1)}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial h_1^{(1)}} \cdot \frac{\partial h_1^{(1)}}{\partial z_1^{(1)}} \cdot \frac{\partial z_1^{(1)}}{\partial b_1^{(1)}}$$

It is already calculated

layer  $\delta_1^{(1)}$  (error signal) of neuron one in first hidden layer  
# no. of neuron

$\delta$  typically tells how much that neuron contributed to overall loss.

$$\frac{\partial L}{\partial b_1^{(1)}} = \delta_1^{(1)}$$

$$\frac{\partial L}{\partial w_{12}^{(1)}} = \delta_1^{(1)} x_2, \quad \frac{\partial L}{\partial w_{11}^{(1)}} = \delta_1^{(1)} x_1$$

So we need gradients of each weight and each bias

Gradients tell you how much changing the weight would change the loss

$$\frac{\partial L}{\partial w_{ij}^{(l)}} = \delta_i^{(l)} a_j^{(l-1)}, \quad \frac{\partial L}{\partial b_i^{(l)}} = \delta_i^{(l)}$$

So, if  $a_j^{(l-1)}$  (input) is large  
and  $\delta_i^{(l)}$  (Sensitivity to loss) is large

$\therefore$  Weight has big effect on loss (should be adjusted)

After calculating the gradient, update the weights using learning rate  $\eta$

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta \frac{\partial L}{\partial w_{ij}^{(l)}}$$

$$b_i^{(l)} \leftarrow b_i^{(l)} - \eta \frac{\partial L}{\partial b_i^{(l)}}$$

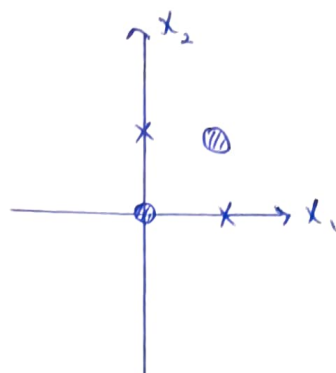
Note: Variants like Momentum, Adam, etc.. modify this basic rule by adding terms that depend on the past gradients. But the core idea remains subtracting (learning rate  $\times$  gradient)

Key points:

- Learning rate  $\eta$  Controls how big is the step
  - Too large  $\eta \rightarrow$  overshoot
  - Too small  $\eta \rightarrow$  slow learning
- You perform these updates after each forward + backward pass (each training example or mini-batch)

For Simplicity, let's look at this dataset [XOR]:

$x_1$	$x_2$	target
0	0	0
0	1	1
1	0	1
1	1	0



● Zero  
 x one

$\therefore$  not linear data  $\therefore$  NN  
 Tabular data  $\rightarrow$  MLP

let's initialize random weights:  $b'_{2 \times 1} = \begin{bmatrix} b'_1 \\ b'_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$

$$w'_{2 \times 2} = \begin{bmatrix} w'_{11} & w'_{12} \\ w'_{21} & w'_{22} \end{bmatrix} = \begin{bmatrix} 1 & 0.5 \\ -0.5 & 1 \end{bmatrix}$$

$$w^2_{1 \times 2} = \begin{bmatrix} w^2_{11} & w^2_{12} \end{bmatrix} = \begin{bmatrix} -1 & 0.5 \end{bmatrix}, b^2_1 = 1$$

$x_1$	$x_2$	$z_1^{(1)}$	$h_1^{(1)}$	$h_2^{(2)}$	$z_1^{(2)}$	$\hat{y}$	$y$	BCE	} Understand the numbers
0	0	1	1	0.5	0.25	0.562	0	0.82	
0	1	1.5	1.5	1.5	0.25	0.562	1	0.57	
1	0	2	2	0	-1	0.27	1	1.3	
1	1	2.5	2.5	1	-1	0.27	0	0.31	

$\rightarrow z_1^{(1)}$  is always +ve

After sample 1)

$$\frac{\partial L}{\partial b_1^{(1)}} = \delta_1^{(1)} = (\hat{y} - y) w_{11}^{(2)} = (0.562 - 0)(-1) = -0.562$$

let  $\eta = 0.1$

$$\frac{\partial L}{\partial w_{11}^{(1)}} = -\delta_1^{(1)} x_1 = -0.562(0) = 0$$

$$b_1 \leftarrow b_1 - \eta \frac{\partial L}{\partial b_1} = 1 - 0.1(-0.562) = 1.056$$

no update for weight

Note: This NN won't yield 100% accuracy (it should since it's XOR), but why?

It's so simple, adding extra hidden layer will solve this problem, but I used this NN for simplicity, now you got the core idea.

Updating weights after each sample isn't a good idea

- Noisy
- Inefficiency
- Convergence Stability

Solution?

Batch or mini-batch gradient descent

$$\frac{\partial L}{\partial b_{11}^{(1)}} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i) w_{11}^{(2)}$$

$$\frac{\partial L}{\partial w_{11}^{(1)}} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i) w_{11}^{(2)} x_{1i}$$

Notice how it updates the weights only once but we subtract each  $y$  from  $\hat{y}$

In other words, we don't average  $\hat{y}_i - y_i$ . But we average the gradient itself.

==