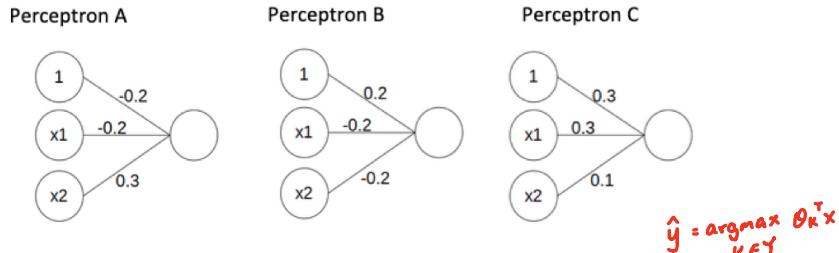


3. Consider a multiclass classification task using the three perceptrons below:



Recall the training procedure for the multi-class perceptron: we predict \hat{y} as the class whose parameters lead to the largest activation. During training, if a wrong class was predicted, we update the parameters activated for the wrong class (decreasing their activation), and the parameters activated for the correct class (increasing their activation).

Following this scheme, do the following:

- a) *Training:* For training example $x_{\text{train}} = (0,1)$, with class **B**, give the updated weights of each perceptron after processing x . Use learning rate $\eta = 0.1$.

$$\left. \begin{array}{l} \text{Perceptron A : } -0.2 + 0.3 = 0.1 \\ \text{Perceptron B : } 0.2 - 0.2 = 0 \\ \text{Perceptron C : } 0.3 + 0.1 = 0.4 \end{array} \right\} \hat{y} = \operatorname{argmax}_{k \in Y} \theta_k^T x$$

\Rightarrow Predict C (not match with "B")

\Rightarrow Update Perceptron ^{true}_B & ^{pred}_C

(B) $\theta_{y_i} \leftarrow \theta_{y_i} + \eta x^i$ move towards predicting y_i for x^i

(C) $\theta_{\hat{y}_i} \leftarrow \theta_{\hat{y}_i} - \eta x^i$ move away from predicting \hat{y}_i for x^i

$$B: \quad \underline{\theta}^{(0)} = \{0.2, -0.2, -0.2\} \quad \eta = 0.1 \quad C: \quad \theta^0 = \{0.3, 0.3, 0.1\}$$

$$\theta_0^{(1)} = 0.2 + 0.1 \times 1 = 0.3 \quad \theta_0^{(1)} = 0.3 - 0.1 \times 1 = 0.2$$

$$x = (0,1) \quad \theta_1^{(1)} = -0.2 + 0.1 \times 0 = -0.2 \quad \theta_1^{(1)} = 0.3 - 0.1 \times 0 = 0.3$$

$$\theta_2^{(1)} = -0.2 + 0.1 \times 1 = -0.1 \quad \theta_2^{(1)} = 0.1 - 0.1 \times 1 = 0$$

- b) *Testing*: Using the new weights from part (a), given test example $x_{\text{test}} = (1, 1)$ with true class **C**, give the class prediction of the group of three perceptrons. Is the prediction correct? Should we change the weights?

$$A: \underline{\theta} = \{-0.2, -0.2, 0.3\}$$

$$\text{Perceptron A : } -0.2 - 0.2 + 0.3 = -0.1$$

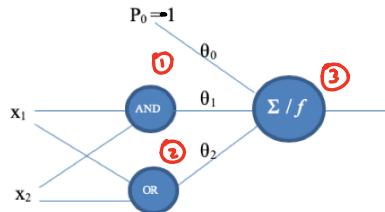
$$\text{Perceptron B : } 0.3 - 0.2 - 0.1 = 0$$

$$\text{Perceptron C : } 0.2 + 0.3 + 0 = 0.5$$

\Rightarrow Predict C (Correct)

Not update (test instance)

1. Consider the two levels deep network illustrated below. It is composed of three perceptrons. The two perceptrons of the first level implement the AND and OR function, respectively.



* Determine the weights θ_1 , θ_2 and threshold θ_0 such that the network implements the XOR function. The initial weights are set to zero, i.e. $\theta_0 = \theta_1 = \theta_2 = 0$, and the learning rate η (eta) is set to 0.1.

Notes:

- The input function for the perceptron on level 2 is the weighted sum (Σ) of its input.
 - The activation function f for the perceptron on level 2 is a step function:
- $$f = \begin{cases} 1 & \text{if } \sum > 0 \\ 0 & \text{otherwise} \end{cases}$$
- Assume that the weights for the perceptrons of the first level are given. (AND, OR)

XOR: for each training example π

$$\left\{ \begin{array}{l} p \leftarrow (p_0, f_{\text{AND}}(\pi), f_{\text{OR}}(\pi)) \quad (\text{level 2}) \\ \hat{y} \leftarrow f(\sum_i \theta_i p_i) = \begin{cases} 1 & \text{if } \sum > 0 \\ 0 & \text{otherwise} \end{cases} \\ y \leftarrow \text{target of } \pi \quad (\text{XOR}) \end{array} \right.$$

train: For $i=1:n$

$$\Delta \theta_i \leftarrow \eta (y - \hat{y}) p_i \quad \left. \begin{array}{l} \text{(update } \theta \text{)} \\ \text{training} \end{array} \right\}$$

$$\theta_i \leftarrow \theta_i + \Delta \theta_i$$

Output of level 1:

instance	x_1	x_2	$f_{\text{AND}}(\pi)$	$f_{\text{OR}}(\pi)$	target y XOR
1	1	0	0	1	1
2	0	1	0	1	1
3	1	1	1	1	0
4	0	0	0	0	0

all possibilities

Inputs for level 2 perceptron: $\langle p_0, p_1, p_2 \rangle = \langle -1, p_1, p_2 \rangle$

Initial params $\theta = \langle \theta_0, \theta_1, \theta_2 \rangle = \langle 0, 0, 0 \rangle$

$$\eta = 0.1$$

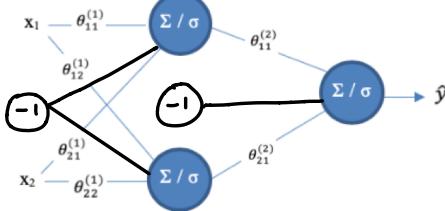
For first epoch (training):

$\langle p_0, p_1, p_2 \rangle$	$\Sigma = \theta_0 p_0 + \theta_1 p_1 + \theta_2 p_2$	$\hat{y} = f(\Sigma)$	y	$\Delta \theta_i = \eta(y - \hat{y}) p_i$
① $\langle -1, 0, 1 \rangle$	0 update θ : $\theta \leftarrow \theta + \Delta \theta$ $\langle 0, 0, 0 \rangle + \langle -0.1, 0, 0.1 \rangle$ $= \langle -0.1, 0, 0.1 \rangle$	0	1	$0.1(1-0) \times \langle -1, 0, 1 \rangle$ $= \langle -0.1, 0, 0.1 \rangle$
② $\langle -1, 0, 1 \rangle$	$-1(-0.1) + 0.1 = 0.2$	1	1	(No update)
③ $\langle -1, 1, 1 \rangle$	$-0.1(-1) + 0.1 \times 1 = 0.2$ update θ : $\langle -0.1, 0, 0.1 \rangle + \langle 0.1, -0.1, -0.1 \rangle$ $= \langle 0, -0.1, 0 \rangle$	1	0	$0.1(0-1) \times \langle -1, 1, 1 \rangle$ $= \langle 0.1, -0.1, -0.1 \rangle$
④ $\langle -1, 0, 0 \rangle$	$-0.1 \times 0 = 0$	0	0	(No update)
...

No changes after 4 epochs \rightarrow Converge

Final $\theta = \langle 0, -0.2, 0.1 \rangle$

2. Consider the following multilayer perceptron.



The network should implement the XOR function. Perform one epoch of *backpropagation* as introduced in the lecture on multilayer perceptrons.

Notes:

- The activation function f for a perceptron is the *sigmoid function*:

$$f(x) = \frac{1}{1 + e^{-x}}$$

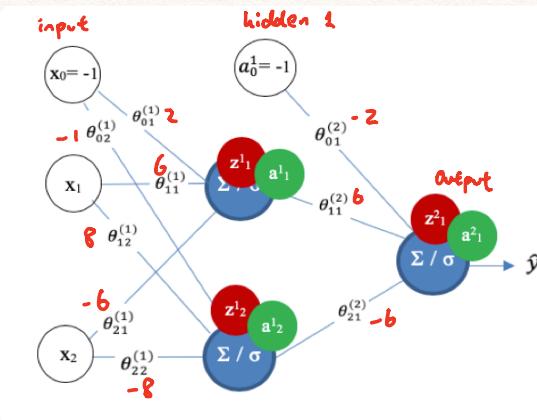
- The ~~thresholds~~ are not shown in the network. The ~~threshold~~ nodes are set to -1.
- Use the following initial parameter values:

$$\theta_{01}^{(1)} = 2 \quad \theta_{02}^{(1)} = -1 \quad \theta_{01}^{(2)} = -2$$

$$\begin{array}{lll} \theta_{11}^{(1)} = 6 & \theta_{12}^{(1)} = 8 & \theta_{11}^{(2)} = 6 \\ \theta_{21}^{(1)} = -6 & \theta_{22}^{(1)} = -8 & \theta_{21}^{(2)} = -6 \end{array}$$

- The learning rate is set to $\eta = 0.7$

- Compute the activations of the hidden and output neurons; $a_1^{(1)}, a_2^{(1)}, a_1^{(2)}$
level 1 level 2
- Compute the error of the network;
- Backpropagate the error to determine $\Delta\theta_{ij}$ for all weights θ_{ij} and updates the weight θ_{ij} .



(i) For level 1:

$$a_1^{(1)} = \sigma(z_1^{(1)}) = \sigma(-1(-1) + 6x_1 - 6x_2) = \sigma(-2 + 6x_1 - 6x_2)$$

$$a_2^{(1)} = \sigma(z_2^{(1)}) = \sigma(-1(-1) + 8x_1 - 8x_2) = \sigma(1 + 8x_1 - 8x_2)$$

For level 2:

$$(\text{Output } \hat{y} =) a_1^{(2)} = \sigma(z_1^{(2)}) = \sigma(-2(-1) + 6a_1^{(1)} - 6a_2^{(1)}) = \sigma(2 + 6a_1^{(1)} - 6a_2^{(1)})$$

(ii) Error:

$$E = \frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}(y - a_1^{(2)})^2$$

(iii) Backpropagation (efficient way of computing partial derivatives of

the error of MLP wrt. each weight $\theta_{jk}^{(l)}$. $\frac{\partial E}{\partial \theta_{jk}^{(l)}}$

1. Forward pass 2. Compute error 3. Backward pass (propagate error terms)

4. Update all θ $\underbrace{\Delta \theta}_{(\text{at once})}$

$$\text{GD: } \theta_i \leftarrow \theta_i - \eta \frac{\partial E}{\partial \theta_i}$$

\nwarrow By back propagation

In MLP:

$$\theta_{jk}^{(l)} \leftarrow \theta_{jk}^{(l)} + \Delta \theta_{jk}^{(l)}$$

$$\Delta \theta_{jk}^{(l)} = -\eta \frac{\partial E}{\partial \theta_{jk}^{(l)}} = \eta \delta_k^{(l)} a_j^{(l-1)}$$

Error term: $\delta_k^{(l)} = (y - a_k^{(l)}) \underline{\sigma'(z_k^{(l)})}$ (lecture)
 (last layer) $= (y - a_k^{(l)}) \underline{\sigma(z_k^{(l)}) (1 - \sigma(z_k^{(l)}))}$

Chain Rule:

$$E = \frac{1}{2}(y - a_k^{(l)})^2$$

$$\begin{aligned} \frac{\partial E}{\partial \theta_{jk}^{(l)}} &= \frac{\partial E}{\partial a_k^{(l)}} \cdot \frac{\partial a_k^{(l)}}{\partial z_k^{(l)}} \cdot \frac{\partial z_k^{(l)}}{\partial \theta_{jk}^{(l)}} \\ &= -(y - a_k^{(l)}) \cdot \underline{\sigma'(z_k^{(l)})} \cdot a_j^{(l)} \\ &= -\delta_k^{(l)} a_j^{(l-1)} \end{aligned}$$

prev layer: $(l-1)$

$$\text{If } \frac{\partial E}{\partial \theta_{jk}^{(l-1)}} = \underbrace{\frac{\partial E}{\partial a_k^{(l)}}}_{\delta_k^{(l)}} \cdot \underbrace{\frac{\partial a_k^{(l)}}{\partial z_k^{(l)}}}_{\theta_{jk}^{(l)}} \cdot \underbrace{\frac{\partial z_k^{(l)}}{\partial a_{ik}^{(l-1)}}}_{\sigma'(z_k^{(l-1)})} \cdot \underbrace{\frac{\partial a_{ik}^{(l-1)}}{\partial z_k^{(l-1)}}}_{a_{ik}^{(l-2)}} \cdot \underbrace{\frac{\partial z_k^{(l-1)}}{\partial \theta_{jk}^{(l-1)}}}_{\text{new } \delta_{ik}^{(l-1)} \text{ at prev layer}}$$

For training instance $\langle 1, 0 \rangle$: $y = 1$

$$a_1^{(1)} = \sigma(-2 + 6x_1 - 6x_2) = \sigma(4) = 0.982$$

$$a_2^{(1)} = \sigma(1 + 8x_1 - 8x_2) = \sigma(9) = 0.999$$

Output: $a_1^{(2)} = \sigma(2 + 6a_1^{(1)} - 6a_2^{(1)}) = \sigma(2 + 6 \times 0.982 - 6 \times 0.999) = \sigma(1.898) = 0.8691$

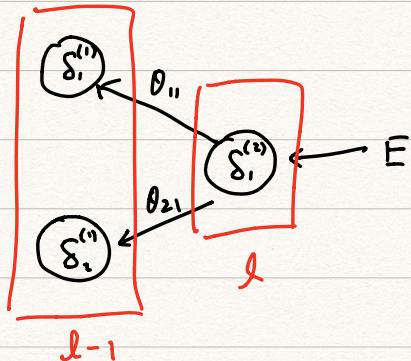
$$E = \frac{1}{2} (y - \hat{y})^2 = \frac{1}{2} (1 - 0.8691)^2 = 0.0086$$

Back propagation: (right to left)

$$\begin{aligned}\delta_1^{(2)} &= (y - a_1^{(2)}) \sigma(z_1^{(2)}) (1 - \sigma(z_1^{(2)})) \\ &= (1 - 0.8691) \sigma(1.898) [1 - \sigma(1.898)] \\ &= 0.0149\end{aligned}$$

$$\begin{aligned}\delta_1^{(1)} &= \theta_{01}^{(2)} \delta_1^{(2)} \sigma(z_1^{(1)}) [1 - \sigma(z_1^{(1)})] \\ &= 6 \times 0.0149 \times \sigma(4) \times [1 - \sigma(4)] \\ &= 0.0016\end{aligned}$$

$$\begin{aligned}\delta_2^{(2)} &= \theta_{12}^{(2)} \delta_1^{(2)} \sigma(z_2^{(2)}) [1 - \sigma(z_2^{(2)})] \\ &= -6 \times 0.0149 \times \sigma(9) \times [1 - \sigma(9)] \\ &= -0.00000103\end{aligned}$$



$\eta = 0.7$, calculate $\Delta \theta$: bias (at once)

$$\Delta \theta_{01}^{(2)} = \eta \delta_1^{(2)} a_0^{(1)} = 0.7 \times 0.0149 \times (-1) = -0.0104$$

$$\Delta \theta_{12}^{(2)} = \eta \delta_1^{(2)} a_1^{(1)} = 0.7 \times 0.0149 \times 0.982 = 0.0102$$

...

...

Update θ :

$$\theta_{01}^{(2)} \leftarrow \theta_{01}^{(2)} + \Delta \theta_{01}^{(2)}$$

$$= -2 + (-0.0104) = -2.0104$$

$$\theta_{11}^{(2)} \leftarrow \theta_{11}^{(2)} + \Delta \theta_{11}^{(2)}$$

$$= 6 + 0.0102 = 6.0102$$

...

This is only one training instance \Rightarrow do this for all training instances
to finish one epoch.