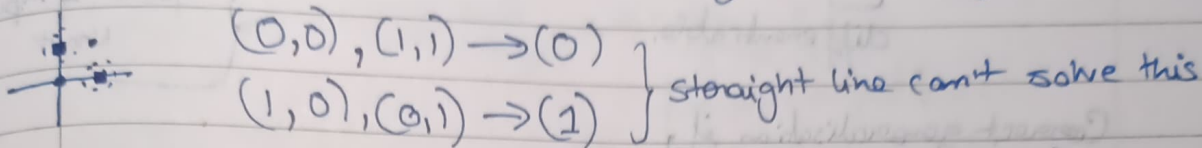
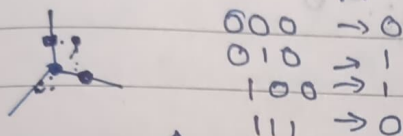


XOR problem (internal working)

- single layer perceptron can only solve linearly separable

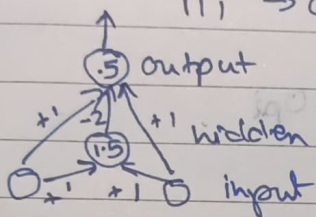


- solvable when mapped to high-D (add one more input)



- hidden layer introduced

$$\begin{pmatrix} \text{feat 1} & A \vee B \\ \text{feat 2} & A \& B \end{pmatrix}$$



$$\begin{aligned} \text{XOR} &= (A \vee B) \wedge \neg(A \wedge B) \\ &= (A \vee B) \wedge \neg(A \wedge B) \end{aligned}$$

Delta Rule - only for single layer (linear) \rightarrow act.

$$w_i = w_i + \eta(d - y)x_i$$

Generalised Δ rule:

$$E_p = \frac{1}{2} \sum_j (t_{pj} - o_{pj})^2$$

$$-\frac{\partial E_p}{\partial w_{ij}} = \delta_{pj} i_p$$

$$= \frac{\partial E_p}{\partial o_{pj}} \times \frac{\partial o_{pj}}{\partial w_{ij}}$$

$$o_{pj} = \sum_i w_{ji} i_p \quad (\text{for linear unit})$$

$$\therefore \Rightarrow -\frac{\partial E_p}{\partial w_{ji}} = \delta_{pj} i_i$$

$$\text{finally } \frac{\partial E}{\partial w_{ji}} = \sum_p \frac{\partial E_p}{\partial w_{ji}}$$

- addition of hidden unit \rightarrow graph no longer concave.

$$\text{net}_{pj} = \sum_i w_{ji} O_{pi} \quad \text{--- (7)}$$

$$O_{pj} = \underbrace{f_j(\text{net}_{pj})}_{\text{differentiable.}}$$

Semilinear func.
 ↓
 whose output
 funct. is diff.
 funcⁿ of input
 unit

Correct generalisation if,

$$\Delta p w_{ji} \propto - \frac{\partial E_p}{\partial w_{ji}} \quad \text{--- (8)}$$

$$\frac{\partial E_p}{\partial w_{ji}} = \frac{\partial E_p}{\partial \text{net}_{pj}} \frac{\partial \text{net}_{pj}}{\partial w_{ji}}$$

$$\text{--- (7)} \rightarrow \frac{\partial \text{net}_{pj}}{\partial w_{ji}} = \frac{\partial}{\partial w_{ji}} \sum_k w_{jk} O_{pk} = O_{pi}$$

$$- \frac{\partial E_p}{\partial w_{ji}} = \delta_{pj} O_{pi}$$

$$\text{--- (8)} \rightarrow \Delta p w_{ji} = \eta \delta_{pj} O_{pi}$$

↳ what should for each unit w_{ji}

$$\text{for } \delta_{pj} = - \frac{\partial E}{\partial \text{net}_{pj}}, \text{ apply partial} = - \frac{\partial E_p}{\partial \text{net}_{pj}} = \frac{\partial E_p}{\partial O_{pj}} \frac{\partial O_{pj}}{\partial \text{net}_{pj}} \quad \text{--- (12)}$$

$$= f'_j(\text{net}_{pj})$$

• Two cases

① y_j is the output of network

$$\frac{\partial E_p}{\partial O_{pj}} = -(t_{pj} - O_{pj}) \quad \text{--- (12)} \rightarrow \delta_{pj} = (t_{pj} - O_{pj}) f'_j(\text{net}_{pj})$$

② Not

$$\sum_k \frac{\partial E_p}{\partial \text{net}_{pk}} \frac{\partial \text{net}_{pk}}{\partial O_{pj}} = \sum_k \frac{\partial E_p}{\partial \text{net}_{pk}} \frac{\partial}{\partial \text{net}_{pk}} \sum_i w_{ki} O_{pi} = \sum_k \frac{\partial E_p}{\partial \text{net}_{pk}} w_{kj} = \sum_k \delta_{pk} w_{kj}$$

$$\delta_{pj} = f'_j(\text{net}_{pj}) \sum_k \delta_{pk} w_{kj}$$

$$\Delta w_{ji} = \eta \delta p_j o_{pi}$$

$$\delta p_j = (t_{pj} - o_{pj}) f'_j(\text{net } p_j) \quad (u \text{ output unit})$$

$$\delta p_j = f'_j(\text{net } p_j) \sum_k \delta p_k w_{kj} \quad (u \text{ not output})$$

* 2 phases → 2nd phase involves composing backward pass.

input propagate to find output o_{pj} for each unit

compared with target → error signal δp_j for each input unit.

- Compute δ for penultimate layer. This propagate the error back one layer, same process for every layer.

Logistic

$$o_{pj} = \frac{1}{1 + e^{-(\sum_i w_{ji} o_{pi} + \theta_j)}} \quad \text{net } p_j = \sum_i w_{ji} o_{pi} + \theta_j$$

$$\frac{do_{pj}}{d \text{net } p_j} = o_{pj}(1 - o_{pj})$$

$$\Rightarrow \delta p_j = (t_{pj} - o_{pj}) o_{pj}(1 - o_{pj}) \quad [\text{output unit } u]$$

$$= o_{pj}(1 - o_{pj}) \sum_k \delta p_k w_{kj} \quad [\text{hidden unit } u]$$

max for $0.5 = o_{pj}$
min for $o_{pj} = 0, 1$

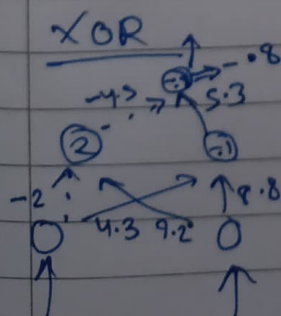
Learning rate

- keep large without oscillation

Momentum → helps

$$\Delta w_{ji}(n+1) = \eta (\delta p_j o_{pi}) + \alpha \Delta w_{ji}(n)$$

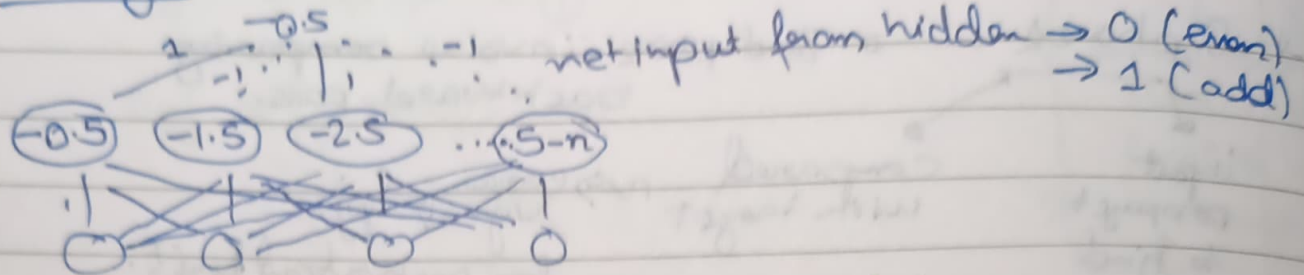
- filter high freq. variations in weight surface.
- sharp curvature tends to cause oscillation
- filter out high curvatures & allows effective weight to be larger.



Symmetry breaking

- start at equal value
- equal error signal
- stuck at local maximum

Polarity



- not readily discovered by unsupervised learning scheme

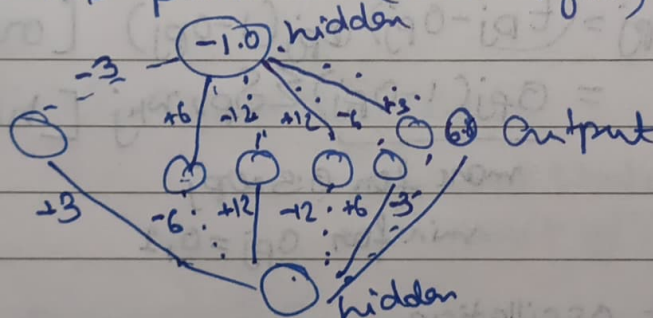
Encoding problem

- $\log_2 N$ hidden unit force to learn in binary.
- tendency of hidden unit to take extreme value.

But when req. to learn \rightarrow can learn \rightarrow (logistic activation?)

Symmetric problem

- hidden unit receives zero act. from input (symm.)
- atleast one positive (non-sym.)



Eqⁿ of backprop

$$\delta^L = \nabla_a C \odot \sigma'(z^L) = \frac{\partial C}{\partial o_j^L} \sigma'(z_j^L) = \frac{\partial C}{\partial z_j^L} = \frac{\partial C}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L}$$

$$\delta^L = ((w^{L+1})^T \delta^{L+1}) \odot \sigma'(z^L)$$

$$\frac{\partial C}{\partial b_j^L} = \delta_j^L$$

$$\frac{\partial C}{\partial w_{jk}^L} = a_k^{L-1} \delta_j^L$$

Proof for BP2 $\rightarrow \delta_j^L = \frac{\partial C}{\partial z_j^L} = \sum_k \frac{\partial C}{\partial z_k^{L+1}} \frac{\partial z_k^{L+1}}{\partial z_j^L}$

$$= \sum_k \frac{\partial z_k^{L+1}}{\partial z_j^L} \delta_k^{L+1}$$

$$z_k^{L+1} = \sum_j w_{kj}^{L+1} a_j^L + b_k^{L+1} = \sum_j w_{kj}^{L+1} \sigma(z_j^L) + b_k^{L+1}$$

$$\frac{\partial z_k^{L+1}}{\partial z_j^L} = w_{kj}^{L+1} \sigma'(z_j^L)$$

finally, $\delta_j^L = \sum_k w_{kj}^{L+1} \delta_k^{L+1} \sigma'(z_j^L)$