

# NNet & Deep Learning

## (model)

Perceptron = node with several binary inputs

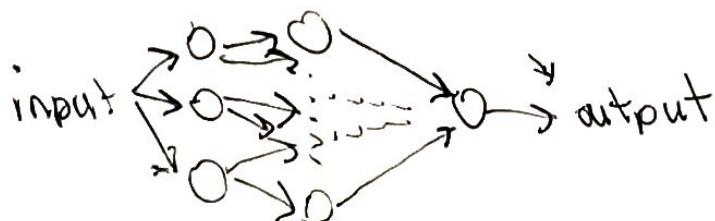
$x_1, x_2, \dots, x_n$  which produces a single binary output  $y$



where each input,  $x_j$ , a uni-directional weighted edge

$$y = \begin{cases} 0 & \sum_i w_i x_i \leq T \\ 1 & \sum_i w_i x_i > T \end{cases} \quad \text{where } T \text{ is the threshold value}$$

∴ a trivial NNet can be modeled as such:



(dot product)

Let  $\sum_j w_j x_j = w \cdot x$  where

$w$  and  $x$  are vectors and

to represent the threshold as a bias where  $b = -\text{threshold}$

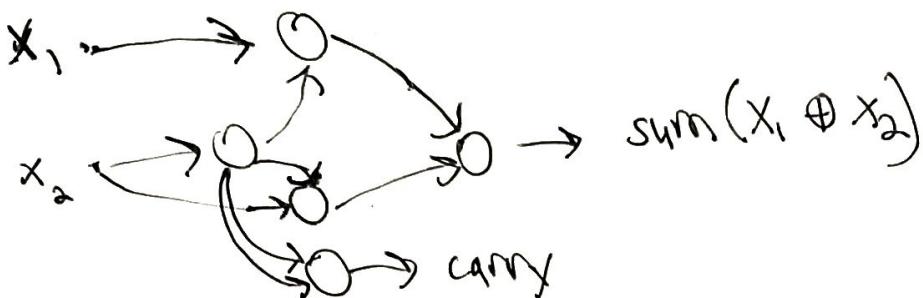
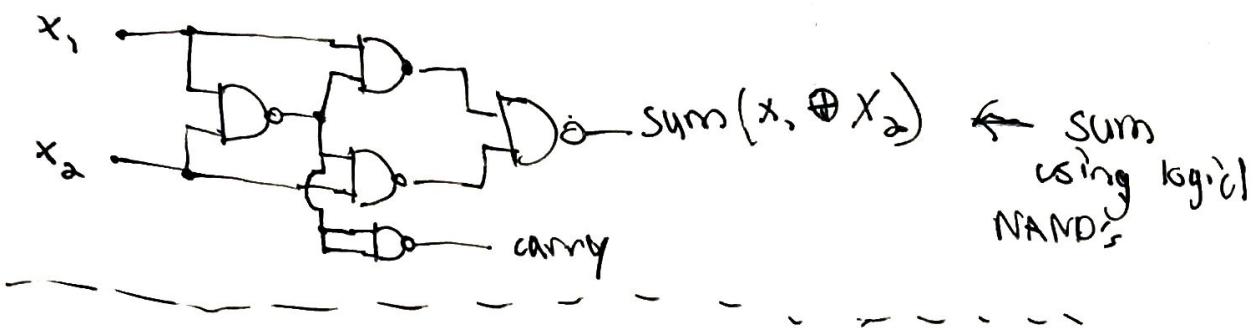
Now

$y$	$=$	$\begin{cases} 0, & w \cdot x + b \leq 0 \\ 1, & w \cdot x + b > 0 \end{cases}$
(output)		X

$$Y = \{y_i \mid y_i = y\}$$

henceforth,  
bias will represent  
-threshold

## NAND gate with perceptrons



each weighted edge value is -2 and a bias of 3

Just as in hardware architecture NAND's are universal,  
so are perceptrons for computation

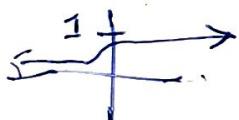
Sigmoid Neurons - same as perceptron except the output

$$x_1 \rightarrow O \rightarrow y \quad \text{where } y = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$z = (w \cdot x + b)$$

$$\Rightarrow \sigma(w \cdot x + b)$$

$\sigma$  is sigmoid function

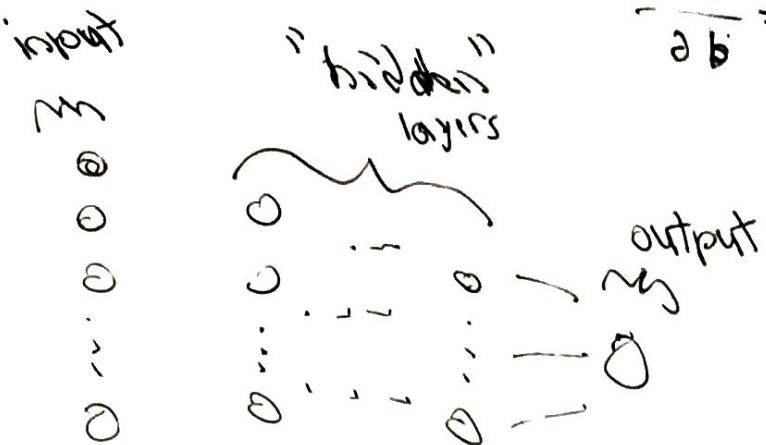


$$\frac{1}{1 + e^{-(\sum_j w_j x_j - b)}} \quad \text{or}$$

i.e. output is always in range of  $[0, 1]$

$$\therefore \Delta \text{output} \approx \sum_j \frac{\partial \text{output}}{\partial w_j} \Delta w_j + \frac{\partial \text{output}}{\partial b} \Delta b$$

$\frac{\partial d}{\partial b}$  = partial derivative of  $d$  with respect to  $b$



"feedforward"-network where output from one layer is used as input for the next

"recurrent" network - loops are allowed

### Gradient Descent

Let  $x = \begin{bmatrix} a_1 \\ \vdots \\ a_m \end{bmatrix}$  where  $x$  represents the vectorized  $28 \times 28$  image and  $a_i$  is the greyscale value of the pixel.

Let  $y = y(x)$  where  $y = \begin{bmatrix} a_1 \\ \vdots \\ a_{10} \end{bmatrix}$  representing the 10-dim. output vector

cost function:

$$C(w, b) = \frac{1}{2n} \sum_x \|y(x) - a\|^2$$

where  $w$  = weights

$b$  = biases

$n$  = # of training inputs

$a$  = output vector containing the correct digit value

$\|v\|$  = length of vector

$C(w, b)$  aka MSE or mean squared error

→ The closer  $C(w, b)$  is to zero the more correct the network is.

∴ goal of any nnet is to minimize the cost function & done via gradient descent

i.e. "learning" means to find the weights and biases that minimize the quadratic cost function

- can't calculate by hand b/c finding derivative of a function with thousands of variables isn't worth it.

- So the nnet will do this by find the slope (or gradient) of the function and determine the min through many iterations

$$\Delta C \approx \frac{\partial C}{\partial v_1} \Delta v_1 + \frac{\partial C}{\partial v_2} \Delta v_2 \quad \text{where } \Delta v = (\Delta v_1, \Delta v_2)^T$$

$$\Delta C = \left( \frac{\partial C}{\partial v_1}, \frac{\partial C}{\partial v_2} \right)^T$$

$$\Delta C \approx \nabla C \cdot \Delta V, \quad \Delta V = -\eta \nabla C$$

$$\Delta C \approx \nabla C \cdot -\eta \nabla C = -\eta \|\nabla C\|^2$$

$\therefore$  the gradient is  $v \rightarrow v' = v - \eta \nabla C$

i.e. compute gradient repeatedly going in the negative direction

Note:  $\eta$  needs to be small enough s.t.  $\Delta C \leq 0$ , otherwise the network will "go uphill" and go away from local min.

$(w_k) \quad (b_i)$

gradient descent - find weights & biases which minimize

$$\text{the cost of } C(w, b) = \frac{1}{2n} \sum_x \|y(x) - a\|^2$$

$$w_k \rightarrow w'_k = w_k - \eta \frac{\partial C}{\partial w_k}$$

$$b_i \rightarrow b'_i = b_i - \eta \frac{\partial C}{\partial b_i}$$

} update rule for  $w, b$

Stochastic Gradient Descent - estimate gradient of  $\nabla C$  by instead computing  $\nabla C_x$  for a small sample of inputs that are randomly chosen. Averaging over this sample will give a good estimate of the true gradient  $\nabla C$ .

for a random sample of inputs  $x_1, x_2, \dots, x_m$  given  $m$  is a large enough sample size

$$\frac{\sum_{j=1}^m \nabla C_{x_j}}{m} \approx \frac{\sum_x \nabla C_x}{n} = \nabla C$$

$$\Rightarrow \nabla C \approx \frac{1}{m} \sum_{j=1}^m \nabla C_{x_j}$$

in stoch. gradient descent

proof:

claim: given a large enough sample size  $m$ , stochastic gradient descent estimates any  $\nabla C_x$  such that

$$\nabla C_x \approx \nabla C$$

∴ Stoch. grad. descent in a neural network will look like: Let  $w_k, b_i$  be the weights & biases

$$w_k \rightarrow w_k - \frac{\eta}{m} \sum_j \frac{\partial C_{x_j}}{\partial w_k}$$

$$b_i \rightarrow b_i - \frac{\eta}{m} \sum_j \frac{\partial C_{x_j}}{\partial b_i}$$

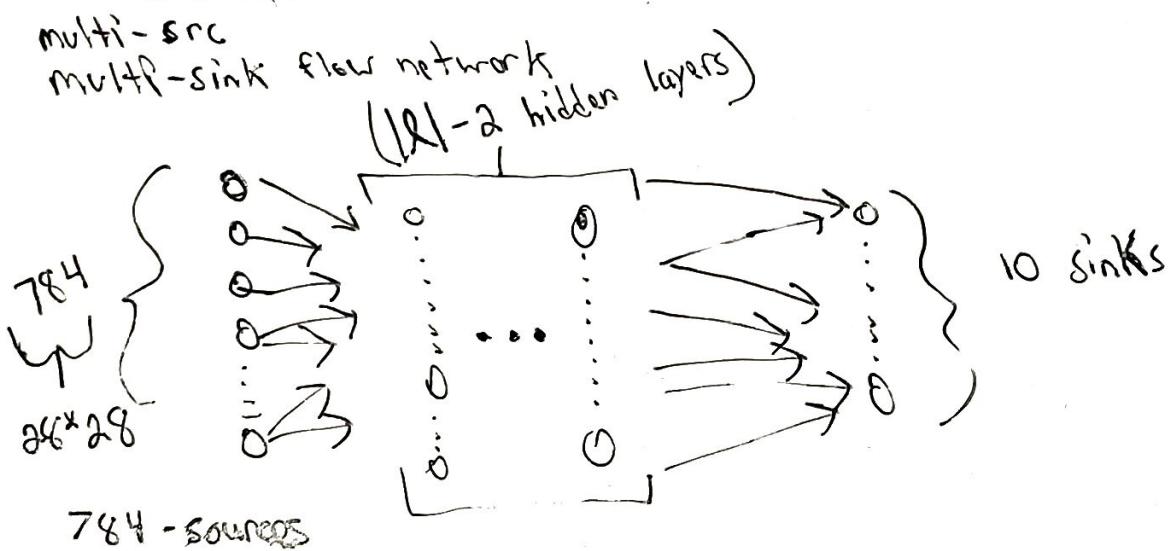
epoch - Let  $I$  be a randomly sorted list of inputs. Perform SGD (stoch. grad. desc.) on all the subsets  $S$  such that  $I_S = I - S$  and  $I_S, S \in I$ . An epoch has occurred when  $|I_S| < |S|$ , therefore,  $I$  needs to be reshuffled.

## Conv. NNet implementation

Let  $l$  be the "layers" of NNet s.t.  $\{l_i | l_i \in l\}$ .

$l_i$  is a partition in a  $|l|$ -partitioned, edge-weighted flow network of sigmoid neuron nodes with no loops or self-loops. Let  $W$  be the matrix of all the weights of the edges s.t.  $W = \begin{bmatrix} w_{11} & \dots & w_{1n} \\ \vdots & \ddots & \vdots \\ w_{m1} & \dots & w_{mn} \end{bmatrix}$ ,  $w_{ij} = \text{edge}(l_i, l_j)$ . Let  $\beta$  be the vector of all the biases of each sigmoid s.t.  $\beta = (b_1, b_2, \dots, b_n)^T$  where  $b_i = -t$ ,  $t$  is the threshold of the sigmoid neuron.

abstract:



initial values: seed each  $w_{ij} \in W$  and  $b_i \in \beta$  with a random scalar from the Gaussian distribution (normal dist).

$$N(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{(x-\mu)^2}{2\sigma^2}}$$

used  $\mu = 0$   
 $\sigma = 1$

$$\therefore \mathcal{W} = \left\{ [w_{ij}] \in \mathcal{W} \mid w_{ij} = \mathcal{N}(w_{ij} \mid 0, 1) \right\}$$

$$\mathcal{B} = \left\{ [b_i] \in \mathcal{B} \mid b_i = \mathcal{N}(b_i \mid 0, 1) \right\}$$

feed forward - Given an input  $\lambda_i$ ,  $ff(\lambda)$  "feeds" the input down the flow network

for every  $\lambda$  in  $\mathcal{L}$  s.t.  $\{\lambda \mid \lambda \in \mathcal{L}\}$

$$ff(\lambda) = \sigma(W(\lambda) \cdot \lambda + B(\lambda))$$

$\sigma$  is sigmoid func.,  $W(\lambda_i)$  is matrix of weights s.t.  $W(\lambda_i) = \left\{ [w_{ij}] \in \mathcal{W} \mid w_{ij} = \text{edge}(\overrightarrow{\lambda_i \lambda_{i+1}}) \right\}$   
 $\lambda$  is layer of nodes,  $B$  is vector of biases of  $\lambda$

Stochastic Gradient Descent - Given,  $TD$ , a list of 2-tuples  $(x, y)$ . corresponding to the input,  $x$ , and desired output,  $y$ .  $bs$  is the size of each subset of  $TD$ .  $\eta$  is learning rate. epoch is number of iterations ("time") to train for.



for  $i$  where  $i \in \{0, 1, 2, \dots, |\text{epoch}|\}$

- randomly order of the sets in TD

Let  $\mathcal{L} = [A_1, A_2, \dots, A_n]$  where  $n = |\text{TD}|$  and

$$a_i = \{a_i \in \mathcal{L} \mid a_i = (x, y)\}$$

$$A = [(x, y)_k, (x, y)_{k+1}, \dots, (x, y)_{k+s}]$$

$$\text{where } k = \{0, 1|m_b|, 2|m_b|, \dots, \frac{n}{m_b}|m_b|\}$$

$\mathcal{L}_e$  is the list of subsets of training data for each epoch

for each  $T$ , subset in  $\mathcal{L}$

| update-batch( $T$ ,  $\eta$ ) where  $\eta$  is learning rate  
end

- if you want to evaluate the net with  $t$ , test data  
| evaluate( $t$ ) for each  $t$   
end

end

Update Batch - Given a subset  $\mathcal{V} \subseteq \mathcal{D}$ , apply gradient descent  
updating  $(\mathbf{w})$  and  $\mathbf{b}$ .  $\eta$  is learning rate

for each  $(x, y)$ -2-tuple  $\in \mathcal{V}$

$$\Delta(\nabla w) = \text{backpropagation}(x)$$

$$\Delta(\nabla b) = \text{backpropagation}(y)$$

for each  $w_i \in \mathbf{w}$  and  $\nabla w_i \in \nabla W$

$$w_i = w_i - \left( \frac{\eta}{|\mathcal{V}|} \right) \cdot \nabla w_i$$

for each  $b_i \in \mathbf{b}$  and  $\nabla b_i \in \nabla \mathbf{b}$

$$b_i = b_i - \left( \frac{\eta}{|\mathcal{V}|} \right) \cdot \nabla b_i$$

## Back propagation

Let  $w_{jk}^l \in W$  s.t.  $w_{jk}^l = \text{weight of edge}(j, k)$  in layer  $L$

Let  $b_j^l \in \beta$  s.t.  $b_j^l = -\text{threshold of } j^{\text{th}} \text{ neuron}$  in layer  $L$

Let  $\alpha_j^l$  be the activation value of the  $j^{\text{th}}$  neuron in layer  $L$  s.t.

$$\alpha_j^l = \sigma \left( \sum_k w_{jk}^l \cdot \alpha_k^{l-1} + b_j^l \right)$$

or in vectorized form  $\uparrow$  remember  $(Wa + \beta)$

$$\alpha^l = \sigma (w^l \alpha^{l-1} + \beta^l)$$

$$\sigma \left( \begin{bmatrix} \vdots & \vdots \\ \vdots & \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ i \end{bmatrix} + \begin{bmatrix} \vdots \\ i \end{bmatrix} \right)$$

Let  $Z^l \equiv w^l \alpha^{l-1} + \beta^l$  be "weight input to the  $L$  layer"

goal of back prop is to compute partial derivatives,

$$\frac{\partial C}{\partial w} \text{ and } \frac{\partial C}{\partial b}$$

$$C = \frac{1}{2n} \sum_x \|y(x) - \alpha^l(x)\|^2$$

$n = \# \text{ of } x$

remember  $C \approx \frac{1}{n} \sum_x C_x$

∴ it can be assumed that computing  $\frac{\partial C_x}{\partial w}$ ,  $\frac{\partial C_x}{\partial b}$  for a single subset of training input,  $\frac{\partial C}{\partial w}$ ,  $\frac{\partial C}{\partial b}$  can be found by

$$C \approx \frac{1}{n} \sum C_x$$

### Hadamard product - (soft)

$$\text{soft} = \begin{bmatrix} i \\ i \end{bmatrix} \odot \begin{bmatrix} m \\ n \end{bmatrix} = \begin{bmatrix} i \cdot m \\ j \cdot n \end{bmatrix}$$

Let  $\delta_j = \frac{\partial C}{\partial z_j^l}$   $\delta_j^l$  is the error in  $j^{th}$  neuron in  $l^{th}$  layer. So instead of computing

- \*  $\rightarrow \sigma(z_j^l)$ , the neuron outputs
- \*  $\rightarrow \sigma(z_j^l + \Delta z_j^l)$  affecting the cost by  $\left( \frac{\partial C_j}{\partial z_j^l} \Delta z_j^l \right)$

- So back prop is a way of calculating  $\delta^l$  for each layer in relation to

$$\frac{\partial C}{\partial w_{jk}^l} \text{ and } \frac{\partial C}{\partial b_j^l}$$

\* remember  $\Rightarrow \delta = \frac{\partial C}{\partial z_j^l}$

$$z^l = w^l \alpha^{l-1} + b^l$$

$$\therefore \alpha_j^l = \sigma(z^l)$$

$$\Rightarrow \delta_j^L = \frac{\partial C}{\partial \alpha_j^L} \sigma'(\alpha_j^L)$$

↓                      ↓

how fast cost is changing as a func. of  $j^{th}$  activation. i.e. how much  $\alpha$  depends on the value of  $C$  to activate

\*  $\sigma'$  is derivative of sigmoid function

in linear algebra, the above, element-wise, operation can be simplified to

$$\boxed{\delta^L = \nabla_{\alpha} C \odot \sigma'(\alpha^L)}$$

$$\nabla_{\alpha} C = \begin{bmatrix} a_0 \\ \vdots \\ a_n \end{bmatrix}, \left\{ a_i \in \nabla_{\alpha} C \mid a_i = \frac{\partial C}{\partial \alpha_i^L} \right\}$$

$$\Rightarrow \delta^L = (\alpha^L - y) \odot \sigma'(\alpha^L)$$

Let  $K = L+1$

$$\delta^L = ((w^K)^T \cdot \delta^K) \odot \sigma'(\alpha^L)$$

given we know  $w^{L+1}$  and  $\delta^{L+1}$ , we can calculate  $\delta^L$ . i.e. "moving the error back through the network" or "back" propagation.

So  $\delta_j^l$  can be used to calculate  $\delta^l$  then  $\delta^l$  will be used for  $\delta^{l-1}$  and so on working backward.

The conclusion can also be drawn that:

$$\frac{\partial C}{\partial b_j^l} = \frac{\partial C}{\partial \delta_j^l} \cdot \sigma'(z_j^l)$$

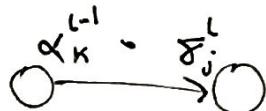
rate of change of cost  
with respect to any bias

$$\Rightarrow \delta_j^l = \frac{\partial C}{\partial b_j^l}$$

$$\frac{\partial C}{\partial w_{jk}^l} = \alpha_k^{l-1} \cdot \delta_j^l$$

rate of change of cost with  
respect to any weight

abstractly:



$\alpha$  is activation of a neuron to the weight  
and  $\delta$  is the error of the neuron's output  
from the weight

- When activation value is closer to zero, the less the gradient will change which means the neuron will "learn" more slowly
- remember also the  $\sigma'$  is derivative of sigmoid  $\sigma$   
so  $\lim_{n \rightarrow 1} \sigma'(n) = 0$ . Therefore, it also "learns" slowly when activation is high ( $\sigma(z) = \text{near one}$ ).  
So when activation is near 0 or 1, neuron will "stop" learning

A further conclusion can be drawn. Any weight that is an input to a "saturated neuron" (a neuron whose activation is near 0 or 1), will also learn slowly.

Corollary

proof: definition of  $\delta^L$

∴ Any weight  $w_j^L$  will "learn slowly" if either input neuron has a near 0 activation or the output neuron is saturated.

So the 4 functions that back propagation hinge on are:

$$1) \delta^L = \nabla_C C \odot \sigma'(z^L)$$

$$2) \delta^L = ((w^{L+1})^T \cdot \delta^{L+1}) \odot \sigma'(z^L)$$

$$3) \frac{\partial C}{\partial b_j} = \delta_j^L$$

$$4) \frac{\partial C}{\partial w_{jk}^L} = \alpha_k^{L-1} \cdot \delta_j^L$$

In essence this is just applying the chain rule in order to approximate closely the gradient of the cost function

Back propagation - compute the gradient of the cost function

- set the activation,  $\alpha$ , for the input layer

- "feed forward" - for each layer,  $l$ ,

$$z^l = w^l \alpha^{l-1} + b^l$$

$$\alpha^l = \sigma(z^l)$$

- "output error"

$$\delta^L = \nabla_{\alpha} C \odot \sigma'(z^L)$$

- back prop the error, for each layer  $l = L-1, L-2, \dots$

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

- gradient descent, update weights and biases s.t.

$$\frac{\partial C}{\partial w_{jk}^l} = \alpha_k^{l-1} \delta_j^l$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l \quad , \text{ for each layer}$$