

Deep Neural Networks

Week 3

Content

Deep Neural Networks

Neural Network Notation

Forward Propagation

Backward Propagation

Hyperparameters

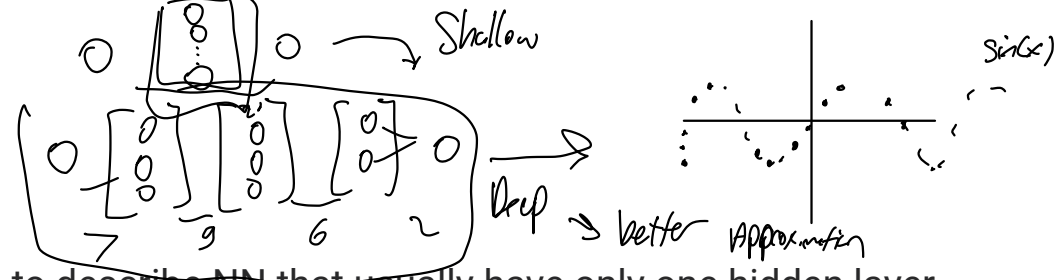
Universal Approximation Theorem



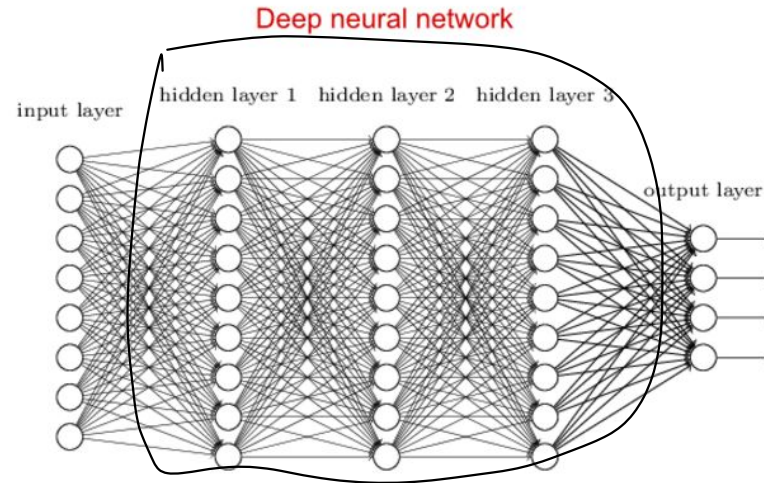
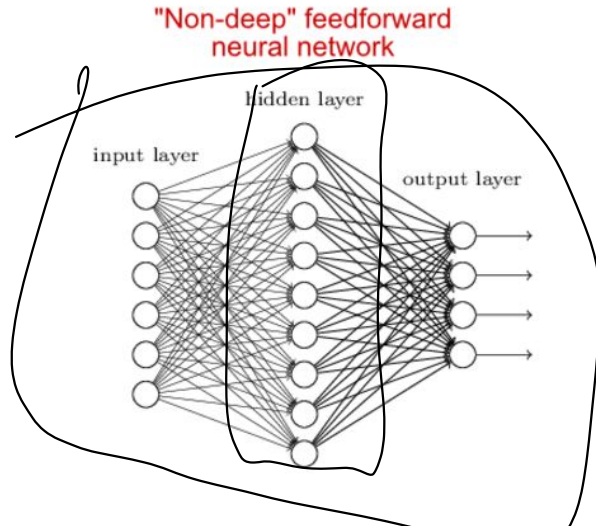
Deep Neural Networks



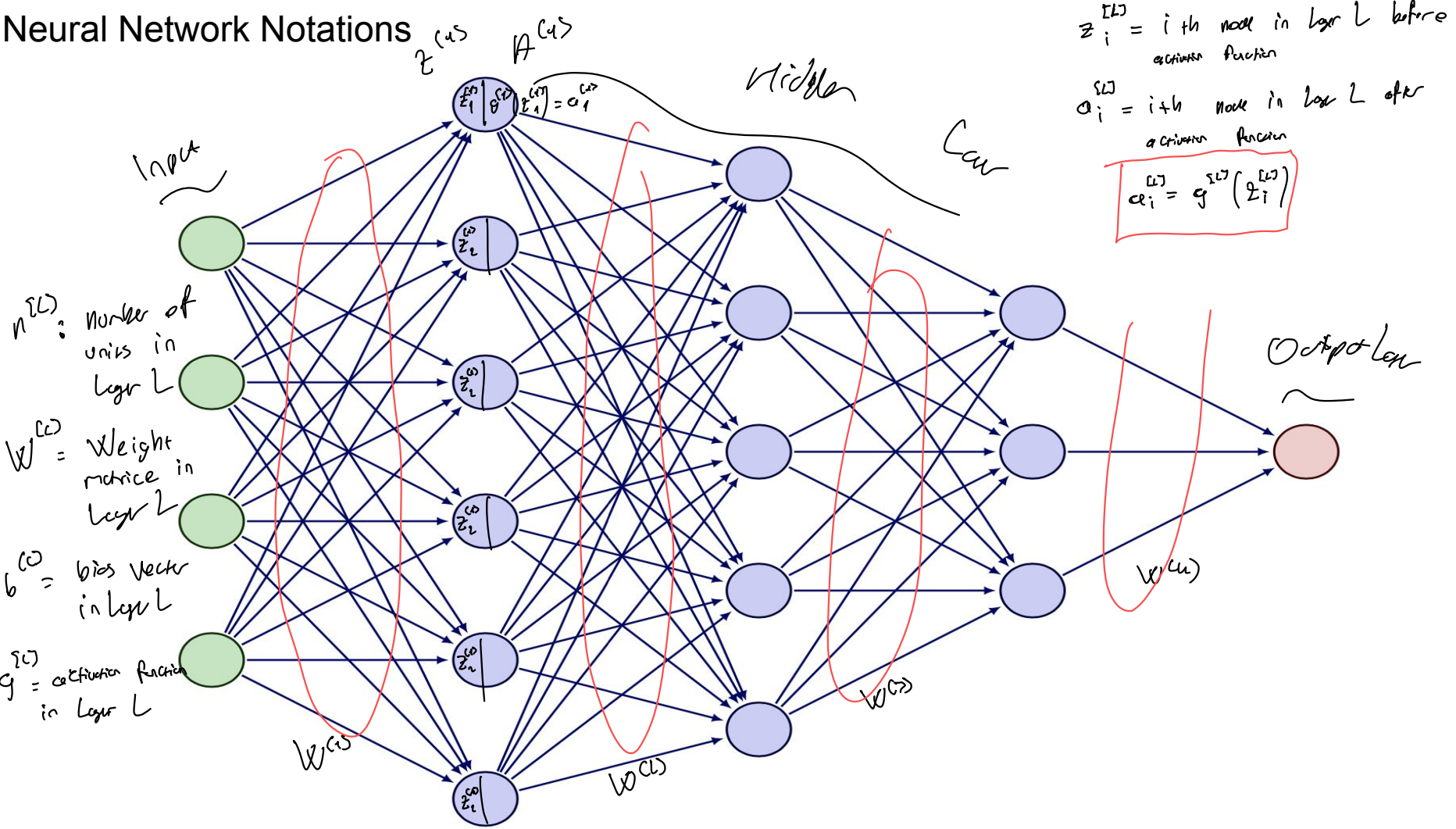
Deep Neural Networks



- Shallow neural networks is a term used to describe NN that usually have only one hidden layer while the term deep neural networks is used to describe NN that have several hidden layers.
- The deep NN with the right architectures achieve better results than shallow ones that have the same computational power.

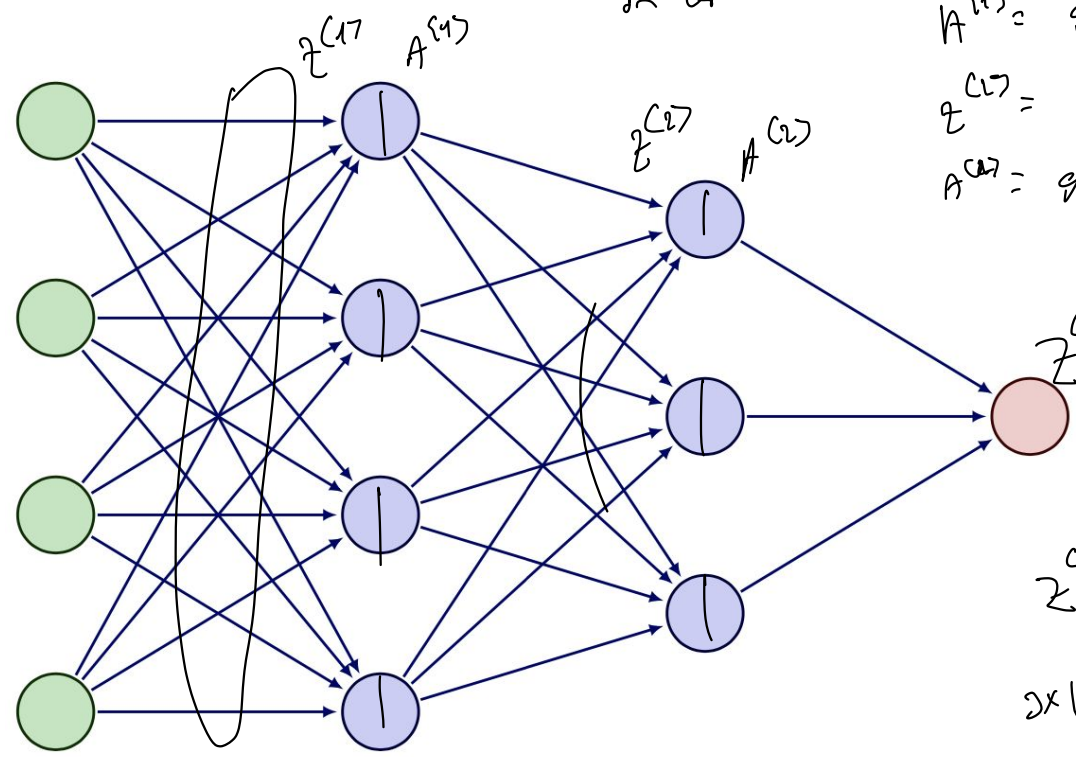


Neural Network Notations



Forward Propagation

$g = \text{activation func on } \text{vec } L$



Assume $A^{(0)}$ Shape = $(4, 1)$
↓
one sample

$$z^{(1)} = W^{(1)} A^{(0)} + b^{(1)}$$
$$A^{(1)} = g(z^{(1)})$$
$$z^{(2)} = W^{(2)} A^{(1)} + b^{(2)}$$
$$A^{(2)} = g(z^{(2)})$$

$W^{(1)} \text{ shape} = \begin{pmatrix} n^{(1)} & n^{(0)} \end{pmatrix}$
 $z^{(1)} \text{ shape} = \begin{pmatrix} n^{(1)} & m \end{pmatrix}$
 $A^{(1)} \text{ shape} = \begin{pmatrix} n^{(1)} & m \end{pmatrix}$
 $b^{(1)} \text{ shape} = \begin{pmatrix} n^{(1)} & 1 \end{pmatrix}$

$$z^{(1)} = W^{(1)} A^{(0)} + b^{(1)}$$

$\begin{matrix} 4 \times 4 & 4 \times 1 & + & 4 \times 1 \\ \text{---} & \text{---} & & \text{---} \\ 4 \times 1 & = & 4 \times 1 & + & 4 \times 1 \end{matrix}$

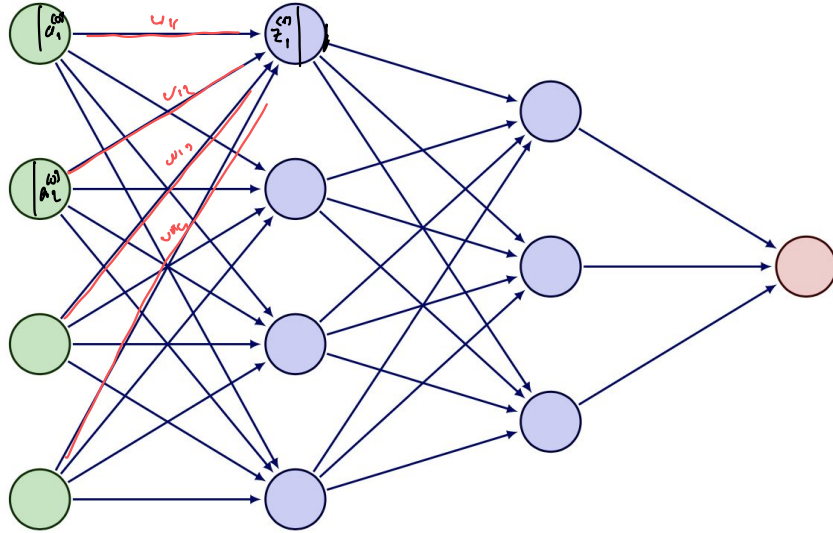
$$z^{(2)} = W^{(2)} A^{(1)} + b^{(2)}$$

$\begin{matrix} 2 \times 4 & 4 \times 1 & + & 2 \times 1 \\ \text{---} & \text{---} & & \text{---} \\ 2 \times 1 & = & 2 \times 1 & + & 2 \times 1 \end{matrix}$

$$z^{(3)} = W^{(3)} A^{(2)} + b^{(3)}$$

$\begin{matrix} 1 \times 2 & 2 \times m & + & 1 \times 1 \\ \text{---} & \text{---} & & \text{---} \\ 1 \times m & = & 1 \times m & + & 1 \times m \end{matrix}$
 $z^{(3)} = 1 \times m$

Forward Propagation



$$[A^{(1)}] =$$

$$\begin{aligned} a_1^{(1)} &= g(z_1^{(1)}) = \left(w_{11}a_1^{(0)} + w_{12}a_2^{(0)} + w_{13}a_3^{(0)} + w_{14}a_4^{(0)} + b_1^{(1)} \right) g^{(1)} \\ a_2^{(1)} &= g(z_2^{(1)}) = \left(w_{21}a_1^{(0)} + w_{22}a_2^{(0)} + w_{23}a_3^{(0)} + w_{24}a_4^{(0)} + b_2^{(1)} \right) g^{(1)} \\ a_3^{(1)} &= g(z_3^{(1)}) = \left(w_{31}a_1^{(0)} + w_{32}a_2^{(0)} + w_{33}a_3^{(0)} + w_{34}a_4^{(0)} + b_3^{(1)} \right) g^{(1)} \\ a_4^{(1)} &= g(z_4^{(1)}) = \left(w_{41}a_1^{(0)} + w_{42}a_2^{(0)} + w_{43}a_3^{(0)} + w_{44}a_4^{(0)} + b_4^{(1)} \right) g^{(1)} \end{aligned}$$

$$w_{11}a_1^{(0)} + w_{12}a_2^{(0)} + w_{13}a_3^{(0)} + w_{14}a_4^{(0)} + b_1^{(1)} = z_1^{(1)}$$

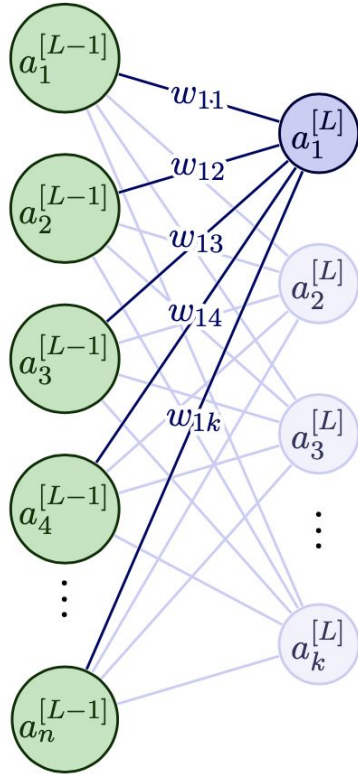
$$w_{21}a_1^{(0)} + w_{22}a_2^{(0)} + w_{23}a_3^{(0)} + w_{24}a_4^{(0)} + b_2^{(1)} = z_2^{(1)}$$

$$w_{31}a_1^{(0)} + w_{32}a_2^{(0)} + w_{33}a_3^{(0)} + w_{34}a_4^{(0)} + b_3^{(1)} = z_3^{(1)}$$

$$w_{41}a_1^{(0)} + w_{42}a_2^{(0)} + w_{43}a_3^{(0)} + w_{44}a_4^{(0)} + b_4^{(1)} = z_4^{(1)}$$

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \end{bmatrix} \begin{bmatrix} a \\ \end{bmatrix} + b^{(1)} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix}$$

Forward Propagation

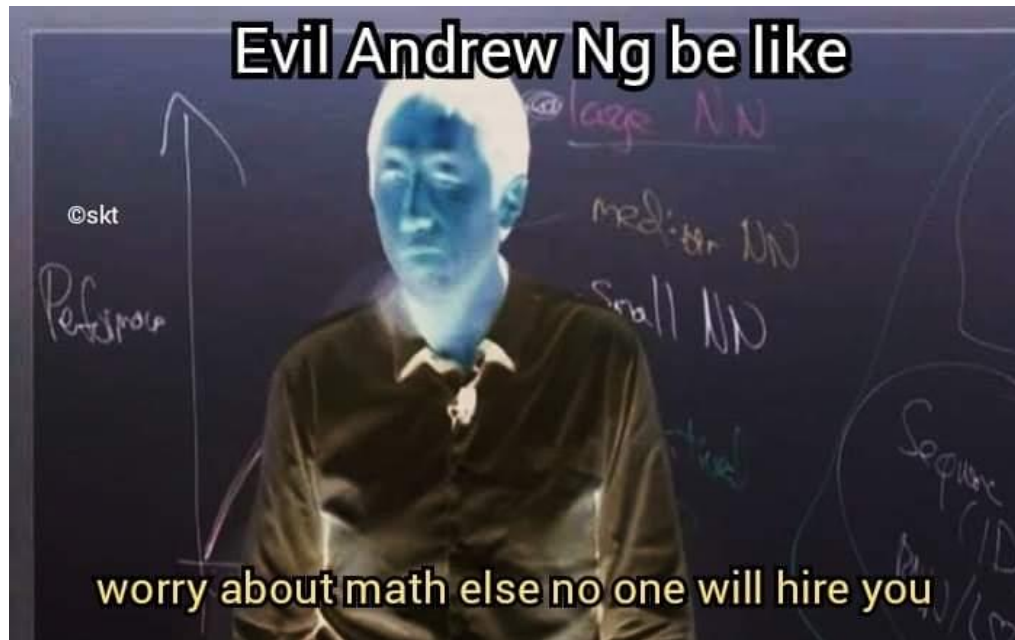


$$a_1^{[L]} = g(z_1^{[L]}) = \left(w_{11}^{[L]} a_1^{[L-1]} + w_{12}^{[L]} a_2^{[L-1]} + \dots + w_{1n}^{[L]} a_n^{[L-1]} \right)$$

$$\begin{aligned} a_k^{[L]} = g(z_k^{[L]}) &= g\left(w_{k1}^{[L]} a_1^{[L-1]} + w_{k2}^{[L]} a_2^{[L-1]} + \dots + w_{kn}^{[L]} a_n^{[L-1]}\right) \\ &= g\left(\sum_{i=1}^n w_{ki}^{[L]} a_i^{[L-1]}\right) \end{aligned}$$

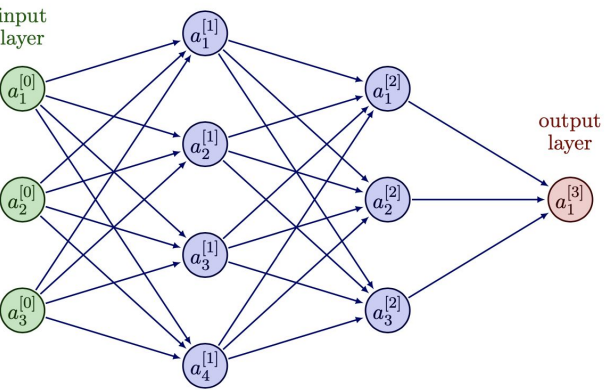
$$\begin{aligned} A^{[L]} &= W^{[L]} A^{[L-1]} + b^{[L]} \\ k \times n &= \underbrace{k \times n}_{\text{matrix}} \times \underbrace{n \times n}_{\text{matrix}} + k \times 1 \end{aligned}$$

$$\begin{aligned} a_k^{[L]} = g(z_k^{[L]}) &= g\left(w_{k1}^{[L]} a_1^{[L-1]} + w_{k2}^{[L]} a_2^{[L-1]} + \dots + w_{kn}^{[L]} a_n^{[L-1]}\right) \\ &= g\left(\sum_{i=1}^n w_{ki}^{[L]} a_i^{[L-1]}\right) \end{aligned}$$



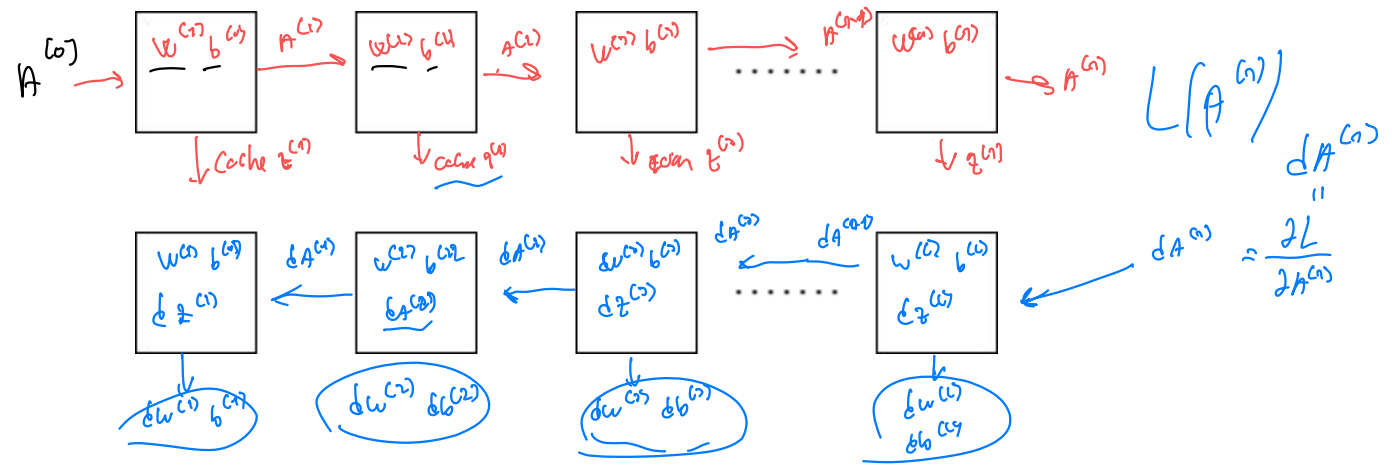
Backward Propagation

Backward Propagation



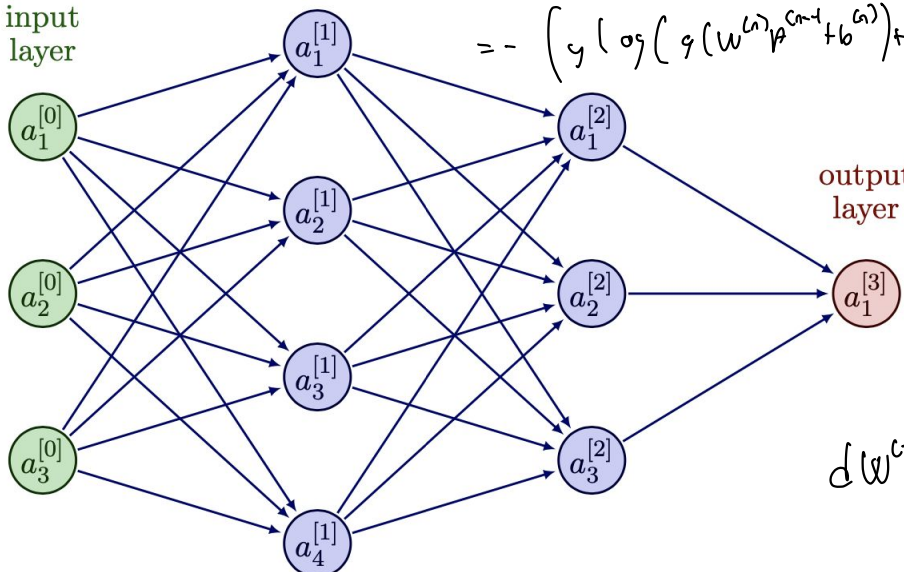
Forward Prop Input : $A^{(L-1)}$
 Output : $A^{(L)}$

Backward Prop Input : $dA^{(L)}$
 cache $\{t\}$
 Output : $dW^{(L)} db^{(L)} d\theta^{(L-1)}$



Backward Propagation

input layer



$$L(a^{(n)}) = - \left(y \log(a^{(n)}) + (1-y) \log(1-a^{(n)}) \right)$$

$$= - \left(y \log(\sigma(z^{(n)})) + (1-y) \log(1-\sigma(z^{(n)})) \right)$$

$$= - \left(y \log(\sigma(W^{(n)} a^{(n-1)} + b^{(n)})) + (1-y) \log(1-\sigma(W^{(n)} a^{(n-1)} + b^{(n)})) \right)$$

output layer

$$\delta z^{(L)} = \delta a^{(L)} \cdot \sigma'(z^{(L)})$$

$$\delta W^{(L)} = \delta z^{(L)} \cdot a^{(L-1)}$$

$$\delta b^{(L)} = \delta z^{(L)}$$

$$-\left(\frac{y}{a} + \frac{(1-y)}{1-a}\right) = \frac{d \sigma(x)}{dx} = \sigma'(x)$$

$$\frac{\partial(\sigma(z^{(n)}))}{\partial(z^{(n)})} = \sigma(z^{(n)}) (1 - \sigma(z^{(n)}))$$

$$\frac{\partial(z^{(n)})}{\partial(z^{(n)})} = \sigma'(z^{(n)})$$

$$\partial(\sigma(z^{(n)}))$$

$$\delta z = \frac{\partial L}{\partial z}$$

$$\delta z^{(n)} = \frac{\partial L}{\partial z^{(n)}} = \frac{\partial L}{\partial A^{(n)}} \cdot \frac{\partial A^{(n)}}{\partial z^{(n)}}$$

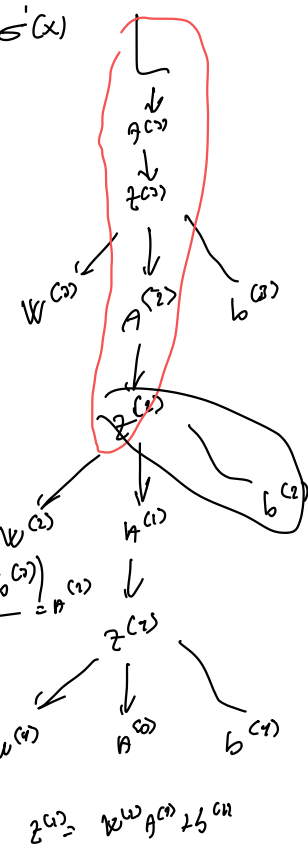
$$\frac{\partial L}{\partial A^{(n)}} \cdot \underbrace{\frac{\partial A^{(n)}}{\partial z^{(n)}}}_{\sigma'(z^{(n)})}$$

$$\delta W^{(n)} = \frac{\partial L}{\partial W^{(n)}} = \frac{\partial L}{\partial A^{(n)}} \cdot \frac{\partial A^{(n)}}{\partial z^{(n)}} \cdot \frac{\partial z^{(n)}}{\partial W^{(n)}}$$

$$\frac{\partial L}{\partial A^{(n)}} \cdot \frac{\partial A^{(n)}}{\partial z^{(n)}} \cdot \frac{\partial z^{(n)}}{\partial W^{(n)}}$$

$$\delta b^{(n)} = \frac{\partial L}{\partial b^{(n)}} = \frac{\partial L}{\partial A^{(n)}} \cdot \frac{\partial A^{(n)}}{\partial z^{(n)}} \cdot \frac{\partial z^{(n)}}{\partial b^{(n)}} \cdot \frac{\partial z^{(n)}}{\partial b^{(n)}}$$

$$\frac{\partial L}{\partial A^{(n)}} \cdot \frac{\partial A^{(n)}}{\partial z^{(n)}} \cdot \frac{\partial z^{(n)}}{\partial b^{(n)}} \cdot \frac{\partial z^{(n)}}{\partial b^{(n)}}$$

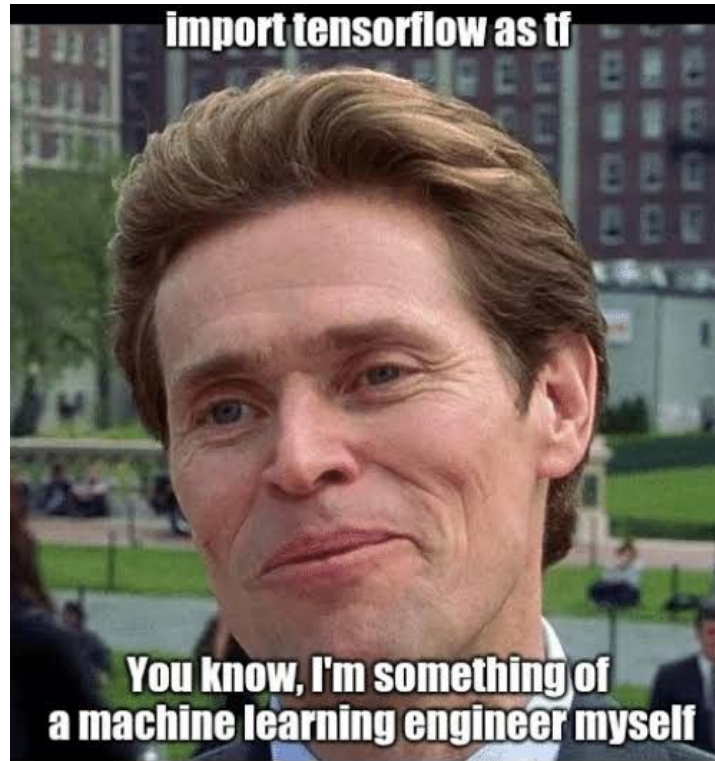


$$z^{(n)} = W^{(n)} a^{(n-1)} + b^{(n)}$$

$$f(x) = x$$

$$f'(x) = 1$$

Hyperparameters



Hyperparameters

Hyperparameters effect parameters

Hyperparameter examples:

- Learning Rate
- #Units
- #Iterations
- #Layers
- Batch size

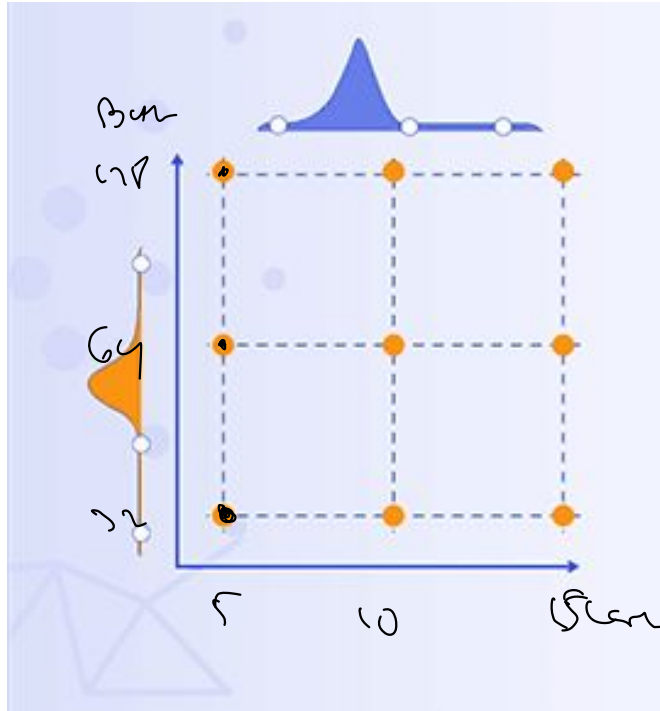
epoch

- data 500
- batch = 500
 - batch = 250
 - batch = 100
- 1 iteration : 1 epoch
 - 2 iterations : 1 epoch
 - 5 iterations : 1 epoch
-

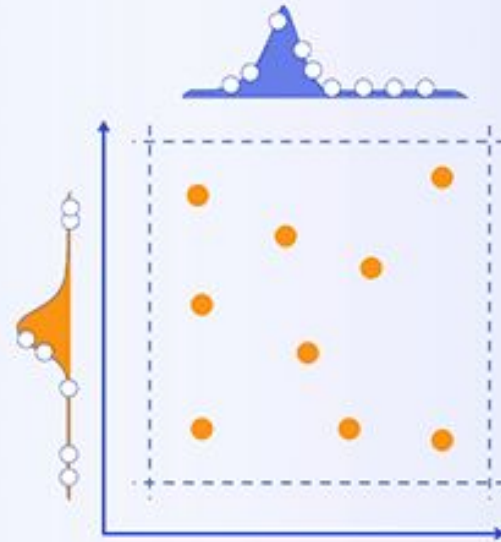
We can select hyperparameters using several methods

Hyperparameter Tuning

Grid Search



Random Search



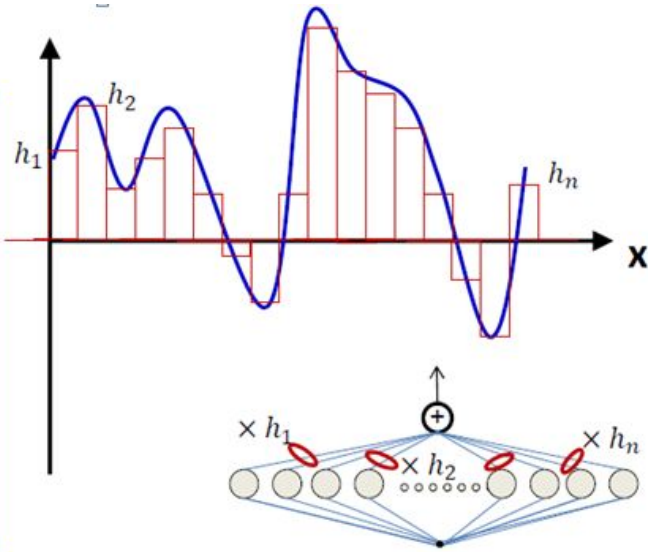
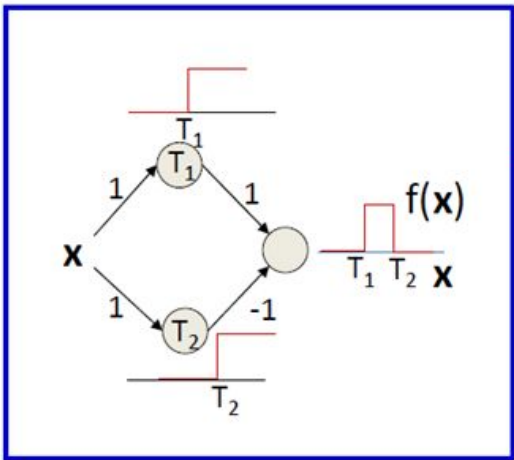
Universal Approximation Theorem

Some conditions \Rightarrow Borel measurable: \mathbb{R}^n
Subst \rightarrow closed \checkmark
Continuous \checkmark
Bounded \checkmark

The Universal Approximation Theorem means that regardless of what function we are trying to learn, we know that a large MLP will be able to represent this function

$$f: \mathcal{N} \rightarrow \mathbb{R}^n \Rightarrow f \text{ is onto}$$

\mathcal{N} : set of neural nets



**inzva: *brings the AI fellows
together***

inzva:

