# Deep Neural Networks

Week 3

# Content

Deep Neural Networks    — Shallow $NN's$

Neural Network Notation    $n^{[i]} =$    $W^{[i]} =$

Forward Propagation

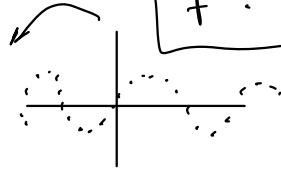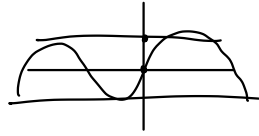Backward Propagation

Hyperparameters

- Universal Approximation Theorem

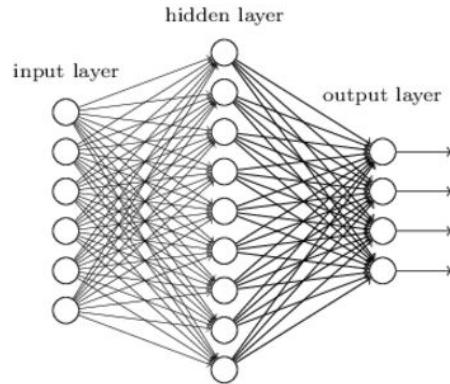# Deep Neural Networks

# Deep Neural Networks
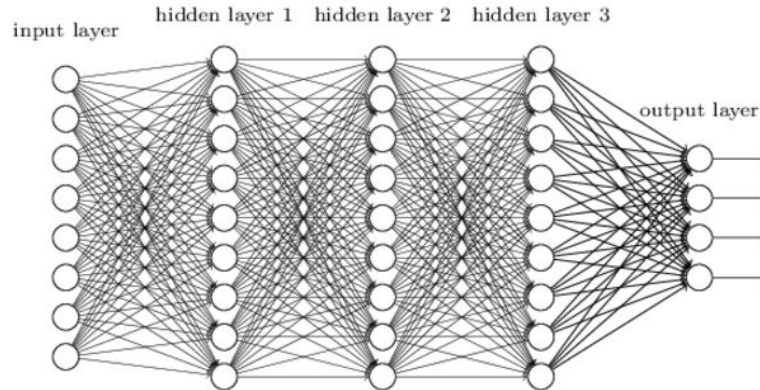
$f : X \longrightarrow y$

$f_3(f_2(f_1(x)))$

fix h

thu    scale

- 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.

"Non-deep" feedforward
neural network

input layer
hidden layer
output layer

Deep neural network

input layer    hidden layer 1   hidden layer 2   hidden layer 3
output layer

# Neural Network Notations



hidden Layer

Input

Output

$n^{[0]} = 3$

$n^{[1]} = 4$

$n^{[2]} = 1$

$n^{[L]} = \#$ of units in Lar L

$x_1$

$x_2$

$x_3$

$z_1^{(1)} \mid a_1^{(1)}$

$a_2^{(1)}$

$a_3^{(1)}$

$z_4^{(1)} \mid a_4^{(1)}$

$A^{(0)}$

$W^{(1)}$

$W^{(2)}$

$z^{(2)} \mid$

$z^{(2)}$

$A^{(2)}$

# Forward Propagation

$g^{[L]}$: action function in layer $L$



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

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

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

$A^{(0)} = (3, m)$

$n^{[0]}$ # sample

$A^{(1)} = g(z^{(1)})$
$4 \times 1$

$$z^{(1)} = W^{(1)} A^{(0)} + b^{(1)}$$
$\underbrace{4 \times 3}_{} \quad \underbrace{3 \times 1}_{} \quad \underbrace{}_{4 \times 1}$

$\boxed{4 \times 1} = 4 \times 1 + 4 \times 1 =$

$$z^{(1)} = W^{(1)} A^{(1)} + b^{(1)}$$
$\downarrow \quad 3 \times 4 \quad 4 \times 1 \quad 3 \times 1$
$3 \times 1 = \quad 3 \times 1 \quad + 3 \times 1$

$W^{[L]}$ . shape $= \left( n^{[L]}, n^{[L-1]} \right)$
$b^{[L]}$ . shape $= \left( n^{[L]}, 1 \right)$
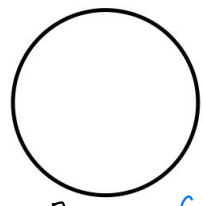$z^{[L]}$ . shape $=$
$A^{[L]}$ . shape $= ( n$

# Forward Propagation



$$Z_1^{[1]} = w_{11}X_1 + w_{12}X_2 + w_{13}X_3 + b_1^{[1]}$$

$$Z_2^{[1]} = w_{21}X_1 + w_{22}X_2 + w_{23}X_3 + b_2^{[1]}$$

$$Z_3^{[1]} = w_{31}X_1 + w_{32}X_2 + w_{33}X_3 + b_3^{[1]}$$

$$\begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} = \begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1n} \\ X_{21} & X_{12} & \cdots & X_{1m} \\ X_{31} & \cdots & \cdots & X_{3n} \end{bmatrix}$$

$(3,1)$

$w^{[1]}.\text{shape} = (4,3)$

$b^{[1]}.\text{shape} = (4,1)$

$$w_{11}X_1 + w_{12}X_2 + w_{13}X_3 + b_1^{[1]} = Z_1^{[1]}$$

$$w_{21}X_1 + w_{22}X_2 + w_{23}X_3 + b_2^{[1]} = Z_2^{[1]}$$

$$w_{31}X_1 + w_{32}X_2 + w_{33}X_3 + b_3^{[1]} = Z_3^{[1]}$$

$$w_{41}X_1 + w_{42}X_2 + w_{43}X_3 + b_4^{[1]} = Z_4^{[1]}$$

$z^{[1]}.\text{shape} = (4,m)$

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{34} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} = \begin{bmatrix} Z \end{bmatrix}$$

$(4,3) \quad (3,m) \quad | \quad (4,1)$

$$(4,m) + (4,1) \quad (4,m)$$

$$A^{[1]} = g\left( \begin{bmatrix} Z \end{bmatrix} \right)$$

$(4,m)$

# Forward Propagation

$L^{[U]} = U$th unit in Layer $L$



$w^{[L]}$. shape $= (n^{[L]}, n^{[L-1]})$

$w_{L^{[U]}, L-1^{[U]}}$

$w_{11} a_1^{[1]} + w_{12} a_2^{[1]} + w_{13} a_3^{[1]} + w_{14} a_4^{[1]} = z_1^{[2]}$

$w_{21} a_1^{[1]} + w_{22} a_2^{[1]} + w_{23} a_3^{[1]} + w_{24} a_4^{[1]} = z_1^{[2]}$

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \end{bmatrix} \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} + \begin{bmatrix} b_1^{[2]} \\ b_2^{[2]} \end{bmatrix} = \begin{bmatrix} z \end{bmatrix}$$
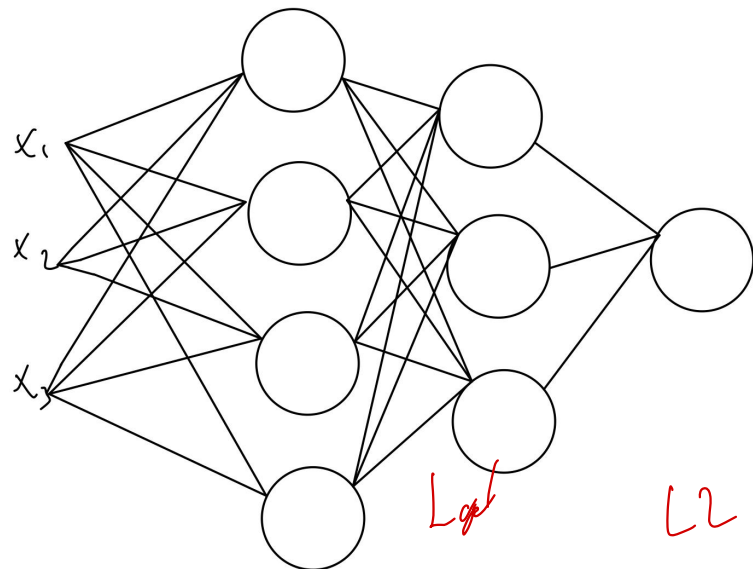
$(2, 4) \times (4, m)$

$(2, m) + (2, 1)$

$(2, m) = \begin{bmatrix} z \end{bmatrix}$ . shape

$w^{[2]}$ . shape $\rightarrow (2, 4)$

$b^{[2]}$ . shape $= (2, 1)$

$z^{[2]}$ . shape $= (2, m)$

Backward Propagation

# Backward Propagation



$$\text{forward} = \text{input } A^{(L-1)} \quad A^{(L)}$$
$$\text{coche}(Z^{(L)})$$

$$\text{backward} = \text{input } dA^{(L)}$$
$$\text{coche}(Z^{(L)}) \quad \text{output}$$
$$dA^{(L-1)}$$
$$dw^{(L)}$$
$$db^{(L)}$$

$$dA^{(L)} = \frac{\partial J}{\partial A^{(L)}}$$

$$A^{(0)} \rightarrow \boxed{w^{(1)} b^{(1)}} \rightarrow A^{(1)} \rightarrow \boxed{w^{(2)} b^{(2)}} \rightarrow A^{(2)} \rightarrow \boxed{w^{(3)} b^{(3)}} \rightarrow A^{(3)} \rightarrow \hat{y}$$

L1    L2    L3

coche $(Z^{(1)})$    coche $(Z^{(2)})$    coche $Z^{(3)}$

$$\frac{\partial (L(A^{(3)}; y))}{\partial A^{(L)}} = dA^{(3)}$$

$$dA^{(0)} \leftarrow \boxed{w^{(1)} b^{(1)}} \xleftarrow{dA^{(1)}} \boxed{w^{(2)} b^{(2)}} \xleftarrow{dA^{(2)}} \boxed{w^{(3)} b^{(3)}} \xleftarrow{dA^{(3)}}$$

$$dw^{(1)} db^{(1)} \quad dw^{(2)} db^{(2)} \quad dw^{(3)} db^{(3)}$$

# Backward Propagation



$$H^{[5]} = g\left(z^{[5]}\right)$$

$$z^{[5]} = W^{[5]} A^{[4]} + b^{[5]}$$

$$L\left(A^{[5]}\right) = \boxed{y \log A^{[5]} + (1-y) \log A^{[5]}}$$

binary cross-entropy

$$A^{[5]} = \hat{y}$$

$$A^{[5]} = \sigma\left(z^{[5]}\right)$$

$$L(\hat{y}) = 1 \cdot \log \frac{a6}{} = \log a6$$

$$\frac{\partial L}{\partial W^{[5]}} = \boxed{\frac{\partial L}{\partial A^{[5]}}} \cdot \left(\boxed{\frac{\partial A^{[5]}}{\partial z^{[5]}} \cdot \frac{\partial z^{[5]}}{\partial W^{[5]}}}\right)$$

$$\frac{1}{06}$$

$$\frac{\partial L}{\partial A^{[5]}} = \frac{\partial \left(\log A^{[5]}\right)}{\partial W^{[5]}} = \frac{1}{A^{[5]}} =$$

$$\frac{\partial A^{[5]}}{\partial z^{[5]}} = \frac{\partial \left(\sigma\left(z^{[5]}\right)\right)}{\partial z^{[5]}} = \sigma\left(z^{[5]}\right)\left(1 - \sigma\left(z^{[5]}\right)\right)$$

$$\frac{\partial z^{[5]}}{\partial W^{[5]}} = \frac{\partial \left(W^{[5]} A^{[4]} + b^{[5]}\right)}{\partial W^{[5]}} = A^{[4]}$$

$$W^{[5]} b^{[5]} + b^{[5]} = A^{[4]}$$

$$z^{[5]} = W b + b$$

$$\int f(\tau) \cdot h(t - \tau) \, d\tau$$

$$\frac{\partial L}{\partial W^{[4]}} = \frac{\partial L}{\partial A^{[5]}} \cdot \frac{\partial A^{[5]}}{\partial z^{[5]}} \cdot \frac{\partial z^{[5]}}{\partial A^{[4]}} \cdot \frac{\partial A^{[4]}}{\partial z^{[4]}} \cdot \frac{\partial z^{[4]}}{\partial W^{[4]}}$$

$$\frac{1}{A^{[5]}} \cdot \sigma'\left(z^{[5]}\right) \cdot W^{[5]} \cdot \sigma'\left(z^{[4]}\right) \cdot A^{[4]}$$

$$\frac{1}{06} \cdot \frac{\partial L}{\partial W^{[5]}} = \frac{1}{A^{[5]}} \cdot \sigma'\left(z^{[5]}\right) \cdot A^{[4]}$$

$$W^{[5]} = W^{[5]} - \alpha \, dW^{[5]}$$

$$\frac{\partial A^{[5]}}{\partial z^{[5]}} = \frac{\sigma\left(z^{[5]}\right)}{z^{[5]}} = \frac{\sigma\left(W^{[5]} A^{[4]} + b^{[5]}\right)}{z^{[5]}} = \boxed{\sigma\left(z^{[5]}\right) \cdot \left(1 - \sigma\left(z^{[5]}\right)\right)} = \frac{\partial A^{[5]}}{\partial z^{[5]}}$$

# Hyperparameters



Don't worry about it if you don't understand

- Andrew Ng

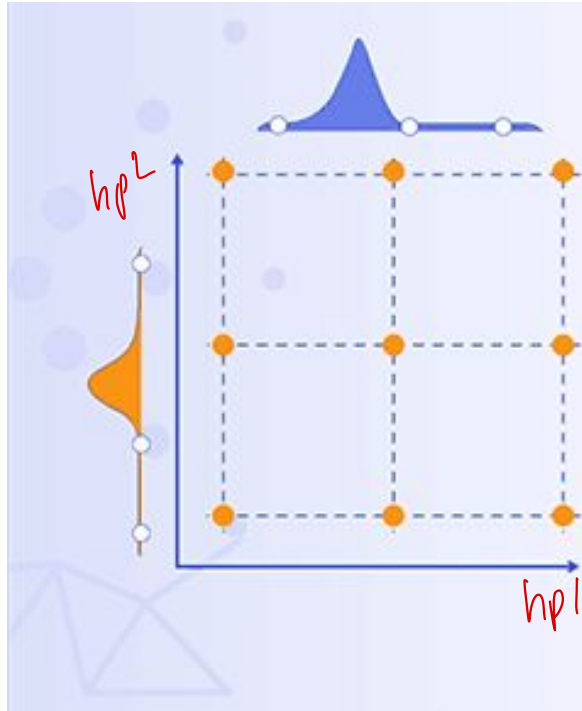# Hyperparameters

Hyperparameters effect parameters

Hyperparameter examples:

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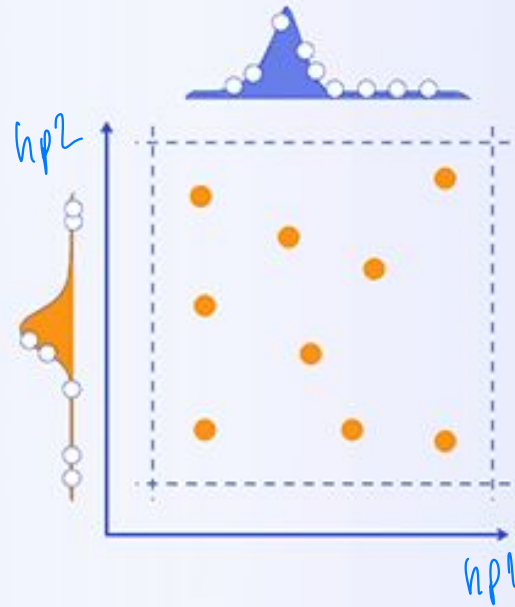- #Units
- #Iterations
- #Layers
- Batch size

We can select hyperparameters using several methods

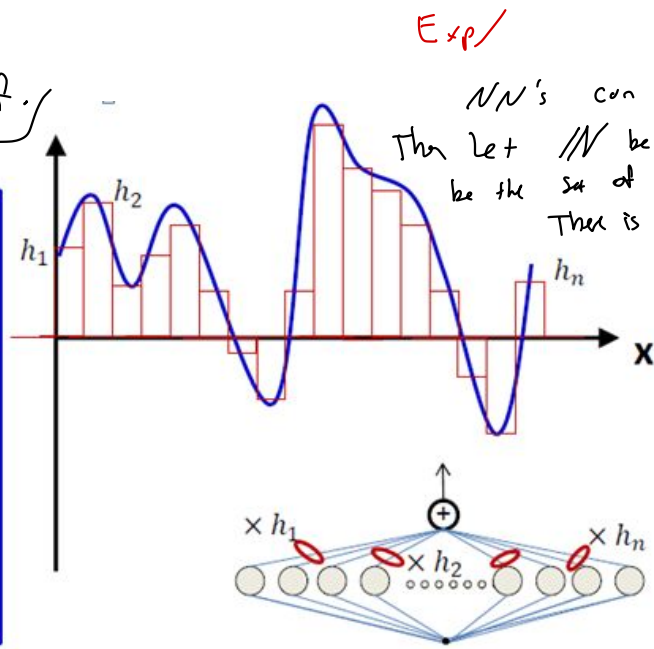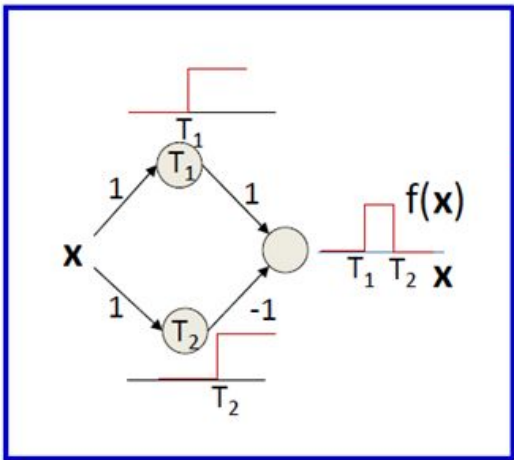# Hyperparameter Tuning



Grid Search

Random Search

# Universal Approximation Theorem

The Universal Approximation Theorem tells us that Neural Networks has a kind of universality no matter what f(x) is, there is a network that can approximately approach the result.

Exp/

NN's can approximate any f.

NN's can be represented as functions.
Then let N be the set of NN's and F be the set of functions on cont. space Then There is a function such that

$$G : N \longrightarrow F$$

G is onto or surjective

inzva: *brings the AI fellows together*

inzva:



Now this is an Avengers level threat

@debo