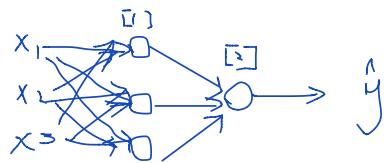


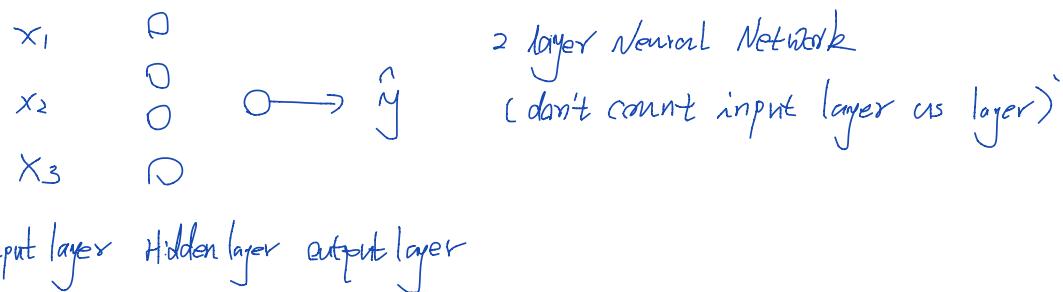
• Neural Network overview

What is a neural Network

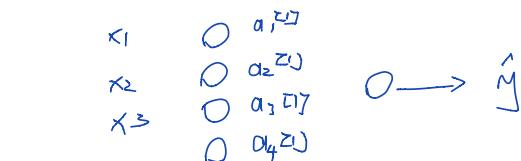
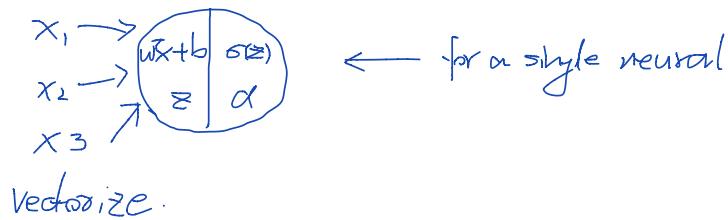


$$x \begin{matrix} \xrightarrow{w^{[1]}} \\ \vdots \\ \xrightarrow{w^{[L]}} \end{matrix} z^{[l]} = w^{[l]}x + b^{[l]} \rightarrow a^{[l]} = \sigma(z^{[l]}) \rightarrow z^{[l+1]} = \dots \rightarrow a^{[L]} = \dots$$

• Neural Network Representation



• Computing a Neural Network's output



$$\underline{z}^{[2]} = \begin{bmatrix} -w_1^{[2]T} \\ w_2^{[2]T} \\ -w_3^{[2]T} \\ w_4^{[2]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_1^{[2]} \\ b_2^{[2]} \\ b_3^{[2]} \\ b_4^{[2]} \end{bmatrix} = \begin{bmatrix} w_1^{[2]T}x + b_1^{[2]} \\ w_2^{[2]T}x + b_2^{[2]} \\ -w_3^{[2]T}x + b_3^{[2]} \\ w_4^{[2]T}x + b_4^{[2]} \end{bmatrix} = \begin{bmatrix} z_1^{[2]} \\ z_2^{[2]} \\ z_3^{[2]} \\ z_4^{[2]} \end{bmatrix}$$

$$a^{[2]} = \sigma(\underline{z}^{[2]})$$

• Vectorize across multiple examples
in training examples.

$$\begin{aligned} x &\rightarrow a^{[2]} = \hat{y} \\ x^{(1)} &\longrightarrow a^{[2](1)} = \hat{y}^{(1)} \\ x^{(2)} &\longrightarrow a^{2} = \hat{y}^{(2)} \\ &\vdots \quad \vdots \quad \vdots \\ x^{(m)} &\longrightarrow a^{[2](m)} = \hat{y}^{(m)} \end{aligned}$$

for $i = 1$ to m :

$$\underline{z}^{[2](i)} = w^{[2]} x^{(i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(\underline{z}^{[2](i)})$$

$$\underline{z}^{[2](i)} = w^{[2]} a^{[2](i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(\underline{z}^{[2](i)})$$

How to vectorize above calculation and get rid of for loop?

$$X = \left[\begin{array}{c|c|c|c} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \end{array} \right]$$

$$z^{[1]} = w^{[1]} X + b^{[1]}$$

$$A^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(z^{[2]})$$

row: different features / hidden neurons

column: different training examples

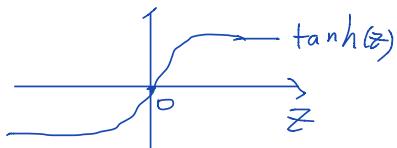
Activation functions

Sigmoid function σ

more generally, the activation function can be g

→ better choice of activation function: $g(z) = \tanh(z)$

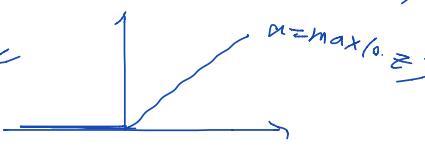
$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



Advantage: compare to sigmoid func, the $\tanh(z)$ gives out output centered around average of 0

[*] Most of times in neural network $\tanh(z)$ is used instead of sigmoid func, but one exception is that when working on a classification problem, the output has to be $y \in \{0, 1\}$, in this situation the output layer uses sigmoid func as the activation function.

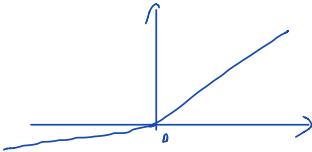
ReLU function: Rectified Linear Unit
 $a = \max(0, z)$
 sigmoid func or tanh func
 slows down gradient descent



General Rule:

- Different layer can have different activation func
- Don't use sigmoid func as activation func unless for the output layer of binary classification
- Most of times people uses ReLU func

Leak ReLU



why do Neural Network need non-linear activation function?

If assume $g(z) = z$ "linear activation function"

then

$$\begin{aligned} a^{[2]} &= w^{[2]}(w^{[1]}x + b^{[1]}) + b^{[2]} \\ &= (w^{[2]}w^{[1]})x + w^{[2]}b^{[1]} + b^{[2]} \\ &= w'x + b' \end{aligned}$$

It's just doing linear regression, hidden unit does nothing!

The only possible case of linear activation function is: all hidden layers use non-linear activation func (e.g. ReLU, tanh...) but the output layer can use linear function (e.g. for linear regression as the last step)

- Derivatives of activation functions

- Sigmoid activation function

$$g(z) = \frac{1}{1 + e^{-z}} \quad \frac{dg(z)}{dz} = g(z)(1 - g(z))$$

- tanh activation function

$$g(z) = \tanh(z) \quad \frac{dg(z)}{dz} = 1 - (\tanh(z))^2$$

- ReLU activation function

$$g(z) = \max(0, z) \quad g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

- Leaky ReLU Activation function

$$g(z) = \max(0.01z, z) \quad g'(z) = \begin{cases} 0.01, & \text{if } z < 0 \\ 1, & \text{if } z \geq 0 \end{cases}$$

| o Gradient Descent for Neural Network: back propagation

for single hidden Units

Parameters:

o Num of units in each layer: $n_x = n^{[0]}$, $n^{[1]}$, $n^{[2]}$. ($n^{[0]} = 1$ for single hidden units)

o $w^{[1]}: n^{[0]} \times n^{[1]}$

$b^{[1]}: n^{[1]} \times 1$

$w^{[2]}: n^{[1]} \times n^{[2]}$

$b^{[2]}: n^{[2]} \times 1$

o Cost function:

$$J(w^{[0]}, b^{[0]}, w^{[1]}, b^{[1]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$$

o Gradient descent

randomly initialize w_s , b_s . (non-zero)

Repeat {

Compute predictions ($\hat{y}^{(i)}$, $i=1, \dots, m$)

$$\delta w^{[1]} = \frac{\partial J}{\partial w^{[1]}}, \quad \delta b^{[1]} = \frac{\partial J}{\partial b^{[1]}}$$

$$w^{[1]} = w^{[1]} - \alpha \delta w^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha \delta b^{[1]}$$

}

o Formula for Computing derivatives

Forward Propagation

$$z^{[1]} = w^{[0]} X + b^{[0]}$$

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

$$z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$a^{[2]} = g(z^{[2]})$$

Back propagation

$$dz^{[2]} = A^{[2]} - Y \quad (\text{when } g = \text{sigmoid}(z))$$

$$dW^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]\top} \quad (\leftarrow \text{broadcasting})$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims=True})$$

$$dz^{[1]} = W^{[2]\top} dz^{[2]} * g'(z^{[2]})$$

↑ element-wise product

$$dW^{[1]} = \frac{1}{m} dz^{[1]} X^T \quad (\text{broadcasting})$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims=True})$$

Random Initialization

• what happens if initialize all weights to 0?

Different hidden unit are identical.

• Have to initialize parameters randomly, to a SMALL value..

$$W^{[1]} = \text{np.random.randn}(2, 2) \times 0.01$$

why small? $z = W^T X + b$ $a = g(z)$

if W to large, gradient will be small, especially for sigmoid functions.