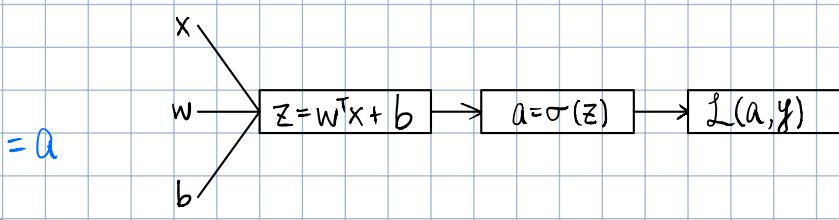
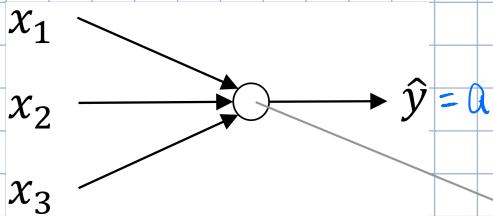
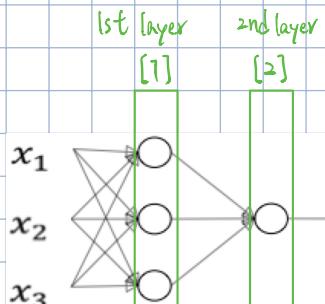


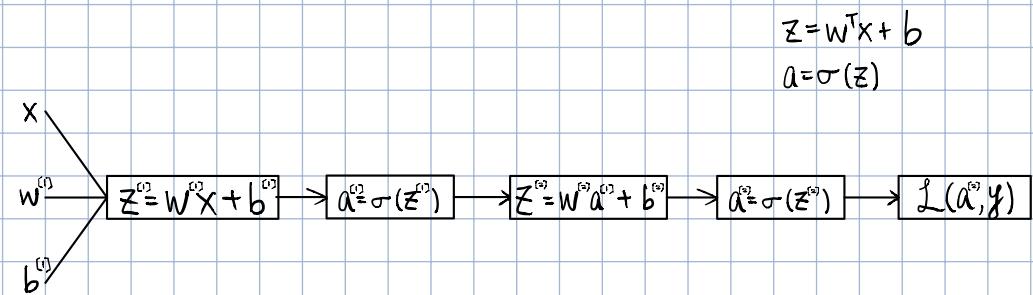
Neural Networks Overview:



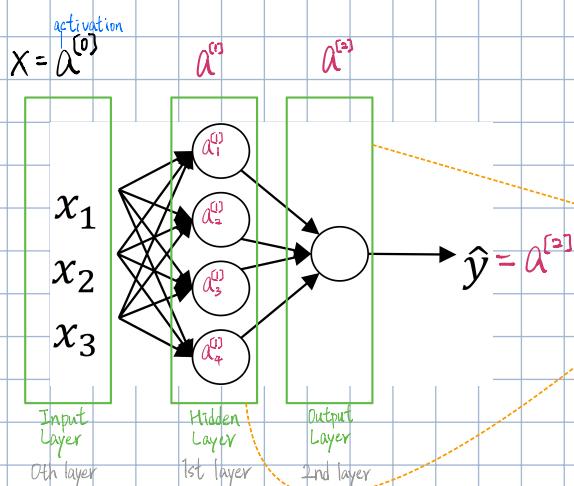
This node does two things: $z = w^T x + b$ and $a = \sigma(z)$



Actually every node in this graph can also do these two things.



Neural Network Representation: One hidden Layer



Conventionally, it's called 2 layer neural network. The input layer is not counted.

The associated parameters for the hidden layer and output layer:
 $W^{(2)}: 1 \times 4, b^{(2)}: 1 \times 1$
 $W^{(1)}: 4 \times 3, b^{(1)}: 4 \times 1$

Standard notations for Deep Learning

This document has the purpose of discussing a new standard for deep learning mathematical notations.

1 Neural Networks Notations.

General comments:

- superscript (i) will denote the i^{th} training example while superscript [l] will denote the l^{th} layer

Sizes:

- m : number of examples in the dataset
 - n_x : input size
 - n_y : output size (or number of classes)
 - $n_h^{[l]}$: number of hidden units of the l^{th} layer
- In a for loop, it is possible to denote $n_x = n_h^{[0]}$ and $n_y = n_h^{[\text{number of layers} + 1]}$.
- L : number of layers in the network.

Objects:

- $X \in \mathbb{R}^{n_x \times m}$ is the input matrix
- $x^{(i)} \in \mathbb{R}^{n_x}$ is the i^{th} example represented as a column vector

· $Y \in \mathbb{R}^{n_y \times m}$ is the label matrix

· $y^{(i)} \in \mathbb{R}^{n_y}$ is the output label for the i^{th} example

· $W^{[l]} \in \mathbb{R}^{\text{number of units in next layer} \times \text{number of units in the previous layer}}$ is the weight matrix,superscript [l] indicates the layer

· $b^{[l]} \in \mathbb{R}^{\text{number of units in next layer}}$ is the bias vector in the l^{th} layer

· $\hat{y} \in \mathbb{R}^{n_y}$ is the predicted output vector. It can also be denoted $a^{[L]}$ where L is the number of layers in the network.

Common forward propagation equation examples:

$a = g^{[l]}(W_x x^{(i)} + b_1) = g^{[l]}(z_1)$ where $g^{[l]}$ denotes the l^{th} layer activation function

$\hat{y}^{(i)} = softmax(W_h h + b_2)$

- General Activation Formula: $a_j^{[l]} = g^{[l]}(\sum_k w_{jk}^{[l]} a_k^{[l-1]} + b_j^{[l]}) = g^{[l]}(z_j^{[l]})$
- $J(x, W, b, y)$ or $J(\hat{y}, y)$ denote the cost function.

Examples of cost function:

· $J_{CE}(\hat{y}, y) = -\sum_{i=0}^m y^{(i)} \log \hat{y}^{(i)}$

· $J_1(\hat{y}, y) = \sum_{i=0}^m |y^{(i)} - \hat{y}^{(i)}|$

2 Deep Learning representations

For representations:

- nodes represent inputs, activations or outputs
- edges represent weights or biases

Here are several examples of Standard deep learning representations

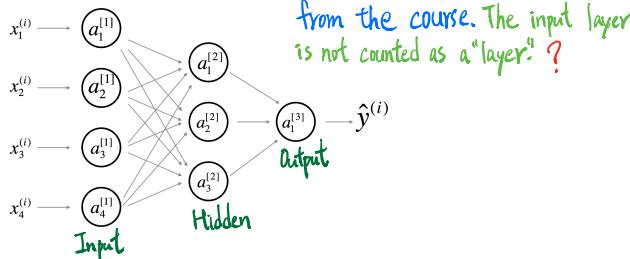


Figure 1: Comprehensive Network: representation commonly used for Neural Networks. For better aesthetic, we omitted the details on the parameters ($w_{ij}^{[l]}$ and $b_i^{[l]}$ etc...) that should appear on the edges

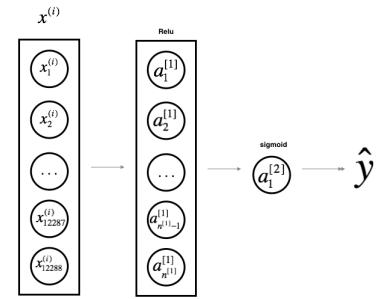
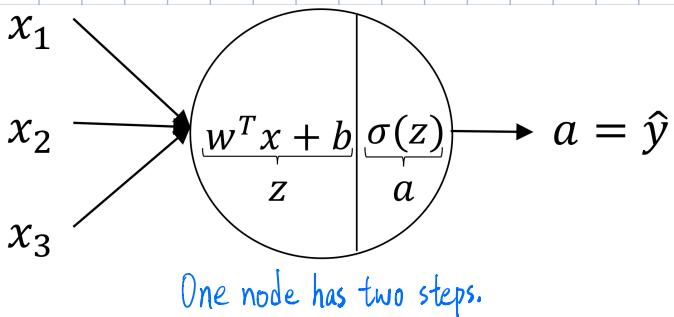


Figure 2: Simplified Network: a simpler representation of a two layer neural network, both are equivalent.

Computing a Neural Network's Output:



$$z = w^T x + b$$

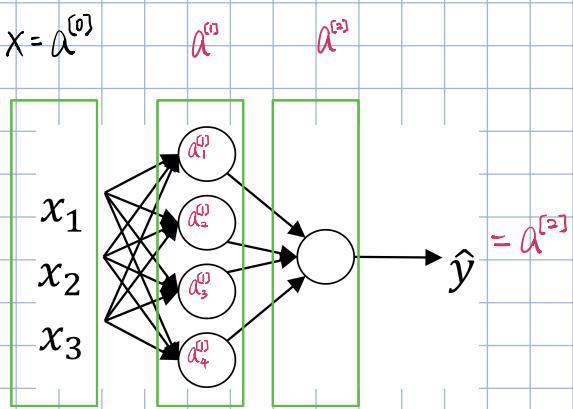
$$a = \sigma(z)$$

The general form is:

$$a_j^{(l)} = g(\vec{w}_j^{(l)} \cdot \vec{a}^{(l-1)} + b_j^{(l)})$$

$a_j^{(l)}$: the activation value of layer l , the j -th unit.

$\vec{w}_j^{(l)}, b_j^{(l)}$: Parameters of layer l , unit j .



$$z_1^{(1)} = W_1^{(1)T} X + b_1^{(1)}, a_1^{(1)} = \sigma(z_1^{(1)})$$

⋮

⋮

$$z_4^{(1)} = W_4^{(1)T} X + b_4^{(1)}, a_4^{(1)} = \sigma(z_4^{(1)})$$

$$\text{Thus, } z = W^T X + b, \quad W^T: 4 \times 3$$

$$X: 3 \times 1$$

$$b: 4 \times 1$$

$$a^{(1)}: 4 \times 1$$

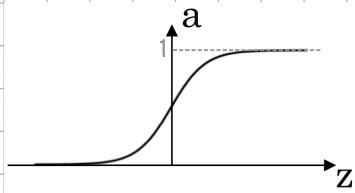
Vectorize across Multiple Samples. **Do Not use for-loop.**
Use vectorization.

$$X^{(1)} \longrightarrow \hat{y}^{(1)} = a^{(2)(1)}$$

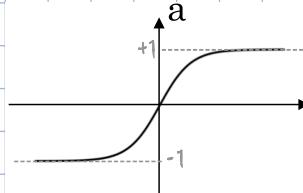
⋮

$$X^{(m)} \longrightarrow \hat{y}^{(m)} = a^{(2)(m)}$$

Activation functions:



$$\text{Sigmoid func: } a = \frac{1}{1+e^{-z}}$$

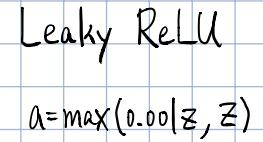
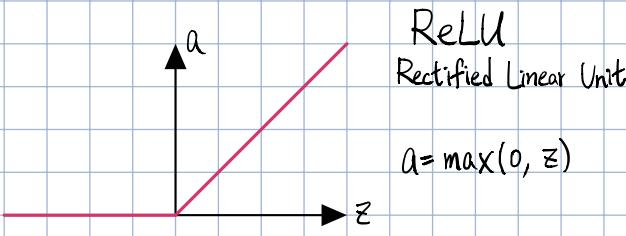


Hyperbolic tangent:

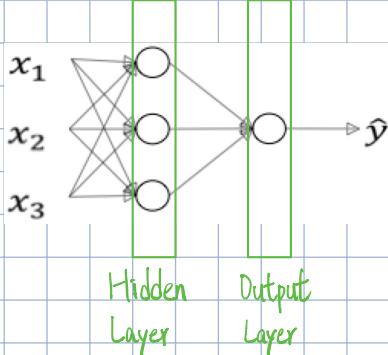
$$a = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

The main drawback of the sigmoid function and hyperbolic tangent is when z is very large or small, then the gradient/slope of sig() and tanh() becomes being close to 0, which slows down gradient descent. Thus, we need ReLU.

Actually, in practice, sigmoid function is almost never used. Except for the output layer of binary classification. Explain Why.



Why Non-Linear Activation functions:



Given the network and X , we have

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

$a^{[1]} = g^{[1]}(z^{[1]})$, $g^{[1]}$ is the activation func of the hidden layer.

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

$a^{[2]} = g^{[2]}(z^{[2]})$, $g^{[2]}$ is the activation func of the output layer.

Now, we select a linear function, an identity function, as the activation function. Thus, $g(z) = z$, $a^{[1]} = z^{[1]}$, $a^{[2]} = z^{[2]}$

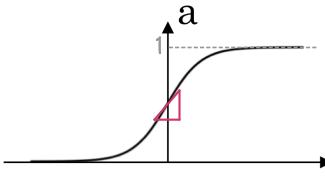
Then we can calculate \hat{y} .

$$\begin{aligned} \hat{y} &= a^{[2]} = z^{[2]} = W^{[2]}(W^{[1]}X + b^{[1]}) + b^{[2]} \\ &= (W^{[2]}W^{[1]})X + W^{[2]}b^{[1]} + b^{[2]} \\ &= W'X + b' \end{aligned}$$

This shows that the neural network is generating a linear function of an input, which is equal to no hidden layers at all.

Now, linear activation function may be used in one case, which is the output layer of a regression problem. E.g., we want to predict house's price $y \in \mathbb{R}$, $0 < y < 10^6$. Or the disparity map of a stereo pair image, $-10 < y < 10$. But all hidden layers should still use non-linear activation function.

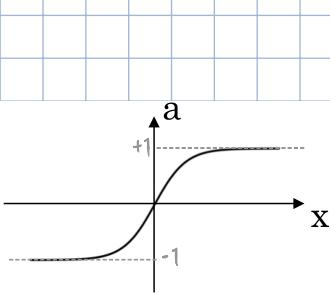
Derivatives of Activation functions:



Sigmoid func:

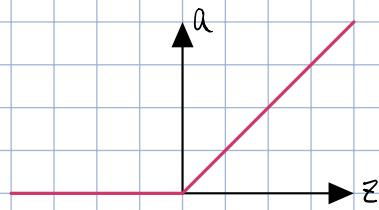
$$g(z) = \frac{1}{1 + e^{-z}}, \quad \frac{d}{dz} g(z) = \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}}\right) = g(z)(1 - g(z))$$

The other notation: $g'(z) = \frac{d}{dz} g(z)$, $g(z) = a$, Thus, $g'(z) = a(1 - a)$



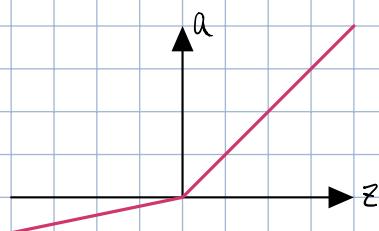
Hyperbolic Tangent: $g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

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



ReLU: $g(z) = \max(0, z)$

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



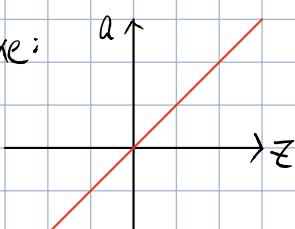
Leaky ReLU: $g(z) = \max(0.01|z|, z)$

$$a = \max(0.01|z|, z)$$

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

There is another activation function called "Linear activation function" which look

like:



$$a = g(z) = z = \vec{w} \cdot \vec{x} + b$$

Q: How do I choose which activation function to use? A: It depends.

For output layer:

It depends on the type of output layer.

Case 1: Binary Classification. Choose Sigmoid. Eg: cat v.s. non-cat

Case 2: Regression. $y \geq 0$ or $y < 0$. Choose Linear. Eg: disparity

Case 3: Regression. $y \geq 0$. Choose ReLU. Eg: house price

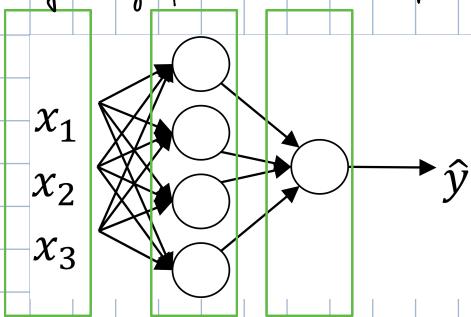
For hidden layer:

most common choice ReLU or Leaky ReLU

Do not use linear activation function in hidden layers. If so, it's equivalent to linear regression or logistic regression.

Gradient Descent for Neural Networks:

Using the graph below as an example:



n_x : Number of nodes in each layer. $n_x: n^{(0)}=3, n^{(1)}=4, n^{(2)}=1$.

Parameters: $W^{(0)}: n^{(0)} \times n^{(1)}, b^{(0)}: n^{(1)} \times 1, W^{(1)}: n^{(1)} \times n^{(2)}, b^{(1)}: n^{(2)} \times 1$

Cost function for binary classification: $J(W^{(0)}, b^{(0)}, W^{(1)}, b^{(1)})$

$$= \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$$

Gradient Descent:

Repeat }

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

$$dW^{(0)} = \frac{dJ}{dW^{(0)}}, db^{(0)} = \frac{dJ}{db^{(0)}}, dW^{(1)} = \frac{dJ}{dW^{(1)}}, db^{(1)} = \frac{dJ}{db^{(1)}}$$

$$W^{(0)} = W^{(0)} - \alpha dW^{(0)}, b^{(0)} = b^{(0)} - \alpha db^{(0)}, W^{(1)} = W^{(1)} - \alpha dW^{(1)}, b^{(1)} = b^{(1)} - \alpha db^{(1)}$$

}

Now, the question is how to find out dW and db ? Follow the formula below: Note that the capital letter denotes the vectorised components: E.g. $Y = [y_0, y_1, \dots, y_m]$

Forward Propagation:

$$Z^{(0)} = W^{(0)} X + b^{(0)}$$

$$A^{(0)} = g^{(0)}(Z^{(0)})$$

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

$$A^{(1)} = g^{(1)}(Z^{(1)}) = \sigma(Z^{(1)})$$

sigmoid function

Backward Propagation:

$$dZ^{(2)} = A^{(2)} - Y$$

$$dW^{(2)} = \frac{1}{m} dZ^{(2)} A^{(1)T}$$

$$db^{(2)} = \frac{1}{m} \text{sum}(dZ^{(2)}, \text{axis}=1)$$

$$dZ^{(1)} = (W^{(1)T} dZ^{(2)}) * (g^{(1)' Z^{(1)}})$$

$n^{(1)} \times m$

$$dW^{(1)} = \frac{1}{m} dZ^{(1)} X^T$$

$$db^{(1)} = \frac{1}{m} \text{sum}(dZ^{(1)}, \text{axis}=1)$$

$n^{(0)} \times 1$

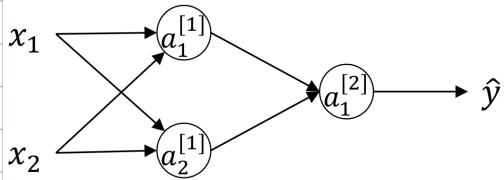
Element-wise product

Please explain/derive the formula above.

Check C1W3L10

Random Initialization of Weights

What happens if I initialize weights to zero?



$$n^{(0)} = 2 \quad n^{(1)} = 2 \quad n^{(2)} = 1$$

$$w^{(0)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, b^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, w^{(1)} = [0 \ 0]$$

$$\rightarrow a_1^{(0)} = a_2^{(0)}$$

$$dZ_1^{(0)} = dZ_2^{(0)}, \quad dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix},$$

When updating w , $w^{(1)} = w^{(0)} - \alpha dw$,

w will end up with the same content for every row.

Here is the right way to initialize weights.

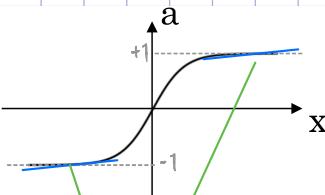
$w^{(0)} = \text{random}(2, 2) \times 0.01$ Why initialize weights a small value?

$b^{(0)} = \text{random}(2, 1)$

\vdots

Recall: $Z^{(0)} = w^{(0)}X + b^{(0)}$

$$a^{(0)} = g^{(0)}(Z^{(0)})$$



The gradient is small or close to zero.

If weights $w^{(0)}$ are large value,

$Z^{(0)}$ will be either very large

value or a very small value.

Either way, it will result to a

small gradient which will make

the gradient descent (learning)

really slow.

The best strategy to initialize weights

is probably Xavir Initializer.