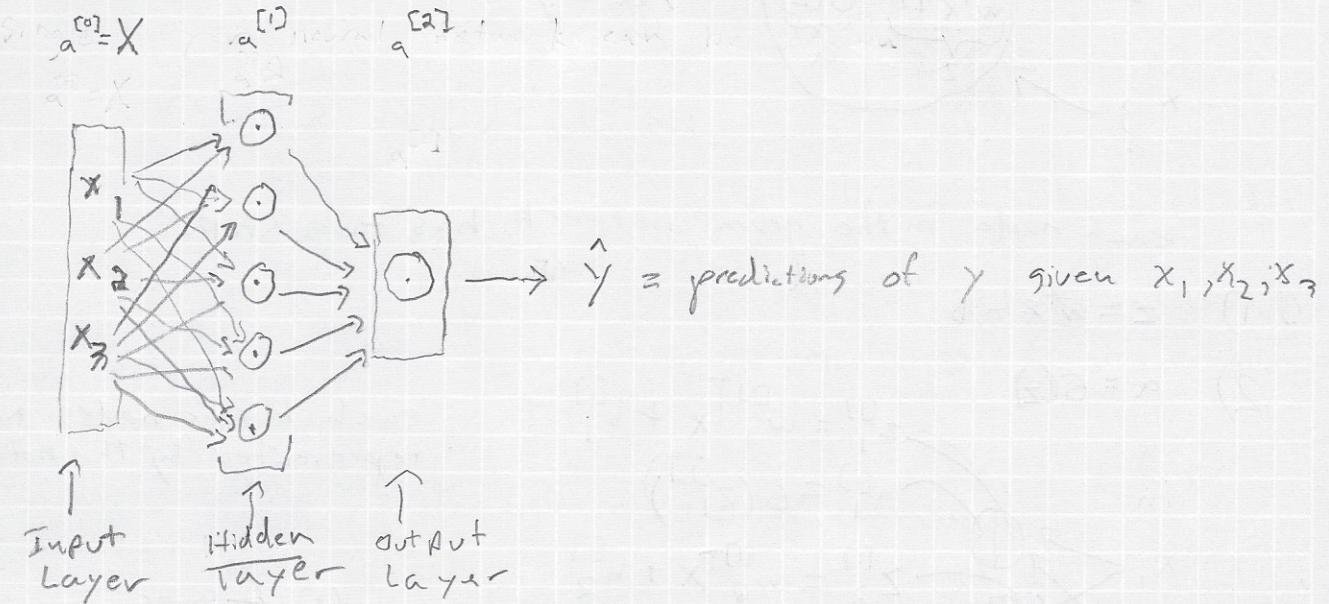


Week 3 - Shallow Neural NetworksNeural Network Representations

The hidden layer means you don't see them in the training set.

In this course what we were calling X now we represent X as $a^{[0]}$.

$$X = a^{[0]}$$

and " $y = a$ "

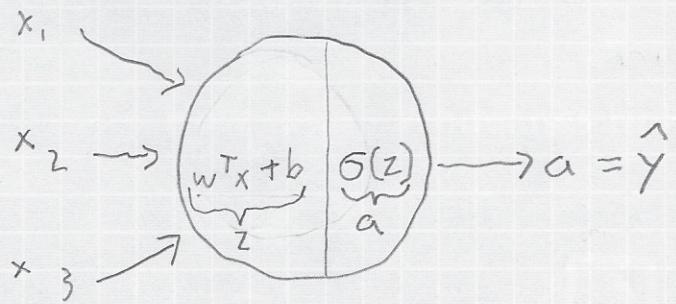
This is called a [2 layer] neural network because we don't count the input layer (by convention).

The hidden layer will also have associated with $w^{[1]}, b^{[1]}$.

In this example $w^{[1]}$ will be a $(4, 3)$ matrix 3 input features
 $b^{[1]}$ will be a $(4, 1)$ vector 4 nodes/hidden units

The outputs $w^{[2]}$ will be a $(1, 4)$ 4 input hidden units,
 and $b^{[2]}$ " " " $(1, 1)$ 1 output

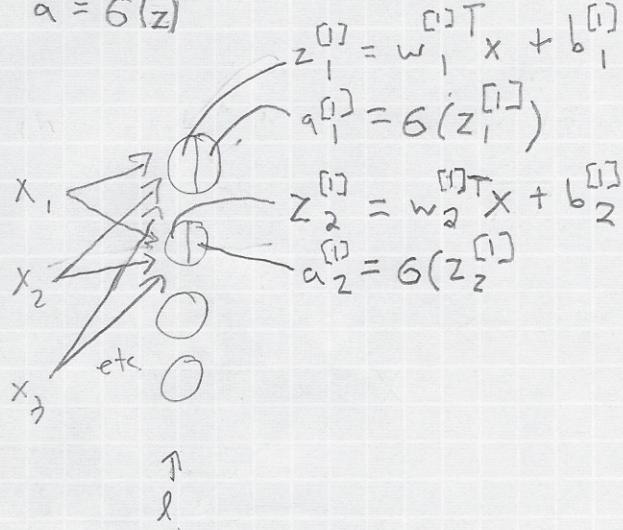
Week 3 - Computing a Neural Network's Output



each node in the neural network has two steps

$$1) z = w^T x + b$$

$$2) a = \sigma(z)$$



each layer's nodes are represented by this notation

$[l]$ ← layer
 $a_i^{[l]}$ ← node in layer

In this example, with 4 nodes we want to compute all z 's, and a 's at the same time, so we transpose our column vectors into rows, and stack them as follows:

$$\begin{matrix} & \begin{bmatrix} w_1^{[1] T} \\ w_2^{[1] T} \\ w_3^{[1] T} \\ w_4^{[1] T} \end{bmatrix} \\ \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} & R \end{matrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

matrix multiplication

We have 4 logistic regression units, and each unit has a corresponding vector w by stacking them you get a 4 by 3 Matrix

$$W^{[1]}$$

$$Z^{[1]}$$

We stack nodes vertically

So, given input x :

$$\begin{aligned} \text{first layer} & \quad z^{[1]} = W^{[1]}x + b \\ & \quad (4,1) \quad (4,3) \quad (3,1) \quad (4,1) \\ & \quad a^{[1]} = \sigma(z^{[1]}) \\ & \quad (4,1) \quad (4,1) \end{aligned}$$

$$\begin{aligned} \text{Second layer} & \quad z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ & \quad (1,1) \quad (1,4) \quad (4,1) \quad (1,1) \\ & \quad a^{[2]} = \sigma(z^{[2]}) \\ & \quad (1,1) \quad (1,1) \end{aligned}$$

to compute the output of a 2 layer nn for a single training example.

Vectorizing Across Multiple Examples

$$\begin{aligned} x & \rightarrow a^{[2]} = y \quad \text{for a single training example} \\ x_1 & \rightarrow a^{[2](1)} = y^{(1)} \\ x_2 & \rightarrow a^{2} = y^{(2)} \\ \vdots & \rightarrow a^{[2](m)} = y^{(m)} \end{aligned}$$

$a^{[2](i)}$ ← example i
layer 2

for $i=1$ to m :

$$\begin{aligned} z^{[1](i)} &= w^{[1]}x^{(i)} + b^{[1]} \\ z^{[1](i)} &= \sigma(z^{[1](i)}) \\ z^{[2](i)} &= w^{[2]}a^{[1](i)} + b^{[2]} \\ z^{[2](i)} &= \sigma(z^{[2](i)}) \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{train on all examples } m$$

but this is an inefficient for loops, let's vectorize it
we stack the x, z , and a vectors into matrices

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \quad (n_x, m)$$

$$Z^{[1]} = \begin{bmatrix} | & | & | \\ z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \\ | & | & | \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} | & | & | \\ a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \\ | & | & | \end{bmatrix}$$

Vertical index corresponds to different nodes in + nn

Horizontally we index across training examples, when you sweep left to right you are scanning through the training set

$$A^{[1]} = \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \end{bmatrix}$$

The value at the topmost corner of the matrix corresponds to the activation of the first hidden unit on the first training example

one value down corresponds to the second hidden unit on the first training example

holds true for Z and X also

scanning down index's into the hidden unit's number
 \rightarrow training examples (M)

Activation Functions

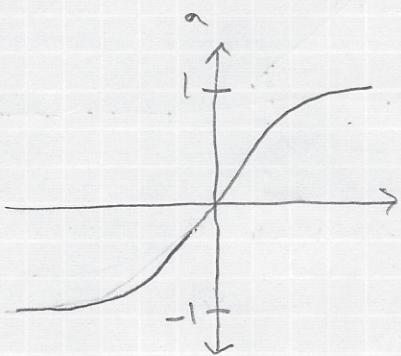
You choose which functions to use for each layer and the output.

The hyperbolic tangent function almost always works better than the sigmoid function.

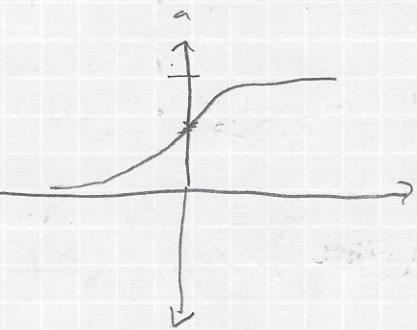
$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

\tanh centers your data around zero which makes learning for the next layer easier

σ
sigmoid



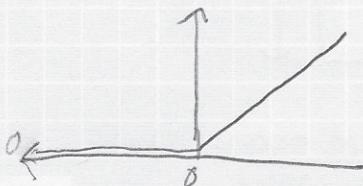
vs



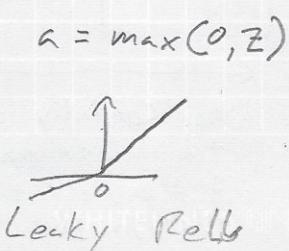
The one exemption would be for the output in Binary Classification.

we use $g(z)$ to signify the activation

If z is very large or very small then the gradient/derivative/slope will be very small so it goes very slowly.
 In these cases we can use $\text{ReLU} = \text{rectified Linear Unit}$



ReLU



Leaky ReLU

very common now in hidden units

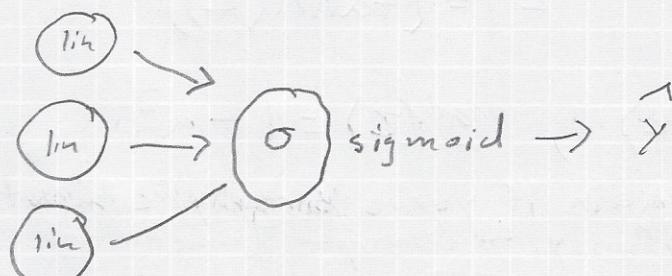
Why do you need non-linear activation functions?

If you use a linear function for $g(z) = z$ ← the "linear activation function"
(or no function)
(could call this "identity activation function")

If you do this it turns out that this model is just computing \hat{y} as a linear function of your input features x .

If you were to use linear activation functions, or alternatively have no activation function then the NN is outputting a linear function of the input. No matter how many layers, all its doing is computing a linear activation function, so you might as well not have any hidden layers.

Even if you have a model:



it is no more expressive than standard logistic regression without any hidden layers. So linear regression hidden unit layer is useless. The composition of 2 linear functions is itself a linear function.

In the single case of doing ML on a regression problem where $y \in \mathbb{R}$
(y is a real number)
e.g. Housing price prediction where $y \in \$0 \dots \$1,000,000$

If y takes on real values then it might be okay to have a linear output activation function. But then the hidden units should not be linear functions (could be ReLU, or tanh, or leaky ReLU, etc.).

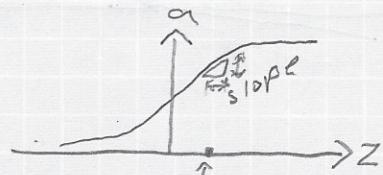
But even in this case you could use a ReLU because housing prices are never negative.

Derivatives of Activation Functions

You need to know this for understanding gradient descent.

Sigmoid activation function derivative

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



$\frac{d}{dz} g(z) = \text{slope of } g(x) \text{ at } z$

look up or
research
calculus
proof

$$\left\{ \begin{aligned} &= \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right) \\ &= g(z)(1-g(z)) \end{aligned} \right.$$

also call $g'(z) = a(1-a)$

Tanh activation function derivative

$$g(z) = \tanh(z)$$

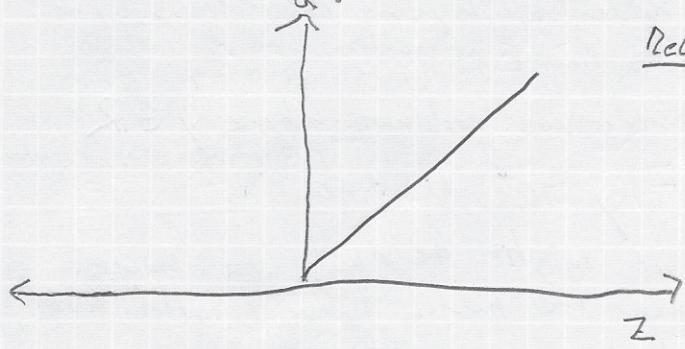
$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\begin{aligned} g'(z) &= \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z \\ &= 1 - (\tanh(z))^2 \end{aligned}$$

$$\text{if } a = g(z), \quad g'(z) = 1 - a^2$$

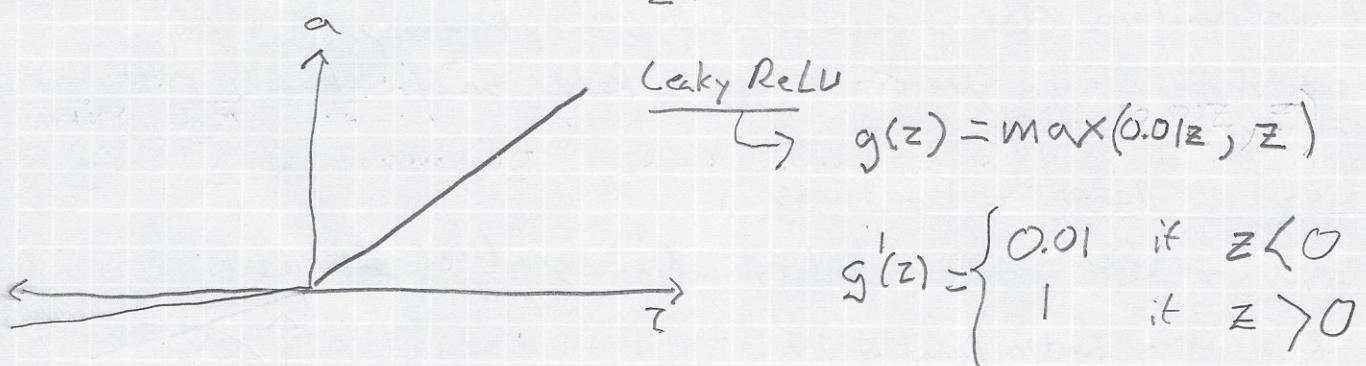
so once again if you've already computed the value of a this is very easy to compute (same thing we saw with the sigmoid function)

ReLU and Leaky ReLU derivatives



$$\underline{\text{ReLU}} \rightarrow g(z) = \max(0, z)$$

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



$$\underline{\text{Leaky ReLU}} \rightarrow$$

$$g(z) = \max(0.01z, z)$$

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

Gradient descent for neural networks

We're just going to give you the equations you need to implement in order to get gradient descent, our back propagation working. In the next section we'll give some intuitions for why these particular equations are the accurate equations, the correct equations, for computing the gradients you need for your NN.

NN with a single hidden layer.

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$

if you have $n_x = n^{[0]}$, $n^{[1]}$, $n^{[2]} = 1$

input features hidden units output units

then

$W^{[1]}$ is matrix $(n^{[1]}, n^{[0]})$

$b^{[1]}$ $(n^{[1]}, 1)$

← an $n^{[1]}$ -dimensional vector

$W^{[2]}$ $(n^{[2]}, n^{[1]})$

$b^{[2]}$ $(n^{[2]}, 1)$

Cost function (assuming we're doing binary classification)

$$J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}, y) \quad \underbrace{a^{[2]}}$$

Gradient descent

This optimizes the parameters with respect to the cost function. In gradient descent it's important to randomize the parameter initialization, rather than to all 0s as we did earlier in this course.

Repeat ↗

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

Compute derivatives $dW^{[1]} = \frac{dJ}{dw^{[1]}}, db^{[1]} = \frac{dJ}{db^{[1]}}, \dots$

g.d. update

$$W^{[1]} = W^{[1]} - \alpha dW^{[1]}$$

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

$$\begin{matrix} W^{[2]} \\ b^{[2]} \end{matrix} \dots \quad \begin{matrix} \uparrow \\ \text{learning rate} \end{matrix}$$

you do this until the parameters look like they are converging. The key you need is how to calculate the partial derivatives.

Forward Propagation

"left to right"

$$\begin{aligned} Z^{[1]} &= W^{[1]} X + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \quad \text{conducted element-wise} \\ Z^{[2]} &= W^{[2]} A^{[1]} + b^{[2]} \quad \text{vectorized across your training set} \\ A^{[2]} &= g^{[2]}(Z^{[2]}) = \sigma(Z^{[2]}) \quad \text{(That's why capital } Z \text{ in } ZWX) \end{aligned}$$

activation function is sigmoid
when doing binary classification

Back Propagation (derivatives) "right to left"

$$\begin{aligned} dZ^{[2]} &= A^{[2]} - Y & Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}] \quad \text{again this is vectorized across all examples.} \\ dW^{[2]} &= \frac{1}{m} dZ^{[2]} A^{[1]T} & \text{prevents output of } (n, 1) \text{ rank one array} \\ db^{[2]} &= \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims=True}) & \uparrow \\ dZ^{[1]} &= W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]}) & (n^{[2]}, 1) \\ & & \text{element-wise product} \\ & (n^{(1)}, m) & (n^{(1)}, m) \end{aligned}$$

$$\begin{aligned} dW^{[1]} &= \frac{1}{m} dZ^{[1]} X^T \\ db^{[1]} &= \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims=True}) \\ & (n^{(1)}, 1) \quad \text{vector} \end{aligned}$$