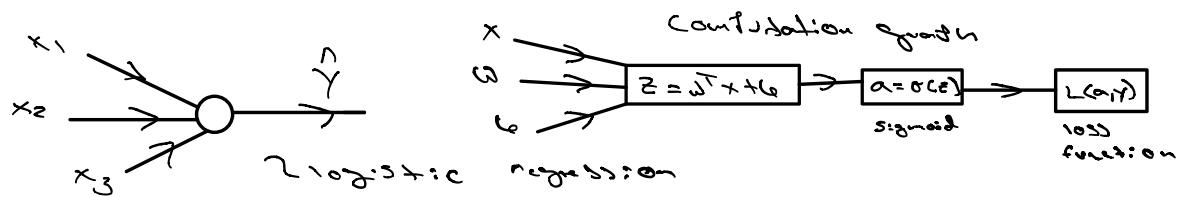
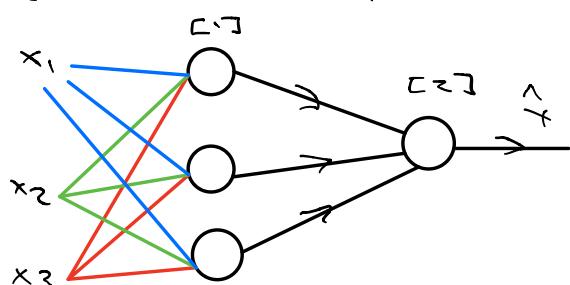


Week 3: Shallow Neural Networks

Overview



neural network

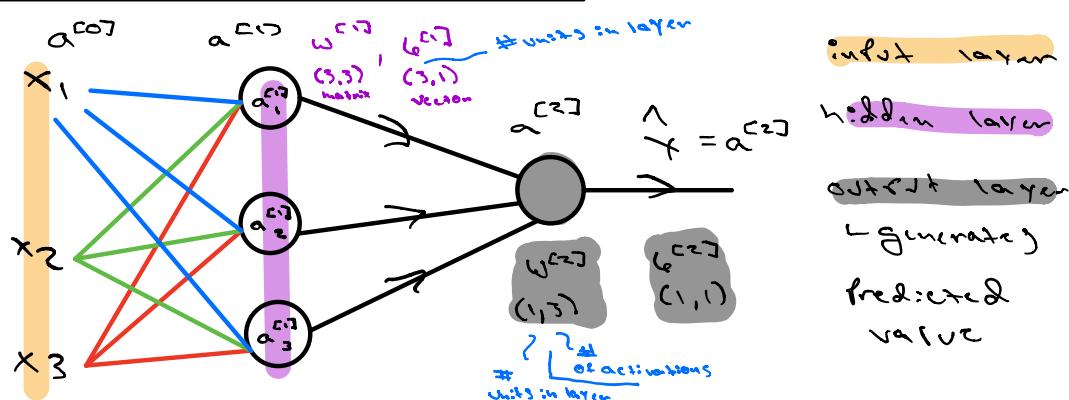


densely connected

z and a like calculations
happened at each node
neuron

notation: $w^{[l]}$; represents the weights for the first layer in the neural net (superscript)
 L round brackets $(\cdot) \rightarrow x^{[l]}$ are used to refer to training examples

neural network representation



input layer

hidden layer

output layer

generates

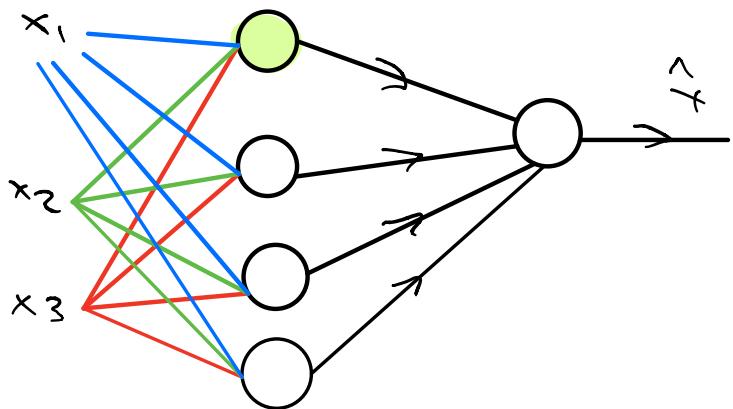
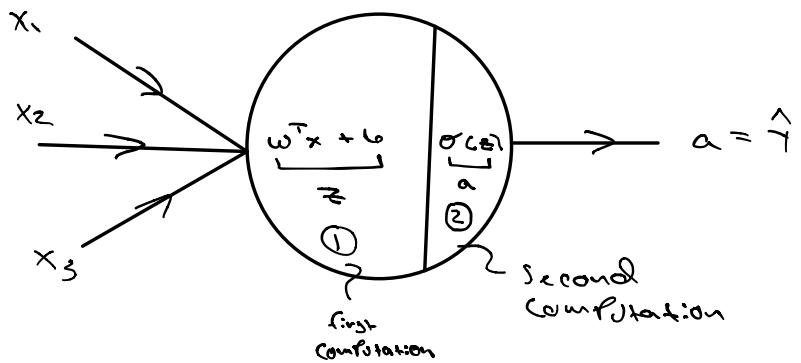
predicted

value

- the name hidden layer comes from the fact that the true value of the nodes in this middle layer are not observed in the training set.
 - ↳ while they are observed for the inputs and outputs.
- input vector representation: vector X or $a^{[0]}$
 - ↳ a stands for activations. These represent the values that layers pass on to subsequent layers.
 - ↳ $a^{[0]}$ is passed to the hidden layer
- the hidden layer generates a set of activations, $a^{[1]}$. Given that we have 3 nodes in this layer, $a^{[1]}$ is a 3×1 dimensional vector.
 - ↳ $a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \end{bmatrix} \approx 3$ hidden units in the layer
- the example above is a 2-layer neural network, as we don't count the input layer.
 - ↳ conventional usage definition

Computing a Neural Network's Output

Computing a neural network is like logistic regression but computed a lot of times.



- Focusing on the first neuron of the first layer

$$\text{Step 1: } z_i^{(1)} = w_i^{(1)T} x + b_i \quad \text{layer } i, \text{ node } i \text{ in layer } i$$

$$\text{Step 2: } a_i^{(1)} = \sigma(z_i^{(1)})$$

- These steps repeat for $z_2^{(1)} \rightarrow a_2^{(1)}$, $z_3^{(1)} \rightarrow a_3^{(1)}$
and $z_4^{(1)} \rightarrow a_4^{(1)}$

↪ this process is repetitive so, how can we vectorize?

ω^T transforms a column vector into a row

$$\begin{array}{c}
 \text{vector} \\
 \sim 4 \times 3 \\
 \left[\begin{array}{c} w_1^{C,T} \\ w_2^{C,T} \\ w_3^{C,T} \\ w_4^{C,T} \end{array} \right] \sim 4 \times 1 \\
 \left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right] \sim 3 \times 1 \\
 + \left[\begin{array}{c} b_1 \\ b_2 \\ b_3 \\ b_4 \end{array} \right] \sim 4 \times 1 \\
 = \left[\begin{array}{c} w_1^{C,T} x + b_1 \\ w_2^{C,T} x + b_2 \\ w_3^{C,T} x + b_3 \\ w_4^{C,T} x + b_4 \end{array} \right] \sim 4 \times 1
 \end{array}$$

- We have 4 logistic regression units in this matrix. Each of these logistic regression units has a parameter vector ω , we end up with a 4×3 matrix. (4 units, 3 activations per unit)

$$\left[\begin{array}{c} z_1 \\ z_2 \\ z_3 \\ z_4 \end{array} \right] = z$$

$$\Theta(z) = \left[\begin{array}{c} \alpha_1 \\ \vdots \\ \alpha_4 \end{array} \right] = \alpha$$

- We need a function that takes z and applies the sigmoid function element-wise.

Simplified: σ or $\sigma^{(0)}$

$$z^{(0)} = \omega^{(0)} x + b^{(0)} \quad \sim \text{first term}$$

$$\alpha^{(0)} = \sigma(z^{(0)})$$

$$z^{(1)} = \omega^{(1)} \alpha^{(0)} + b^{(1)} \quad \sim \text{second term}$$

$$\alpha^{(1)} = \gamma = \sigma(z^{(1)})$$

* this applies to one training example

Vectorizing across multiple examples

- To train with more training examples ($x^{(1)} \dots x^{(m)}$)

we need to repeat the forest summarized above.

$$x \rightarrow \alpha^{[c]} = \hat{y}$$

$$x^{(1)} \rightarrow \alpha^{[c](1)} = \hat{y}^{(1)}$$

$$x^{(2)} \rightarrow \alpha^{[c](2)} = \hat{y}^{(2)}$$

:

$$x^{(m)} \rightarrow \alpha^{[c](m)} = \hat{y}^{(m)}$$

training examples layer

- we would like to vectorize the computation of training on m examples as we don't want to loop per training example applying the same four equations above.

Standard for loop implementation:

for $i=1$ to m :

$$z^{(i)(c)} = w^{[c]} x^{(i)} + b^{[c]}$$

$$\alpha^{[c](i)} = \sigma(z^{(i)(c)})$$

$$\hat{y}^{(i)(c)} = w^{[c]} \alpha^{[c](i)} + b^{[c]}$$

$$y^{(i)(c)} = \alpha^{[c](i)} = \sigma(z^{(i)(c)})$$

To vectorize the training set we can stack them up in columns to create a $3 \times m$ matrix

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots x^{(m)} \\ | & | & | \end{bmatrix} \quad \text{for our example}$$

$$(n \times m) = (\text{features}, \text{examples})$$

The vectorized implementation of the for loop

is never,

$$z^{[1]} = w^{[1]} \times + b^{[1]} \quad * \text{ lowercase } z, w, x,$$

$$A^{[1]} = \sigma(z^{[1]}) \quad A \text{ as they're all matrices.}$$

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

$$A^{[2]} = \sigma(z^{[2]}) \quad b \text{ is a column vector}$$

$$z^{[1]} = \begin{bmatrix} | & | & | \\ z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \\ | & | & | \end{bmatrix} \quad \sim \text{similar for } z^{[2]} ?$$

$$A^{[1]} = \begin{bmatrix} | & | & | \\ a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \\ | & | & | \end{bmatrix} \quad \begin{array}{l} \text{hidden unit} \\ \# \text{ per layer index} \end{array}$$

the horizontal index corresponds to training examples. The vertical index correspond to different nodes in the layer being analyzed.

Justification for vectorized Implementation

$$w^{[1]} = \begin{bmatrix} \parallel & \parallel & \parallel \\ \parallel & \parallel & \parallel \\ \parallel & \parallel & \parallel \end{bmatrix}$$

$$w^{[1]} x^{[1]} = \begin{bmatrix} : \\ : \\ : \end{bmatrix}$$

Each line represents the weights connected to one neuron in the layer

When multiplied by a column vector of input (features), rows \times columns, the first row $w^{[1]}$ is multiplied against the column $x^{[1]}$ and then summed up into one data point

(this is done for all rows in $w^{[1]}$). Hence, each data point $x_1^{[1]}, x_2^{[1]}, x_3^{[1]} \dots x_n^{[1]}$ is scaled by a different $w^{[1]}$ value ($w_{11}, w_{21}, w_{31} \dots w_{n1}$ (neuron in layer 1)).

here, we are scaling the same data point by different weights depending on which neuron of the layer we are evaluating

↳ this applies for any activation value not just $a^{[0]} = x$

*one neuron in
a layer*

- The sample calculation of one $w^{[1]}$ row with one training example is given as:

the sum of a neuron's z with $b=0$

$$(z_1^{[1]} = w_{1-nx} \cdot x_{nx})$$

$$\hookrightarrow w_{11} \cdot x_1 + w_{12} \cdot x_2 + w_{13} \cdot x_3$$

$$(z_2^{[1]} = w_{2-nx} \cdot x_{nx})$$

$$\hookrightarrow w_{21} \cdot x_1 + w_{22} \cdot x_2 + w_{23} \cdot x_3$$

when the example are stacked vertically

such that

$$x = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \\ \dots & & \end{bmatrix}$$

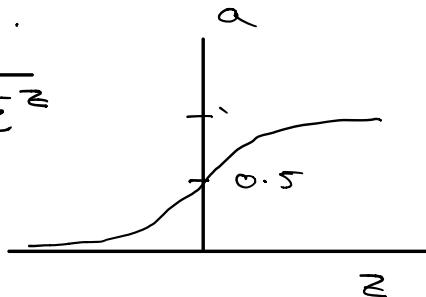
then $w^{(1)} \times w^{(0)}$ be ~ a matrix

where each column is the z for a different training example

Activation Functions

- So far we have been using the Sigmoid activation function, but sometimes, other choices can work much better.

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

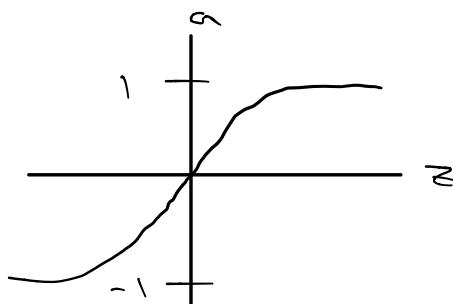


- Instead of $a = \sigma(z)$

we could have $a = g(z)$, where g could be other linear functions which are not the Sigmoid.

- Tanh Function (hyperbolic tan)

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



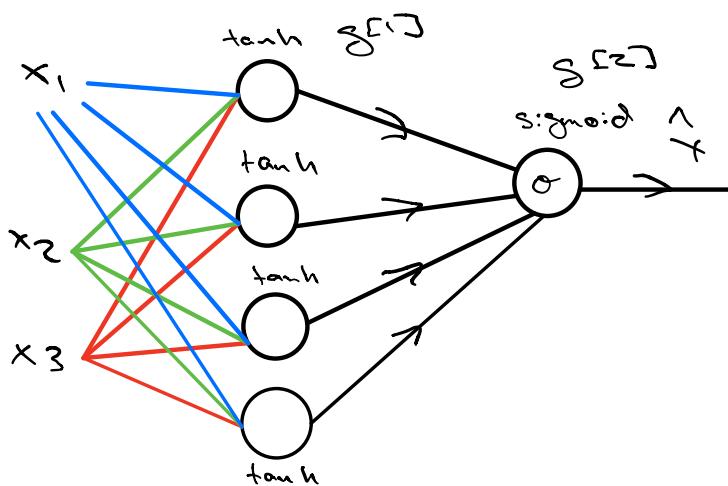
L Using tanh instead of the Sigmoid for the hidden units almost always works better because within values $[1, -1]$, the mean of

The activations that come out of the layer
are closer to zero.

↳ the data is centered around zero
instead of 0.5 making learning
a bit easier.

↳ the one exception for using the sigmoid
function is the output layer because y is
either zero or one and therefore it
makes sense for f to be a number
in this range. (in binary classification)

- to indicate that the activation functions might
be different for different layers, we utilize
another **super script** ($\delta^{(l)}(z^{(l)})$)



- one of the major downsides to both the sigmoid
and tanh is that if z is either very large
or very small, then the gradient of these

functions become very small.

↳ close to zero, slowing down gradient descent.

- Given this drawback, the rectified linear unit is used as an alternative (ReLU)

↳ ReLU: $a = \max(0, z)$

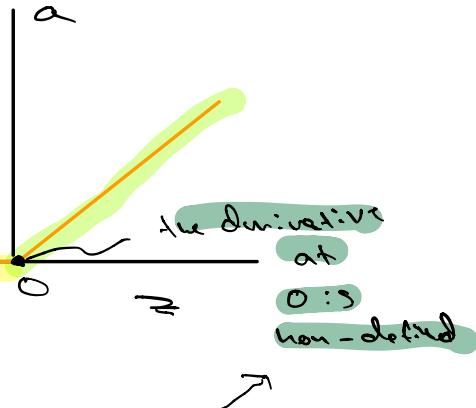
↳ the derivative is 1 when z is positive and

0 when $z \rightarrow$ negative

↳ but we can

use either one

or zero derivative for the zero point



- One disadvantage of the ReLU is that its derivative is 0 for negative. In practice this isn't an issue, but an alternative called the Leaky ReLU is sometimes used.

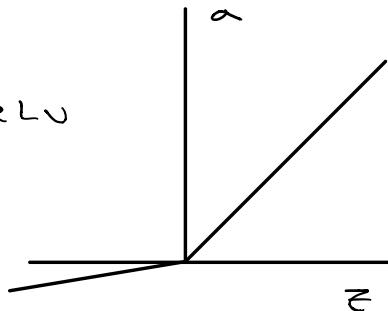
Rules of thumb for choosing Activation function:
∴ Binary classification: Sigmoid at output layer is the natural choice

↳ for all other units the ReLU is used
(in the hidden layer) ... rarely tanh

- Leaky ReLU:

↳ usually works better than ReLU

$$\hookrightarrow a = \max(0.01 \cdot z, z)$$



Why do we need non-linear activation functions?

- to compute interesting functions, we require non-linear activation functions

- if we set $\sigma(z) = z$, we are computing our prediction as a linear function of our input features

$$\hookrightarrow a^{c_1} = z^{c_1} = w^{c_1}x + b^{c_1}$$

$$\hookrightarrow a^{c_2} = z^{c_2} = w^{c_2}a^{c_1} + b^{c_2}$$

$$\begin{aligned} \hookrightarrow a^{c_2} &= w^{c_2} (w^{c_1}x + b^{c_1}) + b^{c_2} \\ &= (\underbrace{w^{c_2}w^{c_1}}_{w'})x + (\underbrace{w^{c_2}b^{c_1} + b^{c_2}}_{b'}) \end{aligned}$$

$$= w'x + b'$$

* using linear activations (identity activations), the neural network will only output a linear function of the input

↳ stacking more layers in a deep neural net won't have any effect as we are always calculating a linear function of the

injust).

↳ a linear hidden unit is useless because the composition of two linear functions is itself a linear function (theorem of superposition).

- We might only use linear activations if we are doing machine learning on a regression problem.

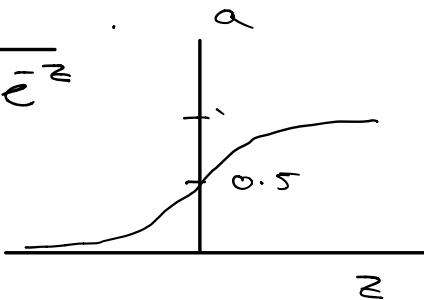
↳ the \hat{y} will be a real number in this case so it's sensible.

↳ but the hidden units must be non-linear (such as ReLU or tanh)

↳ linear units are sometimes used for compression but not discussed in the course.

Derivatives of activation functions

- Sigmoid: $\sigma(z) = \alpha = \frac{1}{1 + e^{-z}}$



$$\frac{d}{dz} \sigma(z) = \sigma(1 - \sigma)$$

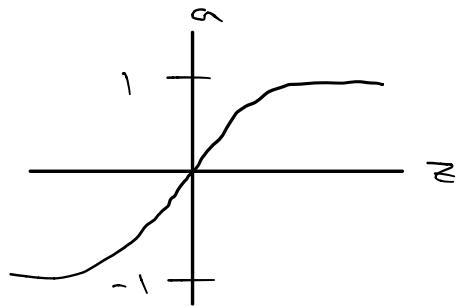
$$= \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}} \right) = \sigma'(z)$$

The advantage of this formula is that if we've already computed the activation (σ),

It's very easy to compute the slope as $a(1-a)$.

- tanh function (hyperbolic tan)

$$\begin{aligned} \hookrightarrow a = \tanh(z) = \\ \frac{e^z - e^{-z}}{e^z + e^{-z}} \end{aligned}$$

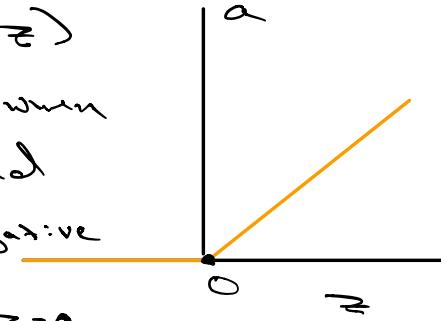


$$\frac{d}{dz} g(z) = g'(z) = 1 - (\tanh(z))^2$$

\hookrightarrow this formula again can be used to quickly compute the derivative of tanh

- RELU: $a = \max(0, z)$

\hookrightarrow the derivative is 1 when z is positive and 0 when z is negative



\hookrightarrow undefined if $z=0$

Gradient Descent for Neural Networks

- a neural network with a 3-layer hidden layer has parameters $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$

\sim input features \sim hidden units \sim output layer
 $n_x = n^{[0]}$ $n^{[1]}$ $n^{[2]} = 1$

$$\begin{aligned} & \hookrightarrow w^{(1)} = (w^{(1)}, b^{(1)}) , \quad b^{(1)} = (b^{(1)}, 1) \\ & \hookrightarrow w^{(2)} = (w^{(2)}, b^{(2)}) , \quad b^{(2)} = (b^{(2)}, 1) \\ & \text{number of } \downarrow \quad \text{connections to previous layer} \\ & \text{neurons in layer} \end{aligned}$$

- For binary classification, we have a cost function
 $J(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}) = \frac{1}{m} \cdot \sum_{i=1}^m L(\hat{y}_i, y_i)$

* $\hat{y} = a^{(2)}$ loss function average.

* When we have a neural network it's better to initialize all the parameters randomly (instead of set 0).

- Gradient Descent:

Repeat {

- Compute predictions ($\hat{y}^{(i)}$, for $i=1$ through m)
- Compute derivatives

$$\hookrightarrow \frac{\partial J}{\partial w^{(1)}} = \frac{\partial J}{\partial w^{(1)}}, \quad \frac{\partial J}{\partial b^{(1)}} = \frac{\partial J}{\partial b^{(1)}} \dots \text{for } w^{(2)} \text{ and } b^{(2)}$$

- Update parameters:

$$\hookrightarrow w^{(1)} = w^{(1)} - \alpha \cdot \frac{\partial J}{\partial w^{(1)}}$$

$$\hookrightarrow b^{(1)} = b^{(1)} - \alpha \cdot \frac{\partial J}{\partial b^{(1)}}$$

$$\hookrightarrow w^{(2)} = w^{(2)} - \alpha \cdot \frac{\partial J}{\partial w^{(2)}}$$

$$\hookrightarrow b^{(2)} = b^{(2)} - \alpha \cdot \frac{\partial J}{\partial b^{(2)}}$$

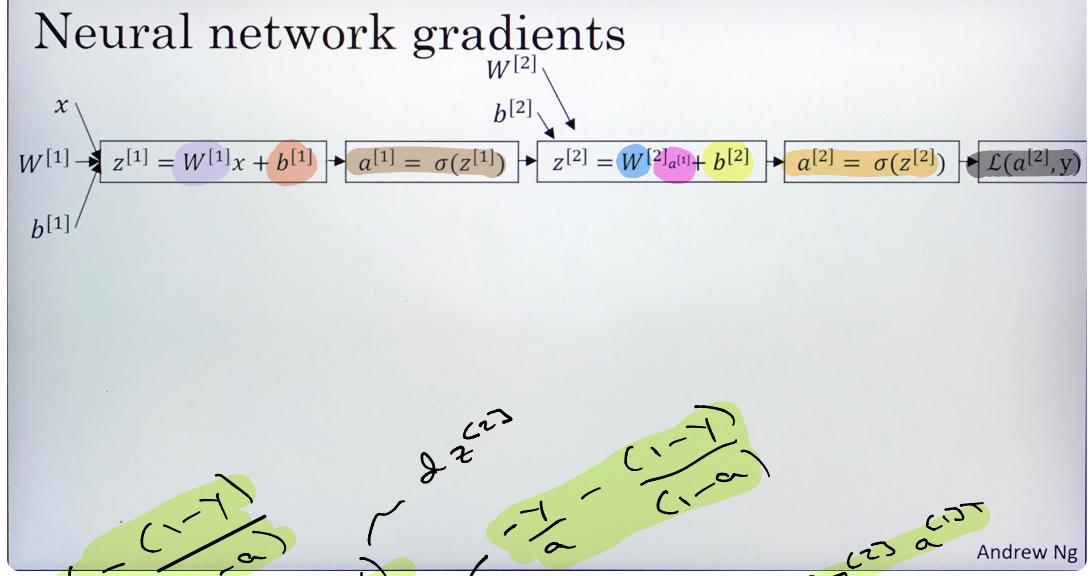
}

- Now do we compute these partial derivatives?

↪ forward propagation:

- Vertical margin
- ↳ $z^{(1)} = w^{(1)} X + b^{(1)}$
 - ↳ $A^{(1)} = g^{(1)}(z^{(1)})$
 - ↳ $z^{(2)} = w^{(2)} A^{(1)} + b^{(2)}$ sigmoid for linear classification
 - ↳ $A^{(2)} = g^{(2)}(z^{(2)}) = \sigma(z^{(2)})$
 - ↳ Backpropagation: g term back in sigmoid func ground truth ~ row vector
 $\left\{ \begin{array}{l} \delta z^{(2)} = A^{(2)} - y \\ \delta w^{(2)} = \frac{1}{m} \cdot \delta z^{(2)} A^{(1)\top} \\ \delta b^{(2)} = \frac{1}{m} \cdot \text{np.sum}(\delta z^{(2)}, \text{axis}=1, \text{keepdims=True}) \\ \quad ? \text{horizontal sum } \sum_{(n^{(2)}, 1)} \end{array} \right.$
 $y(1 \times m) = [y^{(1)} \dots y^{(m)}]$
 - ↳ $\delta z^{(1)} = W^{(2)\top} \delta z^{(2)} * g^{(1)\prime}(z^{(1)}) \quad (n^{(1)}, m)$ element wise product
 - ↳ $\delta w^{(1)} = \frac{1}{m} \cdot \delta z^{(1)} X^\top \sim A^{(0)}$
 - ↳ $\delta b^{(1)} = \frac{1}{m} \cdot \text{np.sum}(\delta z^{(1)}, \text{axis}=1, \text{keepdims=True})$

Computational Graph of 1-layer neural network



$$\begin{aligned}
 \frac{\partial L}{\partial a^{c_2}} &\rightarrow \frac{\partial L}{\partial z^{c_2}} = \frac{\partial L}{\partial a^{c_1}} \cdot \frac{\partial a^{c_1}}{\partial z^{c_2}} \rightarrow \frac{\partial L}{\partial w^{c_2}} = \frac{\partial L}{\partial z^{c_2}} \cdot \frac{\partial z^{c_2}}{\partial w^{c_2}} \\
 \text{and } \frac{\partial L}{\partial b^{c_2}} &= \frac{\partial L}{\partial z^{c_2}} \cdot \frac{\partial z^{c_2}}{\partial b^{c_2}} \rightarrow \frac{\partial L}{\partial a^{c_1}} = \frac{\partial L}{\partial z^{c_2}} \cdot \frac{\partial z^{c_2}}{\partial a^{c_1}} \\
 \rightarrow \frac{\partial L}{\partial z} &= \frac{\partial L}{\partial a^{c_1}} \cdot \frac{\partial a^{c_1}}{\partial z} \rightarrow \frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w} \\
 \text{and } \frac{\partial L}{\partial b} &= \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial b}
 \end{aligned}$$

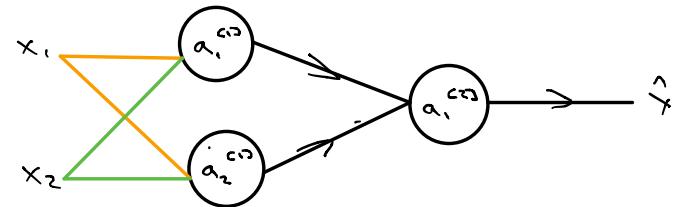
- Vectorized Implementation of Backpropagation
(column vectors)

$$\begin{aligned}
 dZ^{[2]} &= A^{[2]} - Y \\
 dW^{[2]} &= \frac{1}{m} dZ^{[2]} A^{[1]T} \quad \text{← transposed row vector} \\
 db^{[2]} &= \frac{1}{m} np.sum(dZ^{[2]}, axis = 1, keepdims = True) \\
 dZ^{[1]} &= W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]}) \quad \text{elementwise product} \quad \text{← each element with one} \\
 dW^{[1]} &= \frac{1}{m} dZ^{[1]} X^T \quad \text{other element in the} \\
 db^{[1]} &= \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True) \quad \text{sum placed in the other} \\
 &\quad \text{matrix}
 \end{aligned}$$

Andrew Ng

Random Initialization

- in a neural network, if we initialize the parameters to zero and then apply gradient descent, it won't work.



$$n^{(c1)} = 2 \quad h^{(c1)} = 2$$

therefore, $w^{(c1)} = (2, 2)$

this is a problem *this is fine*

$$w^{(c1)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{(c1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

- the problem with the initialization to zero is that for any given example, $a_1^{(c1)}$ and $a_2^{(c1)}$ are equal
 - ↳ both hidden units will compute the same function
 - ↳ when computing local propagation, $\delta z^{(c1)} = \delta z^{(c2)}$ by symmetry.
 - ↳ both hidden units are computing the same function
 - ↳ through induction, it's proved that after multiple iterations of gradient descent, these will be computing the same function
 - ↳ $\delta w = \begin{bmatrix} v & v \\ v & v \end{bmatrix} \rightarrow w^{(c1)} = w^{(c2)} - \alpha \delta w^{(c1)}$
all rows will be equal after each update. One point is having more than

one neuron) (symmetry breaking problem)

- the solution to this \Rightarrow to initialize parameters randomly.

$$\leftarrow w^{(1)} = \text{np.random.rand}(2, 2) * 0.01$$

\leftarrow gaussian random variables

$$\leftarrow b^{(1)} = \text{np.zeros}(2, 1)$$

* we multiply by a small # because we want to have activations that are different when we input the values into our sigmoid or tanh (the sloped area provides better/faster performance in gradient descent).

\leftarrow a saturated sigmoid function slows down learning

Additional

* In general, the # of neurons in the previous layer provides the # of columns, while the # of neurons in the current layer gives us the # of rows of the weight matrix...

for the activation matrix, the row gives the neuron in a specific layer, while the columns give different examples.

The resulting $a^{(l)}$ \Rightarrow a compilation of the activations for the neurons in a layer (rows) with the activation for different examples (columns).