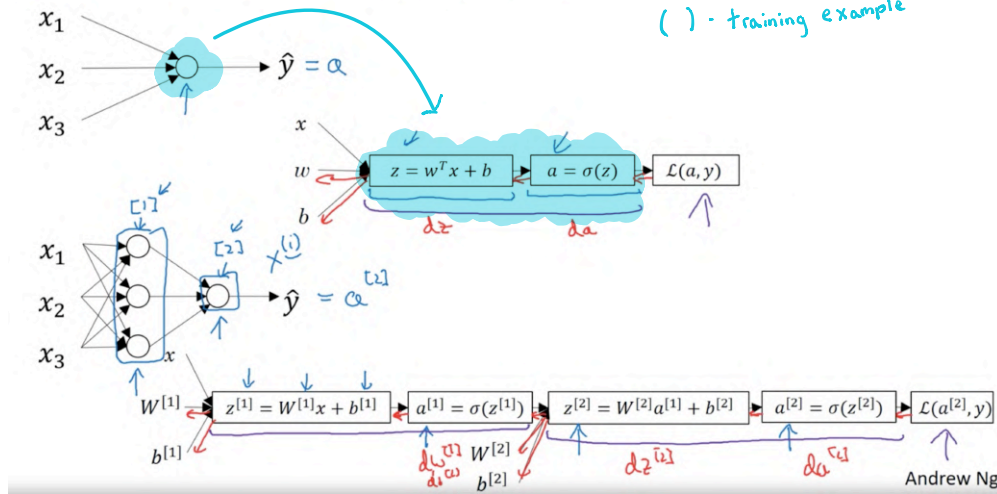
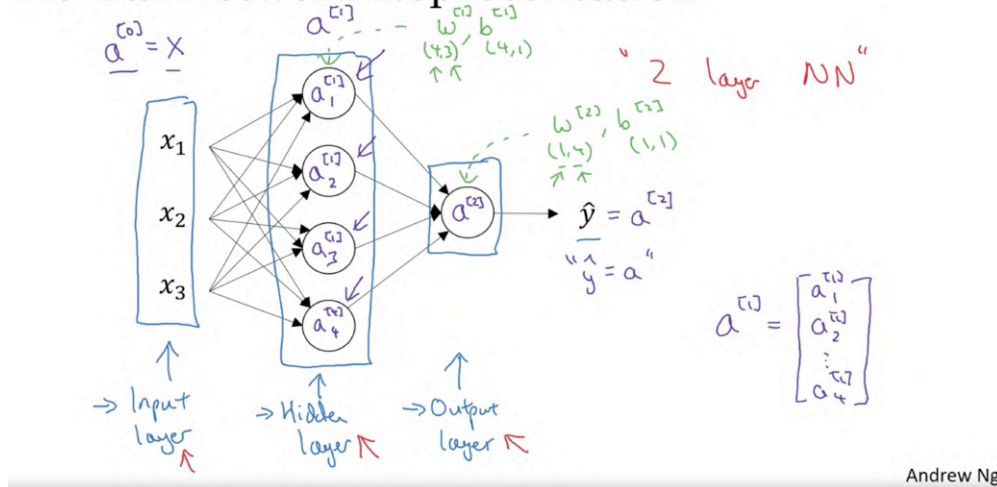


What is a Neural Network?

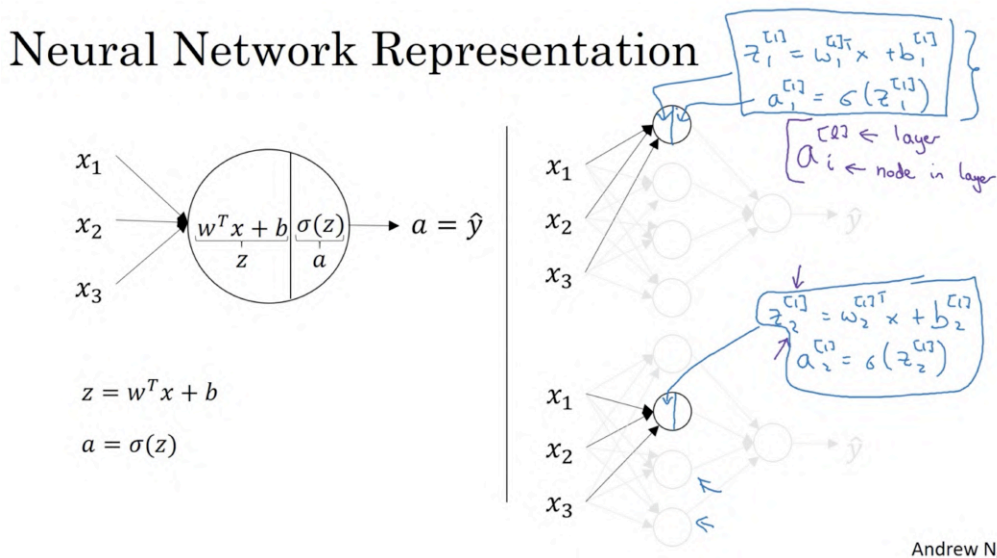
[] - layer
() - training example



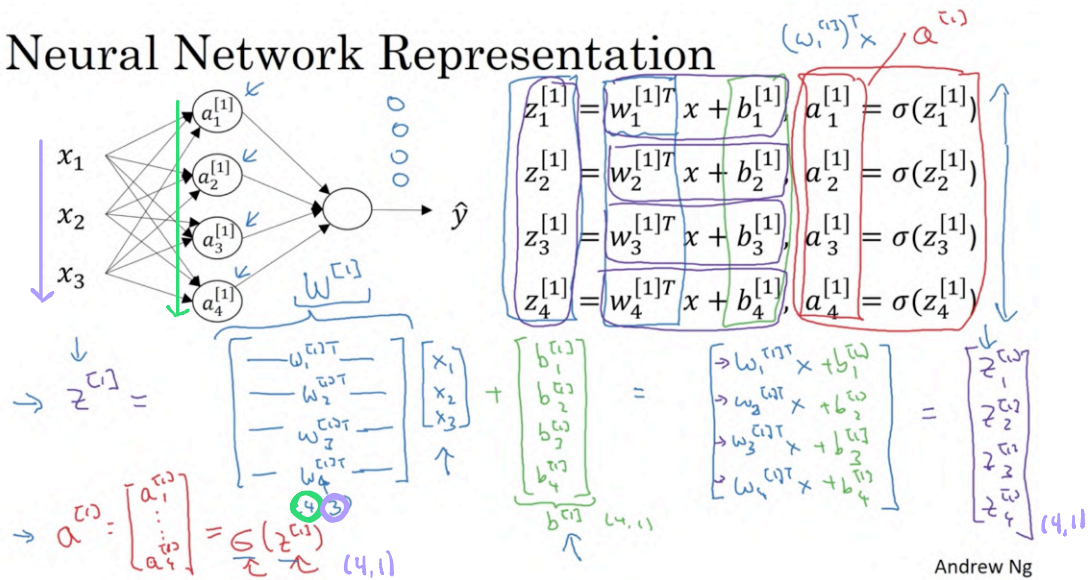
Neural Network Representation



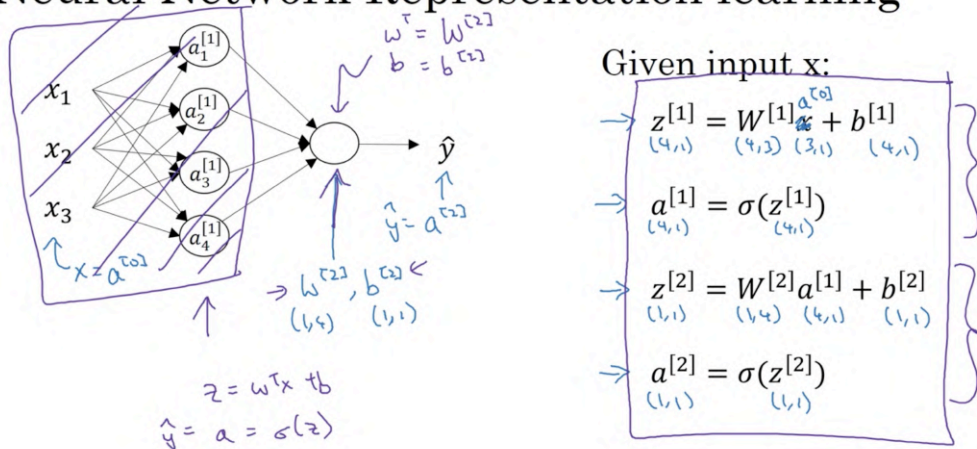
Neural Network Representation



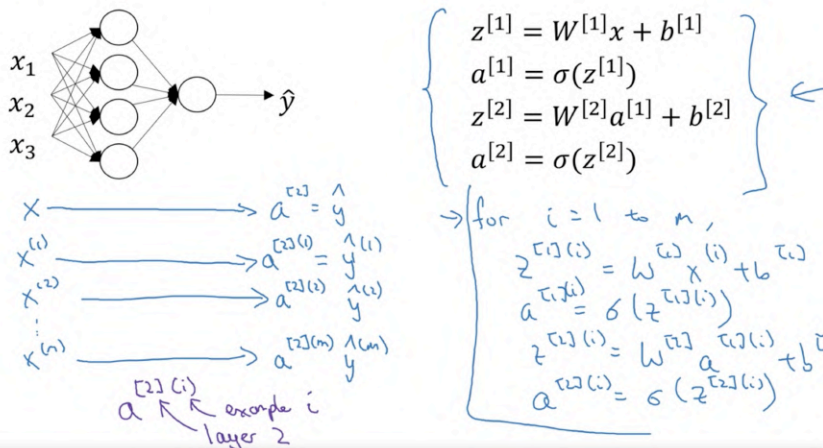
Neural Network Representation



Neural Network Representation learning



Vectorizing across multiple examples



Vectorizing across multiple examples

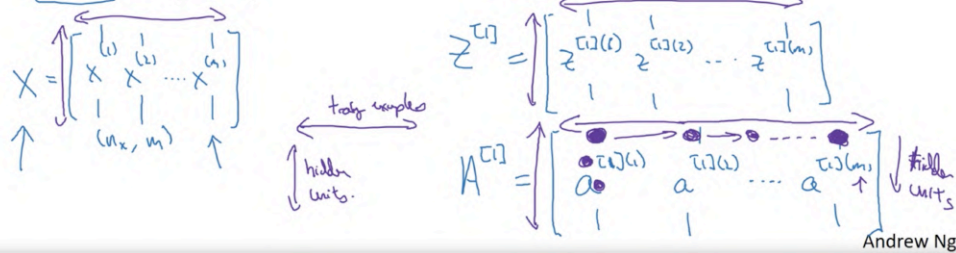
for $i = 1$ to m :

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$

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

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

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



Andrew Ng

Justification for vectorized implementation

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}, \quad z^{[1](2)} = W^{[1]}x^{(2)} + b^{[1]}, \quad z^{[1](m)} = W^{[1]}x^{(m)} + b^{[1]}$$

$$W^{[1]} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}, \quad W^{[1]}x^{(i)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}, \quad W^{[1]}x^{(2)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}, \quad W^{[1]}x^{(m)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

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

Diagram illustrating the justification for vectorized implementation. It shows the matrix multiplication $W^{[1]}X$ resulting in the vectorized output $Z^{[1]}$. The bias vector $b^{[1]}$ is added to each column of $Z^{[1]}$ to produce the final output $Z^{[1]}$.

Andrew Ng

Recap of vectorizing across multiple examples

Diagram illustrating the recap of vectorizing across multiple examples. The input matrix X is shown as a matrix of size (n_x, m) . The hidden layer activation matrix A is shown as a matrix of size (n_h, m) . The output matrix Z is shown as a matrix of size (n_o, m) . The equations for the forward pass are summarized as follows:

$$\text{for } i = 1 \text{ to } m$$

$$\left\{ \begin{array}{l} z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]} \\ a^{[1](i)} = \sigma(z^{[1](i)}) \\ z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]} \\ a^{[2](i)} = \sigma(z^{[2](i)}) \end{array} \right.$$

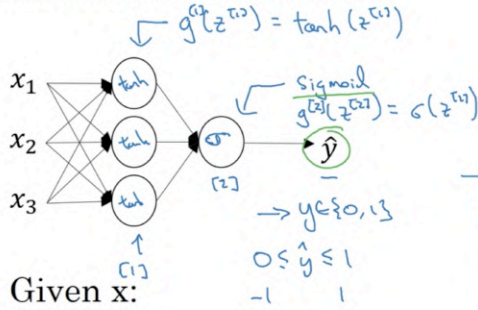
$$\left\{ \begin{array}{l} Z^{[1]} = W^{[1]}X + b^{[1]} \\ A^{[1]} = \sigma(Z^{[1]}) \\ Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} = \sigma(Z^{[2]}) \end{array} \right.$$

The diagram also shows the vectorized equations for the forward pass, including the bias vectors $b^{[1]}$ and $b^{[2]}$, and the activation function σ .

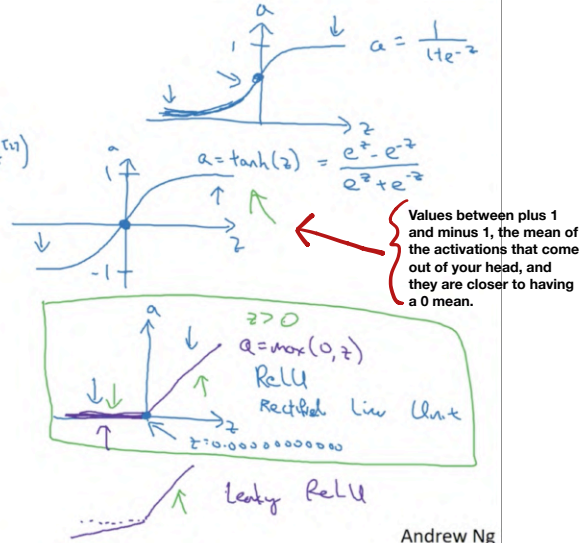
Andrew Ng

Activation functions

Tanh is almost always strictly better than sigmoid. Only exception is for the output layer for binary classification because if y is either 1 or 0, it makes more sense to be between 0 and 1 than 1 and -1.

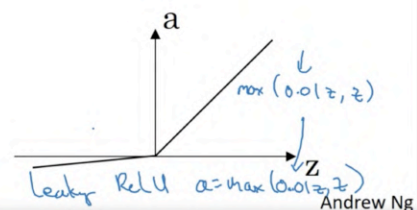
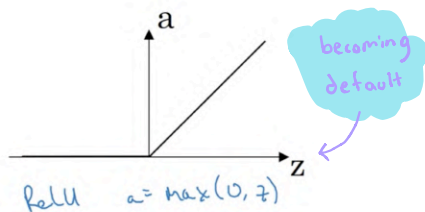
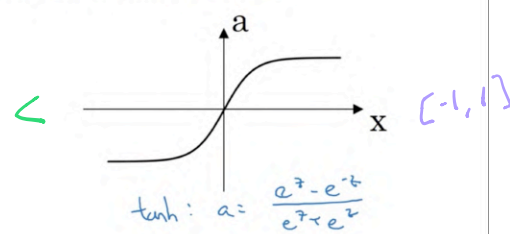
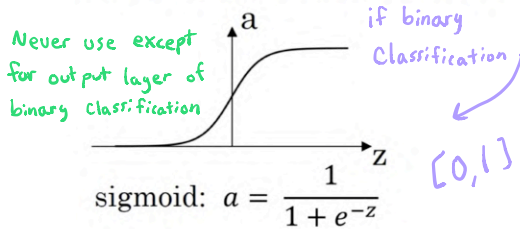


One of the downsides of both the sigmoid function and the tanh function is that if z is either very large or very small, then the gradient or the derivative or the slope of this function becomes very small. So if z is very large or z is very small, the slope of the function ends up being close to 0. And so this can slow down gradient descent.



Andrew Ng

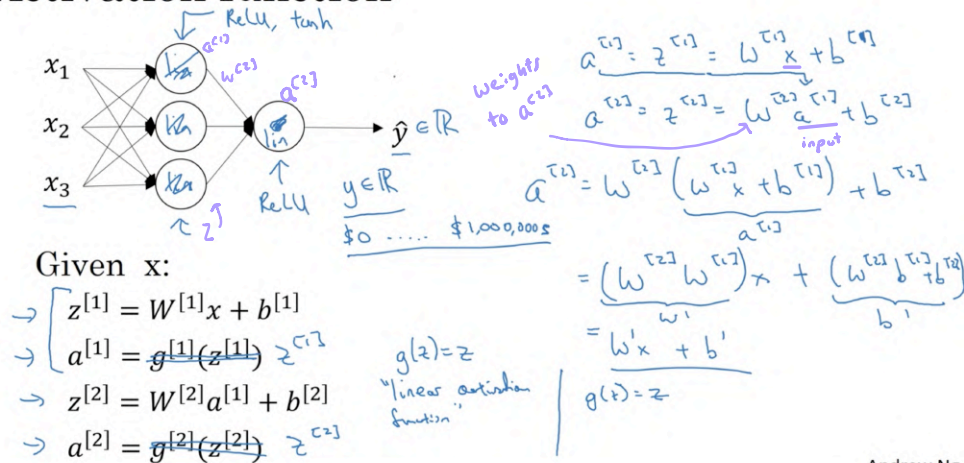
Pros and cons of activation functions



Andrew Ng

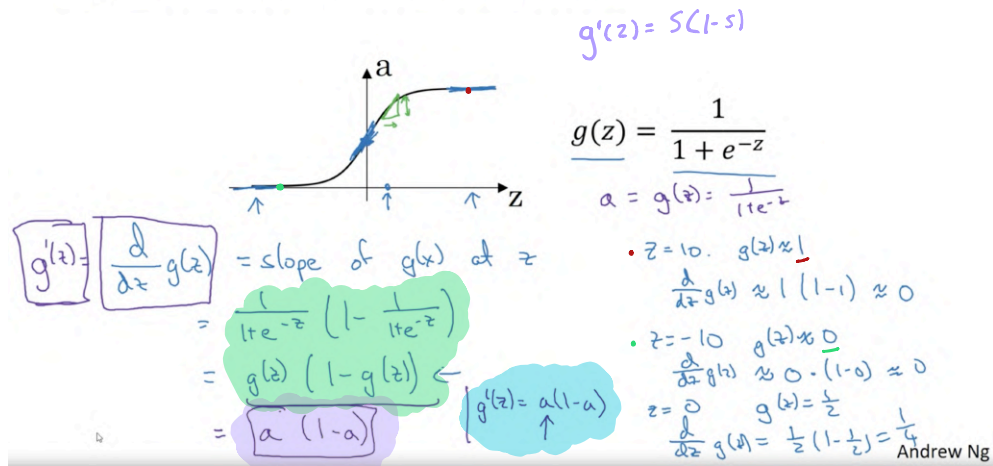
Activation function

Why do we need non-linearity?

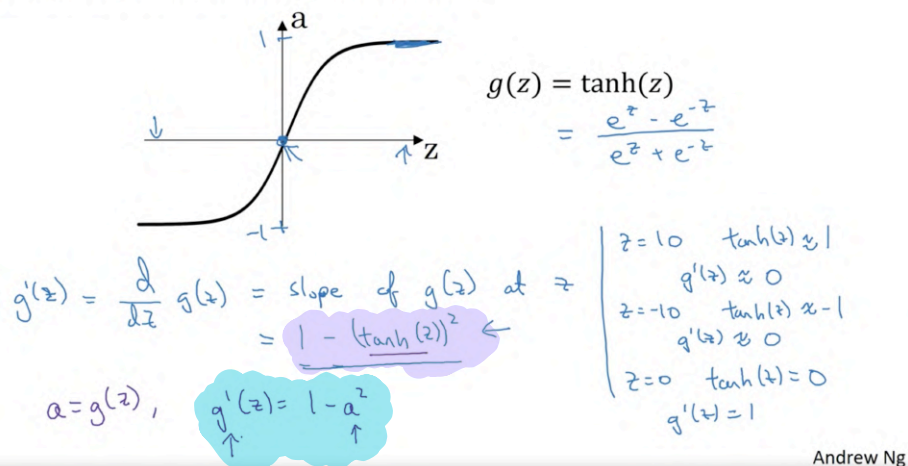


Andrew Ng

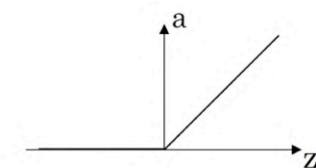
Sigmoid activation function



Tanh activation function



ReLU and Leaky ReLU

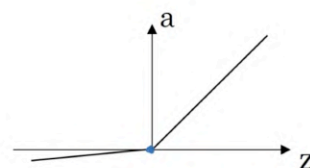


ReLU

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

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

$z = 0.0000000000$



Leaky ReLU

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

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

Gradient descent for neural networks

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$ $n_x = n^{[0]}, n^{[1]}, n^{[2]} = 1$

Cost function: $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})$

Gradient descent:

→ Repeat {
 → Compute gradients $(\hat{y}^{(i)}, i=1, \dots, m)$
 $\frac{dJ}{dW^{[1]}} = \frac{\partial J}{\partial W^{[1]}} \quad \frac{dJ}{db^{[1]}} = \frac{\partial J}{\partial b^{[1]}} \dots$
 $W^{[1]} := W^{[1]} - \alpha \frac{dJ}{dW^{[1]}}$
 $b^{[1]} := b^{[1]} - \alpha \frac{dJ}{db^{[1]}}$

Formulas for computing derivatives

Forward propagation:

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

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

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

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

Back propagation:

$$dz^{[2]} = A^{[2]} - y$$

$$dW^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

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

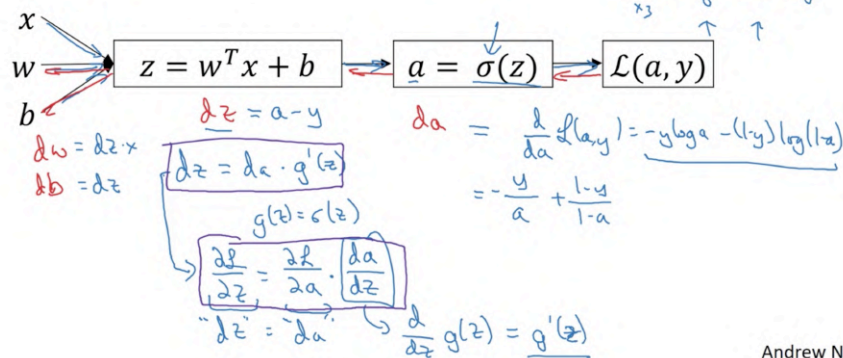
$$dz^{[1]} = \frac{W^{[2]T} dz^{[2]}}{(n^{[2]}, m)} \times \frac{g^{[1]'}(z^{[1]})}{(n^{[1]}, m)}$$

$$dW^{[1]} = \frac{1}{m} dz^{[1]} x^T$$

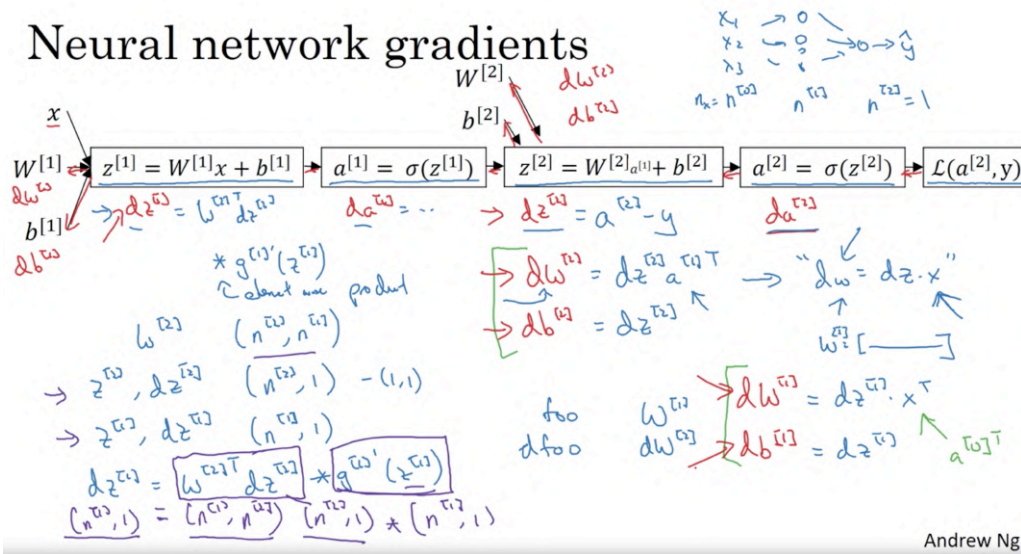
$$db^{[1]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims}=True)$$

Computing gradients

Logistic regression



Neural network gradients



Andrew Ng

Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

Vectorized implementation:

$$\begin{aligned} z^{[2]} &= W^{[2]} x + b^{[2]} \\ a^{[2]} &= g^{[2]}(z^{[2]}) \\ z^{[1]} &= \begin{bmatrix} z^{[1]}(1) & z^{[1]}(2) & \dots & z^{[1]}(n_1) \end{bmatrix} \\ z^{[1]} &= W^{[1]} x + b^{[1]} \\ a^{[1]} &= g^{[1]}(z^{[1]}) \end{aligned}$$

Andrew Ng

Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

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

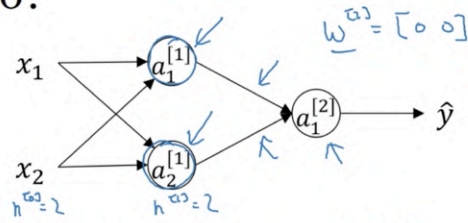
$$dZ^{[1]} = \frac{1}{m} W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

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

Andrew Ng

What happens if you initialize weights to zero?



$$w_k^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

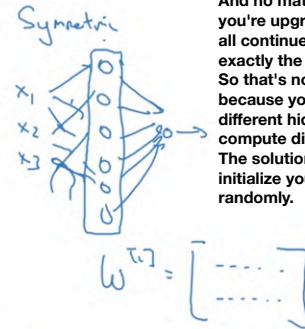
$$b_k^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]}$$

$$dz_1^{[1]} = dz_2^{[1]}$$

$$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

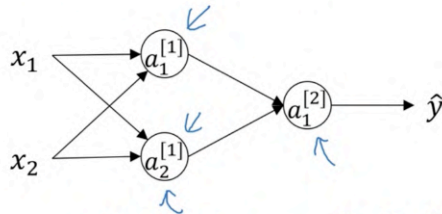
$$w^{[1]} = w^{[0]} - \alpha dw$$



If you initialize the weights to zero, then all of your hidden units are symmetric. And no matter how long you're upgrading the center, all continue to compute exactly the same function. So that's not helpful, because you want the different hidden units to compute different functions. The solution to this is to initialize your parameters randomly.

Andrew Ng

Random initialization



$$w^{[0]} = \text{np.random.randn}(2,2) * 0.01$$

$$b^{[0]} = \text{np.zeros}(2,1)$$

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

$$b^{[1]} = 0$$



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

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

Andrew Ng