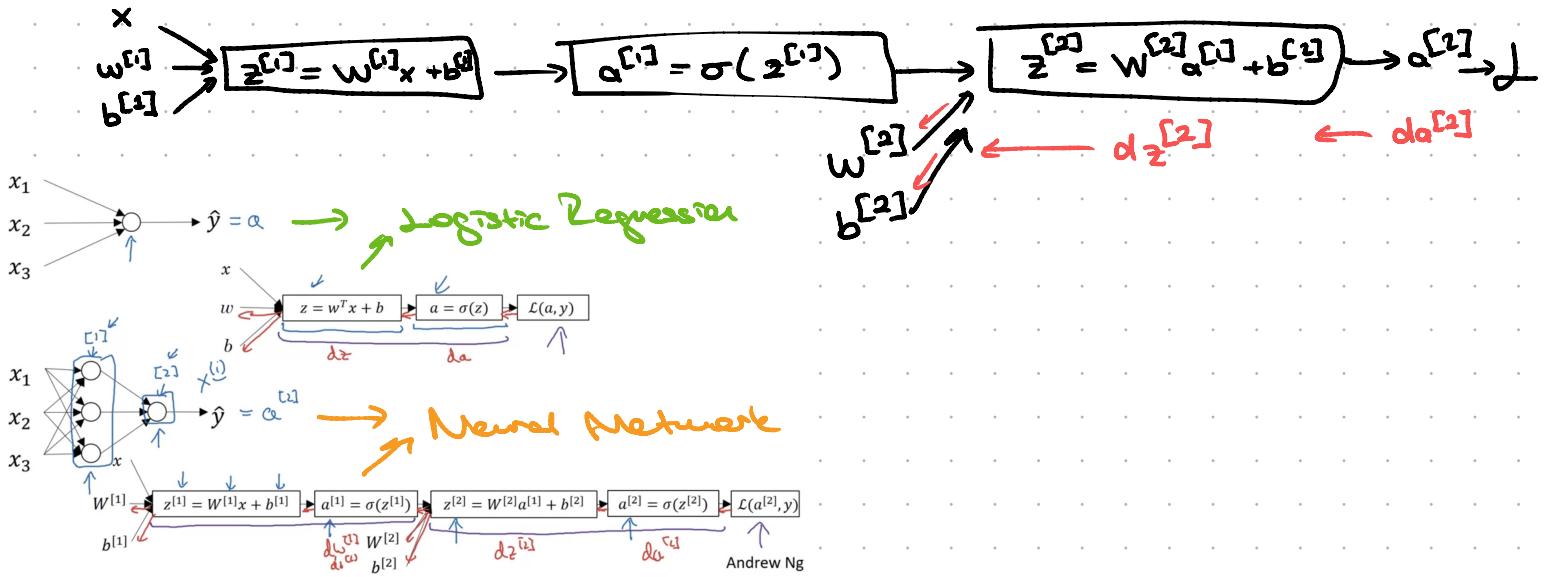
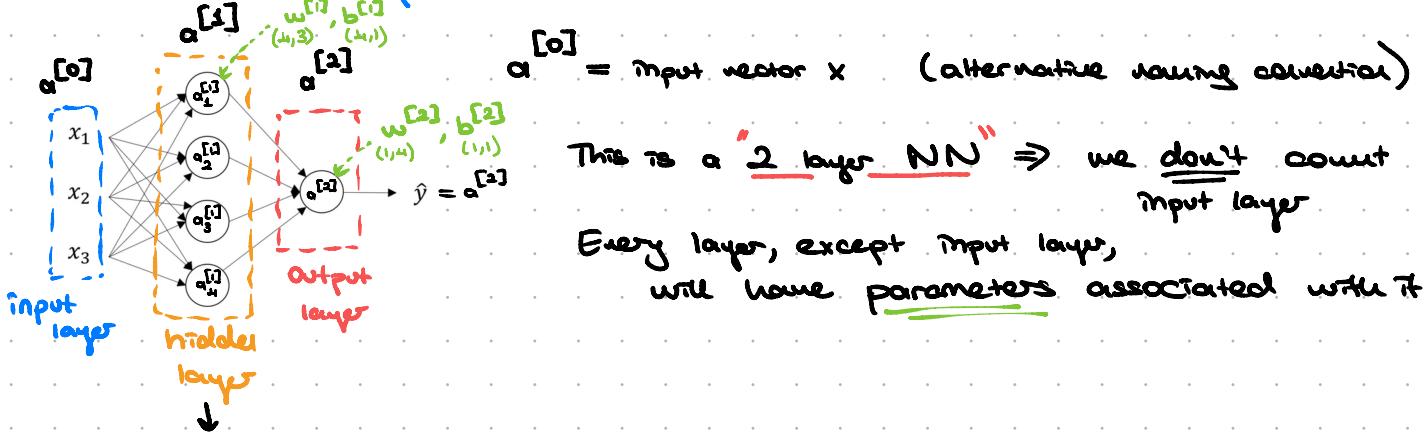


Neural Networks Overview



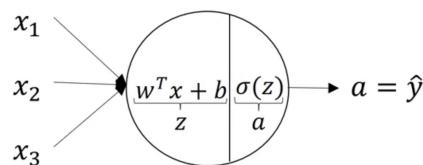
Neural Network Representation



In training set,
we don't see what
hidden layer should be
BUT we see input + output
layer values

Computing a NN Output

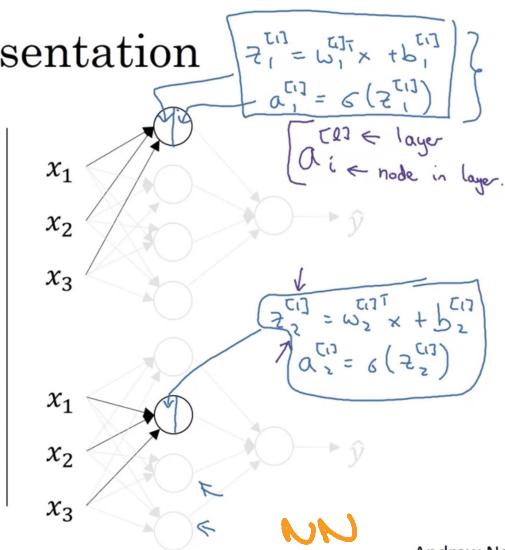
Neural Network Representation



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

$$a = \sigma(z)$$

Linear Regression



Andrew Ng

$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, & a_1^{[1]} &= \sigma(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, & a_2^{[1]} &= \sigma(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, & a_3^{[1]} &= \sigma(z_3^{[1]}) \end{aligned}$$

We need to vectorize this!!

Let $W^{[1]}$ be a matrix such that:

$$W^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ \vdots \\ w_{n_x}^{[1]T} \end{bmatrix}$$

Note

All of this is for a single training example.

Then,

$$Z^{[1]} = \underbrace{\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ \vdots \\ w_{n_x}^{[1]T} \end{bmatrix}}_{W^{[1]} (4,3)} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}}_{b^{[1]} (4,1)} = \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ \vdots \\ w_{n_x}^{[1]T} x + b_{n_x}^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

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

Given input x :

$$z^{[1]} = \sum_{(i,j)}^{[1]} x_i + b^{[1]}$$

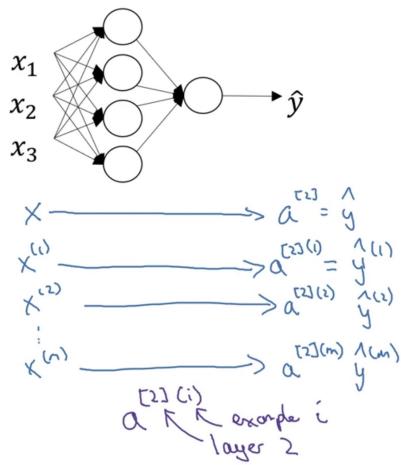
$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = \sum_{(i,j)}^{[2]} a^{[1]}_i + b^{[2]}$$

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

only 1 training example

Vectorizing Across Multiple Examples



without vectorization:

for $i=1$ to m , we want to get rid of this

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

RECALL

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

THEN

$$\begin{aligned} z^{[1]} &= W^{[1]} X + b^{[1]} \\ \text{hidden units} &\quad \text{training example} \\ A^{[1]} &= \sigma(z^{[1]}) \end{aligned}$$

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

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

$$z^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ \vdots \\ w_m^{[1]T} \end{bmatrix} \begin{bmatrix} | & | \\ x^{(1)} & \dots & x^{(m)} \\ | & | \end{bmatrix} + b^{[1]} = \begin{bmatrix} w_1^{[1]x^{(1)}} \\ w_2^{[1]x^{(1)}} \\ \vdots \\ w_m^{[1]x^{(1)}} \end{bmatrix} \quad \begin{matrix} \text{training examples} \\ \text{hidden units} \end{matrix}$$

$$W^{[1]} x^{(2)}, b^{[1]}, \dots, W^{[1]} x^{(m)}, b^{[1]}$$

Activation Functions

Sigmoid function is an activation function

- goes between 0 and 1

tanh \Rightarrow almost always better than sigmoid

- goes between -1 and +1

- mathematically a shifted version of sigmoid

- better as mean of activation is zero, which makes training easier

\hookrightarrow ONE EXCEPTION is the output activation

\rightarrow if output is {0, 1}, sigmoid might be better.

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

NOTE your NN can have different activations at different layers. $\Rightarrow g^{[i]}$: activation at layer i

PROBLEM if z is either very large or very small, gradient of tanh or sigmoid function is very small, which slows down linear regression.



Rectified Linear Unit (ReLU)

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

- problem is derivative is 0 on the left

- advantage: slope is large so train very fast

Leaky ReLU \Rightarrow



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

General Rules

- if output is {0, 1} value, sigmoid is natural choice for output

- ReLU is very popular for hidden units

Why Do You Need Activation Functions?

$g(z) = z$: "linear activation function"

if we use a linear activation function:

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

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

$$= \underbrace{(W^{[2]}W^{[1]})}_{W'} x + \underbrace{(W^{[2]}b^{[1]} + b^{[2]})}_{b'}$$

$$= W'x + b' \rightarrow \text{this is just a } \underline{\text{linear function}}$$

the NN is just outputting a linear function, then

⇒ no matter # of layers we will still calculate a linear function

⇒ might as well have 1 layer. ⇒ USELESS

THE ONLY TIME we may want to use a linear activation function

⇒ if we do ML on a **REGRESSION** problem

Derivatives of Activation Functions

Sigmoid:

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

$$= a \cdot (1-a)$$

if $z=10$ $g(z) \approx 0$

if $z=-10$ $g(z) \approx 0$

(as $g(z)=a$)

tanh

$$g'(z) = 1 - (\tanh(z))^2 = 1 - a^2$$

if we calculated a already
super easy to find $g'(z)$

ReLU

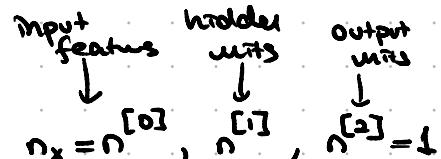
$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undef} & \text{if } z = 0 \end{cases}$$

Leaky ReLU

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

Gradient Descent for Neural Networks

parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$
 $(n^{[1]}, n^{[0]}) \quad (n^{[1]}, 1) \quad (n^{[2]}, n^{[1]}) \quad (n^{[2]}, 1)$



cost function: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{t=1}^m \mathcal{L}(\hat{y}, y)$
 (binary classification)

Gradient Descent

REPEAT {

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

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

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

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

- - - }

FORWARD PROPAGATION

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

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

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

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

BACK PROPAGATION

$$\frac{\partial z^{[2]}}{\partial z^{[1]}} = A^{[2]} - y$$

$$\frac{\partial w^{[2]}}{\partial z^{[1]}} = \frac{1}{m} \frac{\partial z^{[2]}}{\partial z^{[1]}} A^{[1]T}$$

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

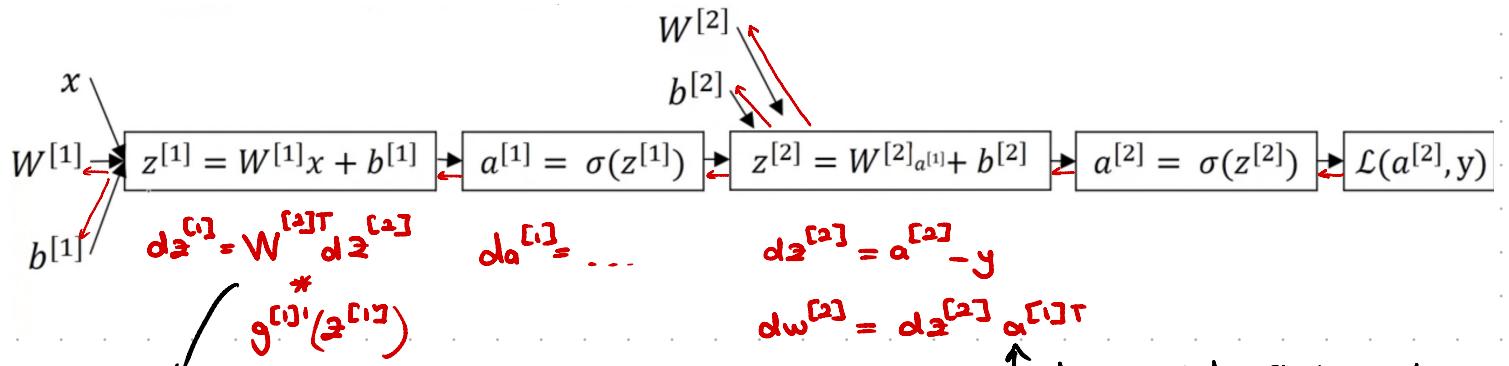
$$\frac{\partial z^{[1]}}{\partial z^{[2]}} = \underbrace{w^{[1]T}}_{(n^{[1]}, m)} \underbrace{\frac{\partial z^{[2]}}{\partial z^{[1]}}}_{(n^{[1]}, m)} * \underbrace{g^{[1]'}(z^{[1]})}_{(n^{[1]}, m)}$$

↑ avoids Rank 1
↑ element-wise product

$$\frac{\partial w^{[1]}}{\partial z^{[2]}} = \frac{1}{m} \frac{\partial z^{[1]}}{\partial z^{[2]}} X^T$$

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

BACKPROPAGATION CALCULATIONS



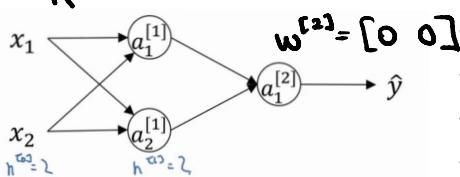
where does this come from?

$$\begin{aligned}
 dz^{[1]} &= \frac{\partial L}{\partial z^{[1]}} \\
 &= \frac{\partial L}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial a^{[1]}} \cdot \frac{\partial a^{[1]}}{\partial z^{[1]}} \\
 &= \frac{\partial L}{\partial a^{[2]}} \cdot \frac{\partial A^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial A^{[1]}} \cdot \frac{\partial A^{[1]}}{\partial z^{[1]}} \\
 &= \frac{dz^{[2]}}{(n^{[2]}, 1)} \cdot \frac{W^{[2]}}{(n^{[2]}, n^{[1]})} \cdot g^{[1]'}(z^{[1]}) \\
 &= W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]}) \\
 &\quad (n^{[2]}, 1)
 \end{aligned}$$

↑ transpose this to column vector
 why? remember that the W is a matrix where each row is $W_j^{[2]T}$ (i.e. a row vector)
 ↑ j th hidden unit.

Random Initiation

what happens if we initialize weights to zero?



$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad (\text{OKAY!})$$

$$a_1^{[1]} = a_2^{[1]} \quad dz_1^{[1]} = dz_2^{[1]}$$

Both hidden units are symmetric

Both compute exactly the same function after gradient descent

If we initialize W 's to zeros, the hidden units will be and will remain to be symmetric

USELESS to have > 1 hidden units

Therefore,

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

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

← Symmetry isn't a problem
for b

we want small values

why do we want small weights?

- if too large, z is too big or too small
- we end up in flat parts of tanh or sigmoid
- gradient is close to zero
- gradient descent is slow