

Neural Network Overview

$$\begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \rightarrow \textcircled{O} \rightarrow \widehat{y}_j = a \quad \dots$$

Input $x, w, b \rightarrow$ compute \hat{y} , using
 $x \quad \downarrow \quad z = w^T x + b \quad \begin{matrix} dz \\ da \end{matrix} \quad \text{compute } a$
 $w \quad \cancel{x} \quad \checkmark \quad a = b(x) \quad \text{Second.}$
 $b \quad \cancel{x} \quad \text{First calculation.} \quad \downarrow \uparrow \quad \text{finally}$
 In logistic regression,
 do backward calculate in order to compute
 L(a, y)
 ing together up all of derivatives.
 we can
 get lost
 function.
 (L)

Can form neural networks by stacking square brackets of sigmoid units.

[1] refer to quantities associated with this stack of node layer.

$x_1 \rightarrow 0$ $x_2 \rightarrow 0$ vertex to quantities untagged
 \downarrow \downarrow another layer.
 $[2] \rightarrow y = a^{[2]}$

~~if~~ \rightarrow 0 this node will correspond
- neural networks - another z and another

this stack of nodes will calculate a_1 & z like that ✓ what to be confused

bracket 1 and 2
of this different
layers. (layer 1,
layer 2)

individual training example.

(subsequent training example.)

$(x^{(i)}): i\text{th training example.}$

Right to left backward calculation.

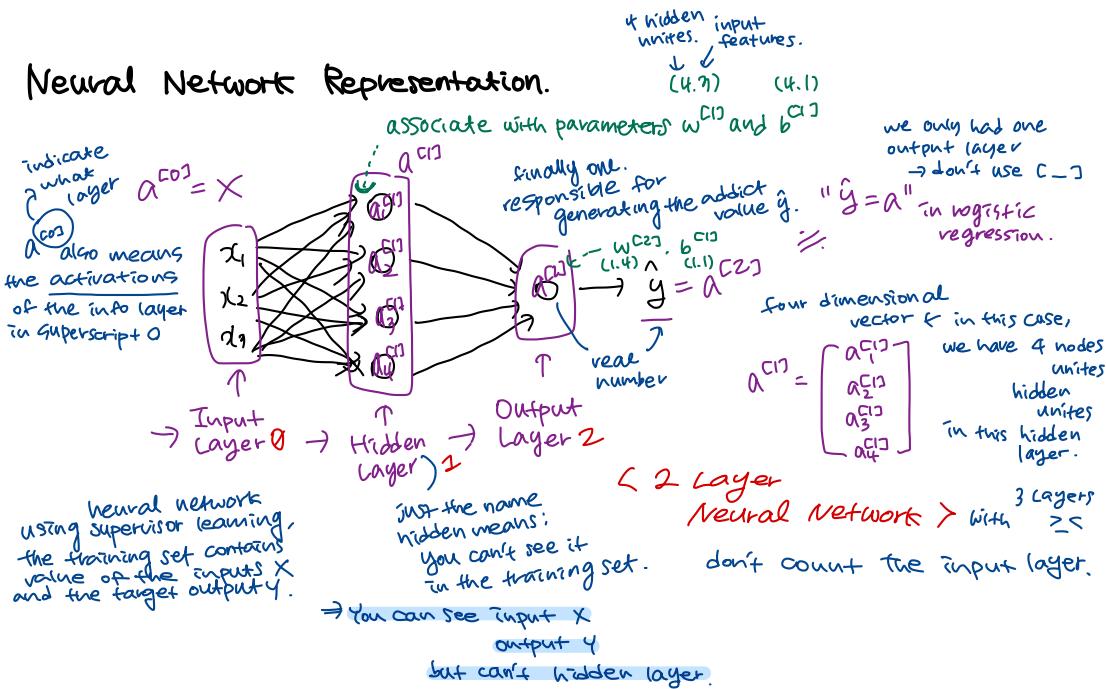
$$\begin{cases} z^{(1)} = w^{(1)} x + b^{(1)} \\ a^{(1)} = \sigma(z^{(1)}) \end{cases}$$

$$\begin{cases} z^{(2)} = w^{(2)} z^{(1)} + b^{(2)} \\ a^{(2)} = \sigma(z^{(2)}) \end{cases}$$

$$\begin{cases} z^{(3)} = w^{(3)} z^{(2)} + b^{(3)} \\ a^{(3)} = \sigma(z^{(3)}) \end{cases}$$

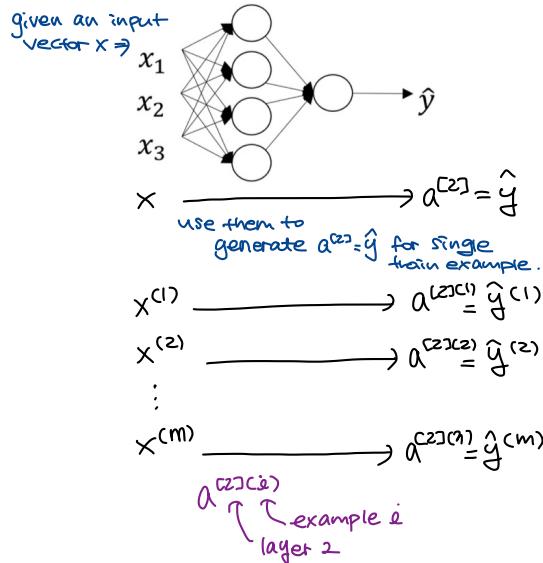
In logistic regression, we have z followed by a calculation in the neural network, we just do it multiple times and finally count a loss.

In logistic regression,
we have π followed by a calculation
in the neural network,
we just do it multiple times
and finally count a loss.



Vectorizing Across Multiple Examples.

How to vectorize multiple examples? Outcome would be quite similar with Logistic Regression.
Stacking up training examples in the different columns of matrices.



for $i = 1$ to m :

$$\begin{aligned} z^{(1)(i)} &= W^{(1)} x^{(i)} + b^{(1)} \\ a^{(1)(i)} &= f(z^{(1)(i)}) \\ z^{(2)(i)} &= W^{(2)} a^{(1)(i)} + b^{(2)} \\ a^{(2)(i)} &= f(z^{(2)(i)}) \end{aligned}$$

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

\uparrow $(n \times m)$

We went from lowercase vector x s to this capital case X matrix by stacking up the lowercase x 's in different columns.

$$\left\{ \begin{array}{l} z^{[1]} = W^{[1]} x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{array} \right. \quad \begin{array}{l} \text{for these 4 eq on top and adding the superscript round bracket is to all the variables that depend on the training example.} \\ \text{for } i = 1 \text{ to } m, \\ z^{[1](i)} = W^{[1]} x^{(i)} + b^{[1]} \\ a^{[1](i)} = f(z^{[1](i)}) \\ z^{[2](i)} = W^{[2]} a^{[1](i)} + b^{[2]} \\ a^{[2](i)} = f(z^{[2](i)}) \end{array}$$

vectorize to remove this for loop..

have to compute

$$\begin{aligned} z^{[1]} &= W^{[1]} x + b^{[1]} \\ A^{[1]} &= f(z^{[1]}) \\ z^{[2]} &= W^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} &= f(z^{[2]}) \end{aligned}$$

hidden units \downarrow \rightarrow training examples

lowercase \downarrow Capital \rightarrow

horizontally we're going to index across training examples - corresponds to different training examples.

activation of the first hidden unit on the first training ex. activation in the second hidden unit on the first training ex... first training ex.

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

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

\Rightarrow similarly $z^{[2]}$ & $A^{[2]}$ vertically, this vertical index corresponds to different notes in neural network.

Explanation for Vectorized Implementation.

<Justification for vectorized implementation>

$$\begin{aligned} z^{1} &= W^{[1]} x^{(1)} + b^{[1]} \\ z^{[1](2)} &= W^{[1]} x^{(2)} + b^{[1]} \\ z^{[1](m)} &= W^{[1]} x^{(m)} + b^{[1]} \end{aligned} \quad \begin{array}{l} \text{(for simplification)} \\ \text{w would be another column vector } (z^{[1](2)}) \end{array}$$

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

$w^{[1]} x^{(i)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$ Column vector.

$w^{[1]} x^{(i)} = z^{[1](i)}$ more stack when training sample is more.

$w^{[1]} x^{(i)} = z^{[1](i)} = z^{1} + b^{[1]}$ featuring examples - you have more examples and dug in more columns

$w^{[1]} x^{(i)} = z^{[1](i)} = z^{1} + b^{[1]} + b^{[2]}$ with python broadcasting, add $b^{[i]}$ individually.

Stack them up in different columns \rightarrow corresponding result is that

you end up with the \otimes also stack in columns.

Can convince yourself if you want \otimes with python broadcast if you add back in these values of B the values are still correct.

$z^{[1]} = W^{[1]} x + b^{[1]}$: correct vectorization in the first step of the 4 steps in the previous slide.

input columns
 \rightarrow output also stacked up in columns.

한 번에 하나의 훈련 샘플에 대해
정방행 전파를 한 다음 끝까지
모든 실행해야 함

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

equal to $A^{[1]}$
input feature vector
 $X = A^{[0]}$
 $x^{(i)} = a^{[0](i)}$

help vectorize
 $W^{[1]} A^{[0]} + b^{[1]}$
symmetric \Rightarrow doing the same thing / same computation over and over

Activation functions

What activation functions to use in hidden layers and output unit.

Sigmoid - Sometimes, other function is better.

- indicate $g^{(1)} \neq g^{(2)}$

Activation functions

$\text{almost always work}$

better than sigmoid function cause activation

$\hat{y} \pm 1 \rightarrow \text{mean}$

Sigmoid $a = \tanh \frac{z}{\sqrt{2}}$

$y \in \{0, 1\} = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

Given x: $0 \leq y \leq 1$

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

use Sigmoid function

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

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

\uparrow
called (2) $\rightarrow a^{[2]} =$
"activation function"

Rectified Linear Unit.

could be a nonlinear function Andrew Ng
that may not be sigmoid function.

You might center data and have your data has 0 mean using a tanh instead of a sigmoid function

- Centering your data \rightarrow the mean of the data is close to the zero
 \rightarrow actually makes learning for the next layer a little bit easier.

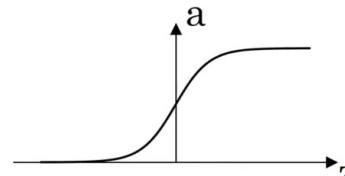
- one exception we use Sigmoid function: binary classification.
 $\rightarrow \tanh$ for hidden layer, sigmoid for output layer.

& ReLU for all the Rectified Linear Unit.

→ Not sure what to use in hidden unit, Just use
ReLU "ReLU function!"
work better, but not uses many practice.

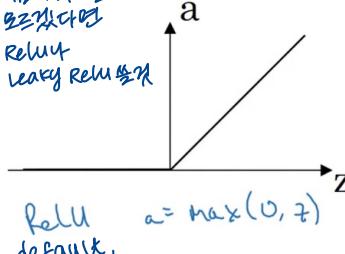
neural networks would be much faster when we use ReLU than using sigmoid or tanh. (most part are not 0) practice.

Pros and cons of activation functions



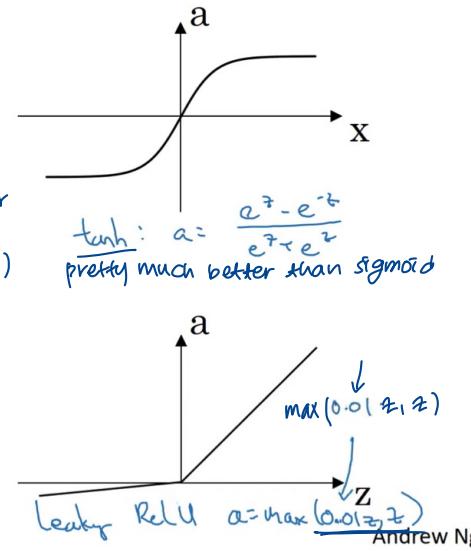
$$\text{sigmoid: } a = \frac{1}{1 + e^{-z}}$$

never use except for
"Output Layer"
(binary classification)



$$P_{\text{fail}} = \alpha = \max(0, \gamma)$$

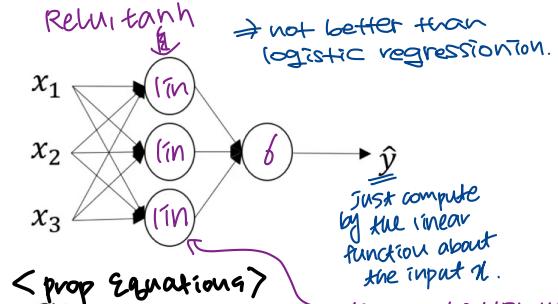
when you don't know which to use → try it all and see!!



함수의 기울기가 0이
가까워지면 학습이 느려짐
ReLU는 기울기가 0에 가까워지는
것을 막으므로 다른 활성화 함수
이나 바꾸어.

Why Non-Linear Activation functions

Activation functions



$$\begin{aligned}a^{[1]} &= z^{[1]} = w^{[1]} x + b^{[1]} \\a^{[2]} &= z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \\a^{[3]} &= w^{[3]} (w^{[1]} x + b^{[1]}) + b^{[3]} \\&= (w^{[3]} w^{[1]}) x + (w^{[3]} b^{[1]} + b^{[3]}) \\&= w^{[3]} x + b^{[3]}\end{aligned}$$

< prop equations >

Given x :

$$\begin{aligned}z^{[1]} &= W^{[1]} x + b^{[1]} \\a^{[1]} &= \cancel{z^{[1]}} \quad g(z) = z \\z^{[2]} &= W^{[2]} a^{[1]} + b^{[2]} \quad \text{"Linear Equation functions"} \\a^{[2]} &= \cancel{z^{[2]}}\end{aligned}$$

$$g(z) = z$$

"Linear Equation functions"
≡ identity activation functions.

: composition of 2 linear functions → itself a linear function.
⇒ not using non-linearly function, you can't compute more even as you go deeper.

Andrew Ng

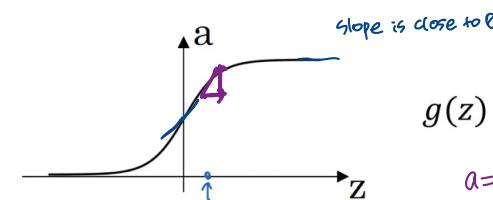
if you use a linear activation function or alternatively don't have an activation function, then no matter how many layers your neural network has, all those doing is just computing a linear activation function. ≡ not have any hidden layer

$g(z) = z$ is useful when you do machine learning regression program.
only used at out layer.

It might be okay to use linear activation when y ∈ ℝ and g ∈ ℝ
But the hidden unit, don't use linear function and use Relu, tanh, leaky Relu or other non-linear function.

Derivatives of Activation functions.

Sigmoid activation function



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

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

$$\text{i)} z=0, g(z) \approx 1$$

$$\frac{d}{dz} g(z) \approx 1/(1-1) \approx 0$$

$$\text{ii)} z=-10, g(z) \approx 0$$

$$\frac{d}{dz} g(z) \approx 0 \cdot (1-0) \approx 0$$

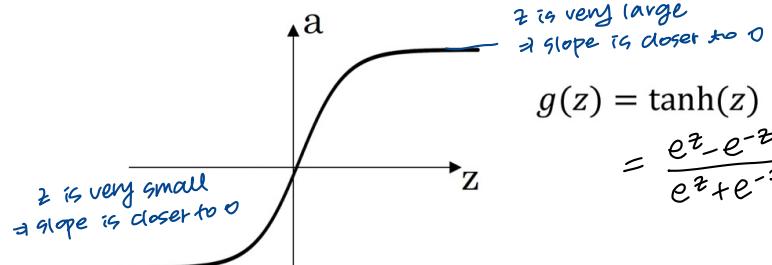
$$\text{iii)} z=0, g(z) = \frac{1}{2}$$

$$\frac{d}{dz} g(z) = \frac{1}{2}(1-\frac{1}{2}) = \frac{1}{4}$$

advantage:
if you already know a ,
you can know $g'(z)$
(slope of $g(z)$) easily.

Andrew Ng

Tanh activation function



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

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

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

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

→ quickly calculate

$$\text{i)} z=0, \tanh(z) \approx 1$$

$$g'(z) \approx 0$$

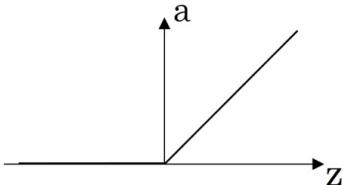
$$\text{ii)} z=-10, \tanh(z) \approx -1$$

$$g'(z) \approx -1$$

$$\text{iii)} z=0, \tanh(z) = 0$$

$$g'(z) = 1$$

ReLU and Leaky ReLU



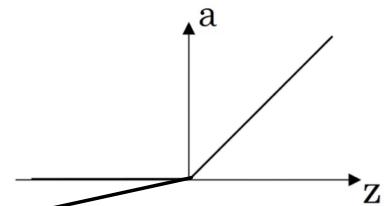
ReLU

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

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

~~unbiased if $z=0$~~

doesn't matter to think $z=0$ \dots
this value to 0 or 1. (either.. of my mind)



Leaky ReLU

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

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

no matter to set either 0 or 1.

Andrew Ng

Gradient Descent for Neural Networks

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

$n_x = n^{[0]}$, input
 $n^{[1]}$, hidden
 $n^{[2]} = 1$, output
example.

$$\text{Cost function: } J(w^{[0]}, b^{[0]}, w^{[1]}, b^{[1]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$$

Gradient descent:

initialize the parameters randomly
rather than to all 0 $\stackrel{?}{=} \text{not to 0!!}$

repeat

repeat this until parameters converges.

Compute Predicts $(\hat{y}^{(i)}, i=1 \dots, m)$
 $\frac{\partial J}{\partial w^{[1]}} = \frac{\partial J}{\partial w^{[1]}} = \frac{\partial J}{\partial b^{[1]}} = \dots$
same with $w^{[2]}, b^{[2]}$

 $w^{[1]} = w^{[1]} - \alpha \cdot \frac{\partial J}{\partial w^{[1]}}$
 $b^{[1]} = b^{[1]} - \alpha \cdot \frac{\partial J}{\partial b^{[1]}}$
 \vdots same with $w^{[2]}, b^{[2]}$

< Formulas for Computing Derivatives >

Forward Propagation:

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

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

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

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

binary classification in here..

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

Back Propagation:

$$dz^{[2]} = A^{[2]} - Y \quad \text{true value}$$

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

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

Sum across one dimension
of a matrix
 \rightarrow this case, Sum horizontally

$(n,)$
Time \rightarrow $(n, 1)$ $\forall i$ function
 $dz^{[2]} = W^{[2]T} dz^{[2]} \otimes g^{[2]'}(z^{[2]})$
 $(n^{[1]}, m)$ element wise product

$$dw^{[1]} = \frac{1}{m} dz^{[2]} X^T \quad (n^{[1]}, m)$$

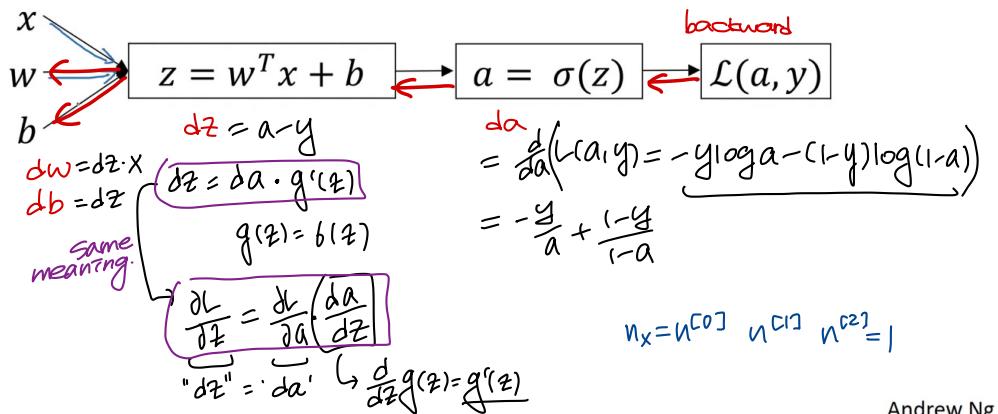
$$db^{[1]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, (n^{[1]}, 1))$$

keepdims = True

Backpropagation Intuition

Computing gradients

Logistic regression



Andrew Ng

Summary of gradient descent >

$$\begin{aligned} dz^{[2]} &= a^{[2]} - y \\ dw^{[2]} &= dz^{[2]} a^{[1]T} \\ db^{[2]} &= dz^{[2]} \end{aligned}$$

$$\begin{aligned} dz^{[1]} &= W^{[2]T} dz^{[2]} * g'(z^{[1]}) \\ dw^{[1]} &= dz^{[1]} X^T \\ db^{[1]} &= dz^{[1]} \end{aligned}$$

Vectorized Implementation -

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

capital.

$$\begin{aligned} \vec{z}^{[1]} &= W^{[1]} X + b^{[1]} \\ A^{[1]} &= g^{[1]}(\vec{z}^{[1]}) \end{aligned}$$

Forward pass:

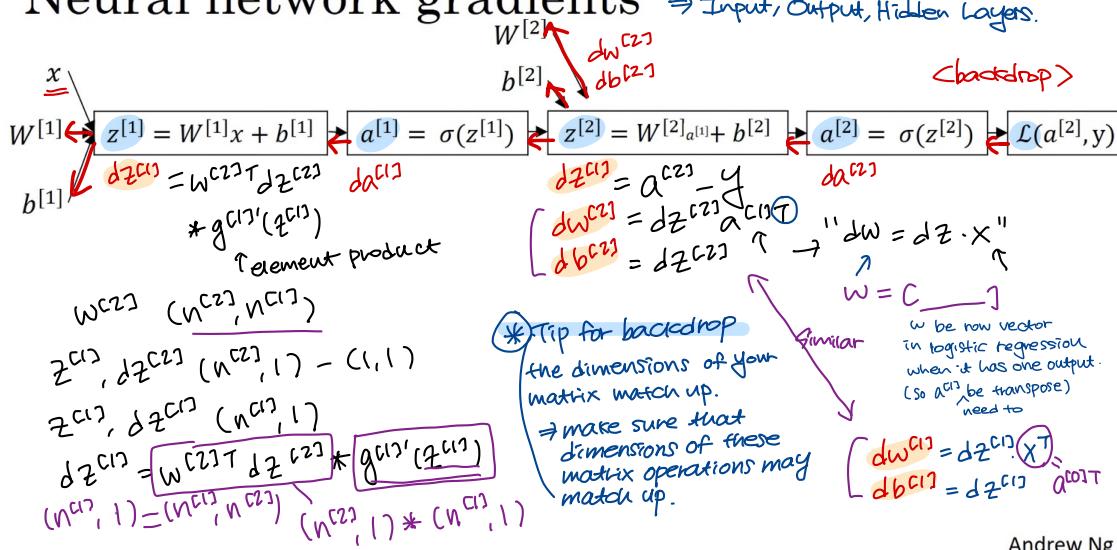
$$\begin{aligned} \vec{z}^{[2]} &= A^{[2]} - y \\ \vec{w}^{[2]} &= \frac{1}{m} \vec{d}z^{[2]} A^{[1]T} Y \\ \vec{b}^{[2]} &= \frac{1}{m} \text{np. sum}(\vec{d}z^{[2]}, \text{axis}=1, \text{keepdims=True}) \end{aligned}$$

Backward pass:

$$\begin{aligned} \vec{d}z^{[2]} &= W^{[2]T} \vec{d}z^{[2]} * g^{[2]}(z^{[1]}) \\ dw^{[2]} &= \frac{1}{m} \vec{d}z^{[2]} X^T \\ db^{[2]} &= \frac{1}{m} \text{np. sum}(\vec{d}z^{[2]}, \text{axis}=1, \text{keepdims=True}) \end{aligned}$$

elemental product.

Neural network gradients



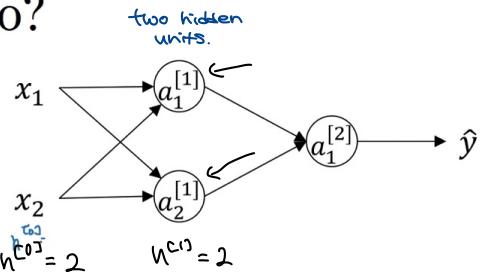
Andrew Ng

Random Initialization

It's important to initialize number randomly.

In Logistic Regression, it is okay to initialize θ . But for Neural Networks, To initialize the numbers to 0 and use gradient descent
 \Rightarrow Don't work.

What happens if you initialize weights to zero?



if) $w^{(1)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$ $b^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$
 \uparrow \uparrow
also $w^{(1)} = 0 \dots$ okay

can make problem.

$$a_1^{(1)} = a_2^{(1)} \quad dZ_1^{(1)} = dZ_2^{(1)}$$

$$w^{(2)} = [0 \ 0] \quad \therefore \text{hidden unit } a_1^{(1)} \text{ and } a_2^{(1)} \text{ be completely same.}$$

$$dW = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \quad w^{(1)} = w^{(1)} - dW$$

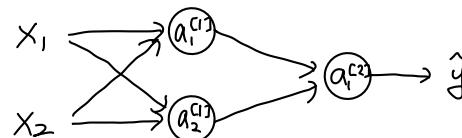
$$w^{(1)} = \begin{bmatrix} \underline{\hspace{2cm}} \\ \underline{\hspace{2cm}} \end{bmatrix} \quad \text{first row equal to second row.}$$

Symmetric
 \Rightarrow compute exactly same function.

$$\begin{array}{l} x_1 \\ x_2 \\ x_3 \\ x_3 \end{array} \quad \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \text{"Same"}$$

if you initialize to 2,
all these hidden units are
symmetric
 \Rightarrow how long you use gradient
descent, all units would always
compute same function.
(not helpful)

<Random Initialization>



$\Rightarrow w^{(1)} = \text{np.random.rand}(2, 2) * 0.01$
 $b^{(1)} = \text{np.zeros}(2, 1)$
 $w^{(2)} = \text{np.random.rand}(1, 2) * 0.01$
 $b^{(2)} = 0$

any very
small number

: Symmetry breaking problem,
it's okay to initialize b to 0
 $Z^{(1)} = w^{(1)}x + b^{(1)}$
 $a^{(1)} = g^{(1)}(Z^{(1)})$
if the weight is too large,
activation value \rightarrow slow gradient.
slope is very small,
gradient descent is
very slow
learning is slow.