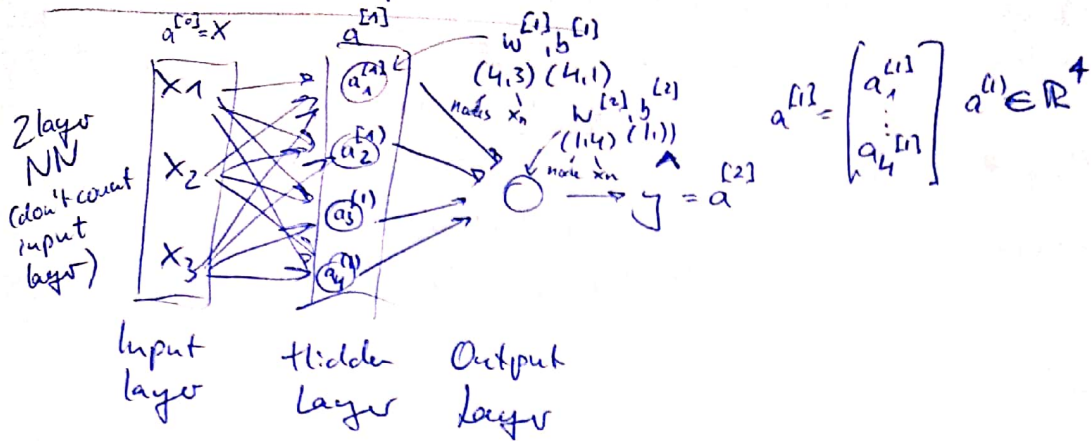
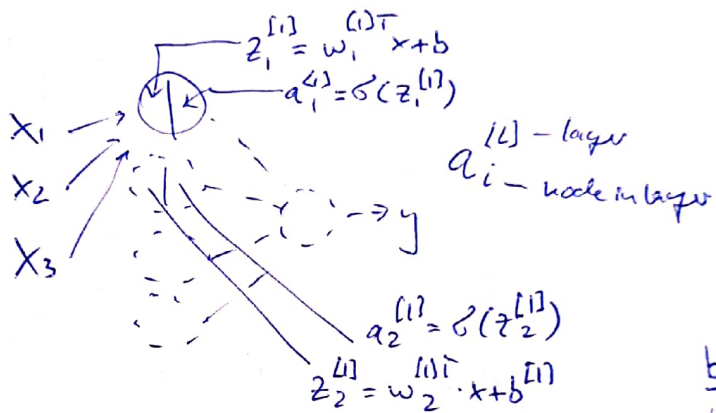
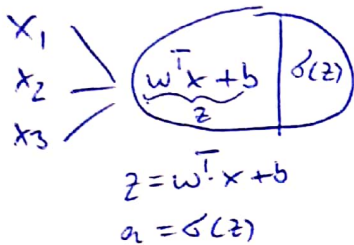


Neural Network Representation



Computing NN's output



calculate $\begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix}$ in for loop very inefficient!

$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{bmatrix} = \sigma(z^{(1)})$

Use Vectorization!

$\underline{W}^{(1)} = \begin{bmatrix} -w_1^{(1)T} \\ -w_2^{(1)T} \\ -w_3^{(1)T} \\ -w_4^{(1)T} \end{bmatrix} \in \mathbb{R}^4$

shape $(4, 3)$

$\underline{b}^{(1)} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix} \in \mathbb{R}^4$

$\underline{z}^{(1)} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix}$

$\underline{a}^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{bmatrix}$

Given input x :

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

$(4, 1) \quad (4, 3) \quad (3, 1) \quad (4, 1)$

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

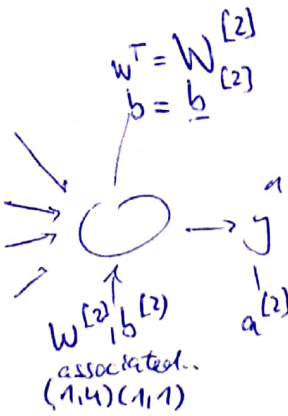
$(4, 1) \quad (4, 1)$

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

$(1, 1) \quad (1, 4) \quad (4, 1) \quad (1, 1)$

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

$(1, 1) \quad (1, 1)$



4 lines of code to compute output

Vectorizing across multiple examples

$$\begin{aligned}
 x &\longrightarrow a^{[2]} = y \\
 x^{(1)} &\longrightarrow a^{[2](1)} = y^{(1)} \\
 x^{(2)} &\longrightarrow a^{2} = y^{(2)} \\
 &\vdots \\
 x^{(m)} &\longrightarrow a^{[2](m)} = y^{(m)}
 \end{aligned}$$

$a^{[2](i)}$
 a | i -th example
 layer 2

to compute for multiple examples:

for $i=1$ to m ($\sum_{i=1}^m$)

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

want to vectorize:

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

features

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

Vectorized implementation

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

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

corresponds to node in NN (hidden units)

Justification for vectorized implementation:

$$W^{[1]} = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix} \quad W^{[1]} x^{(1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad W^{[1]} x^{(2)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad W^{[1]} x^{(3)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$W^{[1]} \cdot \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots \\ | & | & | & | \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} = \begin{bmatrix} | & | & | & | \\ z^{1} & z^{[1](2)} & z^{[1](3)} & \dots \\ | & | & | & | \end{bmatrix} = z^{[1]}$$

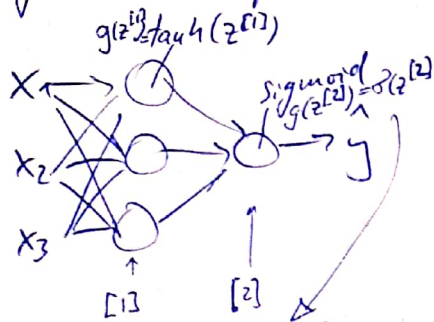
X

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

does correctly with Python broadcasting

Activation functions

When building NN architecture you need to choose activation functions for hidden ~~layer~~ units and output unit



so far:

Given x :

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

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

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

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

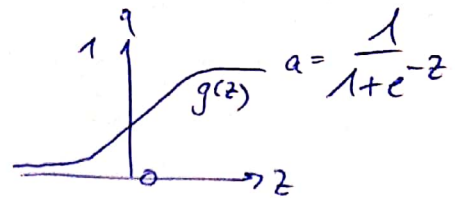
Almost always use tanh

exception: binary classification
 ↳ output unit should give $\hat{y} \in \{0, 1\}$ rather than $-1 < \hat{y} < 1$ Sigmoid!

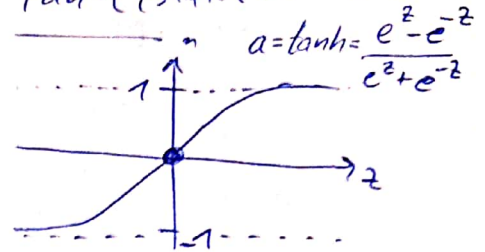
One downside for sigmoid and tanh:

• If z is either very large or very small the gradient is really small! Slope of function is close to zero and slows down Gradient Descent

Sigmoid



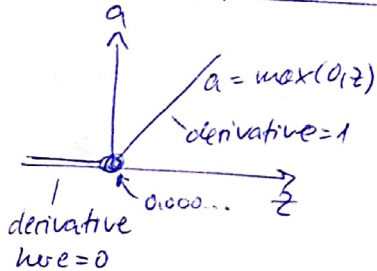
Tanh (shifted version of Sigmoid)



(almost always work better than Sigmoid)
 values between -1 and 1
 the mean of the activation is closer to 0.
 When learning an algorithm you want to center the data

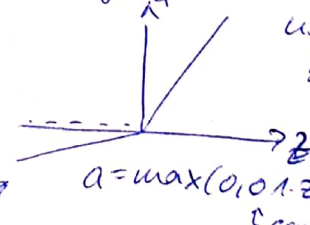
Centering
 ↳ mean closer to 0, rather than 0.5 (Sigmoid), makes learning for next layer easier

ReLU - Rectified Linear Unit



↳ disadvantage that derivative is 0 when z is negative

Leaky ReLU



usually works better but not often used in practice

common advantage: for a large space of z the derivative is very different from 0
 ↳ NN will learn much faster

Pros and Cons of activation functions

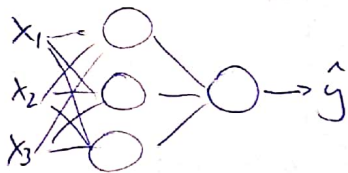
Sigmoid: never use this except for output layer when doing binary classification

tanh: superior to sigmoid

new default: ReLU

or Leaky ReLU

Why do we need non-linear activation functions?



↳ Neural Network would just compute a linear response of the input no matter how many hidden layers we have

A linear hidden layer is useless!

Given x :

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

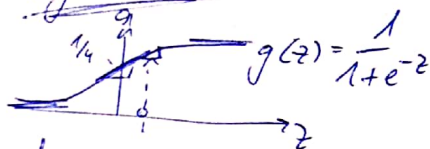
$a^{[1]} = g^{[1]}(z^{[1]}) \rightarrow z^{[1]} \rightarrow g(z) = z$ would be called
 "linear activation function"
 $z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$
 $a^{[2]} = g^{[2]}(z^{[2]}) \rightarrow z^{[2]}$
 or "identity function"

$$z^{[2]} = W_a^{[2]} [1] + b^{[2]} \quad (2)$$

$$a^{[2]} = g^{[2]}(z^{[2]}) \rightarrow z^{[3]}$$

Derivatives of activation functions

Sigmoid



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

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

$$g'(z) = g(z)(1 - g(z))$$

with $q = g(z)$

$$G'_{f(A)} = a(1-a)$$

Sanity Check!

when z is large

$$\rightarrow \frac{dg(z)}{dz} \approx 1(1-1) = 0 \checkmark$$

when z is small

$$\frac{dg(z)}{dz} \approx 0 \cdot (1-0) \approx 0 \checkmark$$

when $z=0$

$$\frac{dg(z)}{dz} = \frac{1}{2} \left(1 - \frac{1}{z} \right) = \frac{1}{4}$$

Safety Check

when z is large

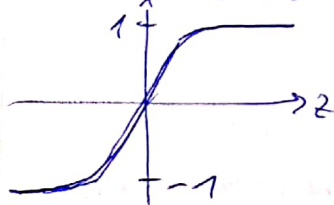
$$\tanh(z) \approx 1 \vee g'(z) \approx 0 \vee$$

When 7 is over

$$\text{Len } h(z) \approx -1 \quad g^{(z)} \approx 0V$$

when $z=0$

$$\begin{aligned} \tanh(z) &= 0 \\ g'(z) &= 1 \end{aligned}$$

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


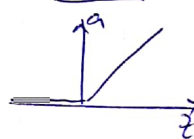
$$\frac{d g(z)}{dz} = 1 - (\tanh(z))^2 \quad \tanh(z) = 0$$

$$g'(z) = 1$$

$$a \equiv g(z)$$

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

Relly

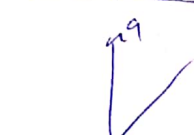


$$g(z) = \max(0, z)$$

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

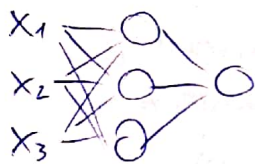
(undefined if $z=0$)

Leaky ReLU



$$g(z) = \max(0, 0.1 \cdot z, z)$$

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

Gradient Descent for Neural Networks

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$
 dimensions: $(n^{[1]}, n^{[2]})$ | $(n^{[2]}, n^{[3]})$
 $(n^{[1]}, n^{[1]})$ | $(n^{[2]}, n^{[2]})$

$n_x = n^{[0]}, n^{[1]}, n^{[2]}$
 / | /
 input hidden output
 features

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

to train parameters:

Gradient Descent

• initialize parameters randomly!

• compute predictions $(\hat{y}^{(1)}, \dots, \hat{y}^{(m)})$

• partial derivatives $\frac{\partial J}{\partial w^{[1]}}, \frac{\partial J}{\partial b^{[1]}}$

• GD update = $w^{[1]} = w^{[1]} - \alpha \frac{\partial J}{\partial w^{[1]}}$
 $b^{[1]} = b^{[1]} - \alpha \frac{\partial J}{\partial b^{[1]}}$

Forward propagation

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

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

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

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

can be
sigmoid here...

→ obtain loss $\mathcal{L}(A^{[2]}, y)$

Back propagation:

$$dz^{[2]} = A^{[2]} - Y \quad Y = [y^{(1)}, y^{(m)}]$$

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

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

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

element
wise
product
of 2
matrices

derivative

sum up
horizontally

prevents
python from
"rank-1-array"
can also
be done
with reshape!

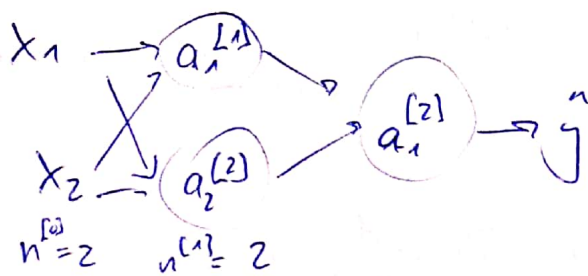
$$dw^{[1]} = \frac{1}{m} dz^{[1]} x^T$$

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

$(n^{[1]}, 1)$ ← ! avoid "rank-1-arrays!"

Random initialization

→ initialize parameters (w, b) randomly!



$$W^{(1)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

(2,2) ↑

problem!

OK!

→ $a_1^{(1)}$ will be equal to $a_2^{(1)}$!

$$a_1^{(1)} = a_2^{(1)} \text{ (circled)}, \text{ so } dz_1^{(1)} = dz_2^{(1)} \text{ (circled)}$$

→ hidden units in 1st layer are "symmetric" → compute exactly the same function!

after iterations still compute the same function

$$dW = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \quad \text{when perform weight update: } W^{(1)} = W^{(1)} - \alpha dW$$

$$\Rightarrow W^{(1)} = \begin{bmatrix} - & - & - \\ - & - & - \end{bmatrix} \quad \left. \begin{array}{l} \text{1st row} \\ \text{2nd row} \end{array} \right\} \text{ equal to}$$

Solution: initialize W randomly!

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

initialize to small values

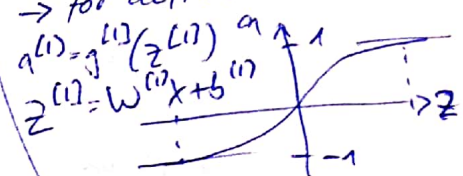
$$b^{(1)} = np.zeros((2,1))$$

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

$$b^{(2)} = 0$$

Why use 0.01 and not 1000?

→ for activations z will be large!



→ we will end up at flat parts of tanh or sigmoid and have low slope and therefore slow convergence for Gradient Descent