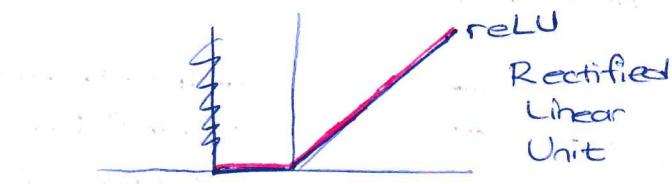
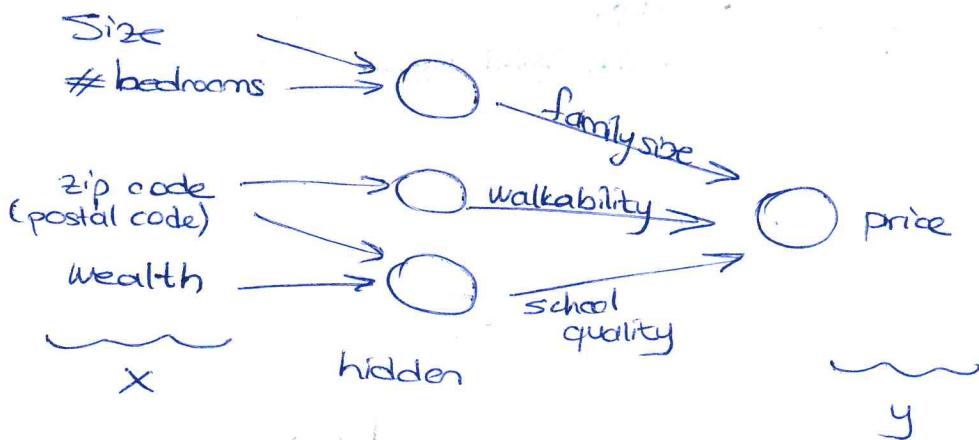


Neural Networks and Deep Learning

WEEK 1
1



Ex:



Supervised Learning:

<u>Input (x)</u>	<u>Output (y)</u>	<u>Application</u>	
Home features	Price	Real Estate	standard NN
Ad, user info	Click on ad? (0/1)	Online Advertising	
Image	Object (1, ..., 1000)	Photo tagging	CNN
Audio	Text transcript	Speech recognition	
English	Chinese	Machine Translation	RNN
Image, Radar Info	Position of other cars	Autonomous driving	
			custom/hybrid

CNN are often used for image data

RNNs are ~~not~~ very good for this type of one-dimensional sequence data that has maybe a temporal component.

Structured Data

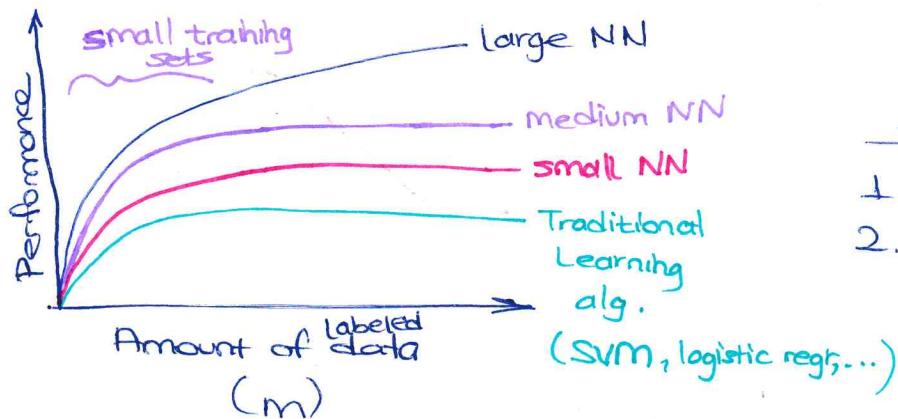
databases of data

Unstructured Data

audio, raw audio, image, text

Scale drives deep learning Progress

WEEK 1
2



2 options to get better results

1. big neural network
2. large data

Scale drives deep learning progress

- Data
- Computation
- Algorithms

$$\sigma(x)$$

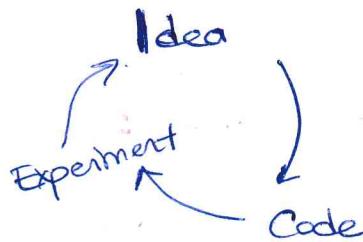
Sigmoid

$$\begin{cases} 0 & \text{if } x \leq 0 \\ 1 & \text{if } x > 0 \end{cases}$$

Value funct.

ReLU

gradient is very slow
so learning is
very slow



Outline of this course

Week 1.

- || 2.
- || 3.
- || 4.

Courses Resources

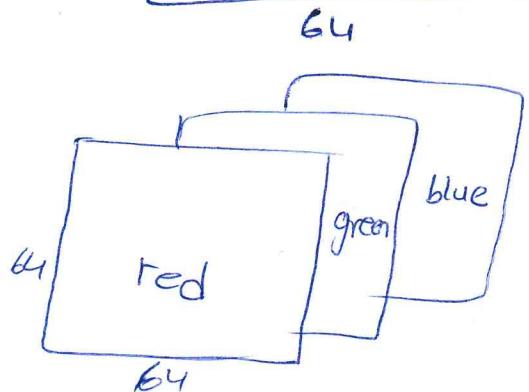
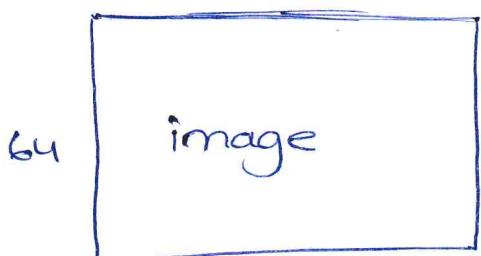
Discussion forum

Leaky ReLU

Binary Classification

WEEK 2

3



$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 234 \\ \vdots \\ 1 \end{bmatrix} \quad \begin{array}{l} \text{red pixels} \\ \text{green pixels} \\ \text{blue pixels} \end{array}$$

$$64 \times 64 \times 3 = 12288$$

representation of input dimension
 $n = n_x = 12288$

$$x \rightarrow y$$

$$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

m training examples : $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$$m = M_{\text{train}} \quad M_{\text{test}} = \# \text{ test examples}$$

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \end{bmatrix} \quad \begin{array}{c} \uparrow \\ n_x \end{array}$$

$\xleftarrow{m} \qquad \qquad \qquad \xrightarrow{m}$

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

$$Y \in \mathbb{R}^{1 \times m}$$

$$X \in \mathbb{R}^{n_x \times m}$$

$$X.\text{shape} = (n_x, m)$$

$$Y.\text{shape} = (1, m)$$

Logistic Regression

Given x

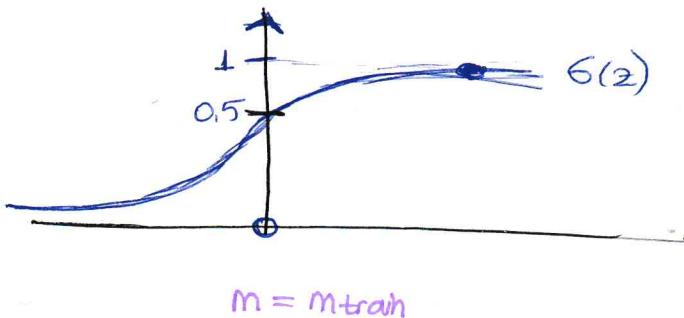
$$\hat{y} = P(y=1|x)$$

$$0 \leq \hat{y} \leq 1$$

$$x \in \mathbb{R}^{n_x}$$

parameters : $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

$$\text{output } \hat{y} = \sigma(w^T x + b)$$



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

$$\text{If } z \text{ large } \sigma(z) \approx \frac{1}{1+0} = 1$$

If z large ~~is~~ negative number

$$\sigma(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+\text{big number}} \approx 0$$

Logistic Regression Cost Function:

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \text{ where } \sigma(z) = \frac{1}{1+e^{-z}}$$

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$

Loss (error) function : $L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

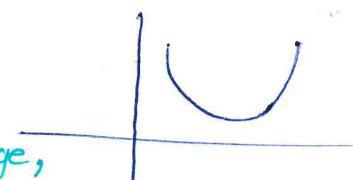
↓
(least square)

$x^{(i)}$ $y^{(i)}$ $z^{(i)}$
i-th example

many local opt.

gradient descent may not find the global optimum.

We need global optimum



→ we need this

want $\log 1-\hat{y}$ large
want \hat{y} small

want $\log \hat{y}$ large,
want \hat{y} large

*Finally, the loss function was defined to a single training example

Cost function which measures how well you're doing on an entire training set.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

the loss function for training example and overall cost function for the parameters

can be viewed as a very very small neural network of your algorithm.

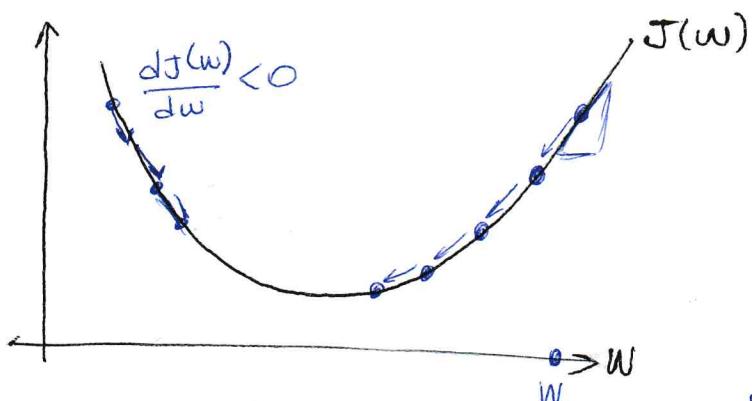
Gradient Descent

WEEK 2

(5)

Recap

Want to find w, b that minimize $J(w, b)$



Repeat {

$$w := w - \alpha \frac{dJ(w)}{dw}$$

}

$$w := w - \alpha dw$$

Learning rate
↓
 $\frac{dJ(w)}{dw}$
dw

$$J(w, b)$$

$$w := w - \alpha \cdot \frac{dJ(w, b)}{dw}$$

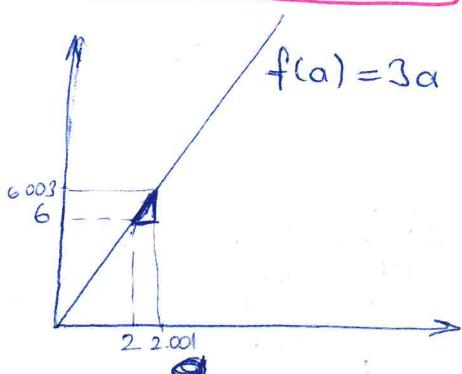
$$b := b - \alpha \cdot \frac{dJ(w, b)}{db}$$

∂ "partial derivative"

d ← J(~~one variable~~)

$$\begin{array}{c} \boxed{\frac{\partial J(w, b)}{\partial w}} \\ \downarrow \\ \boxed{\frac{\partial J(w, b)}{\partial b}} \end{array}$$

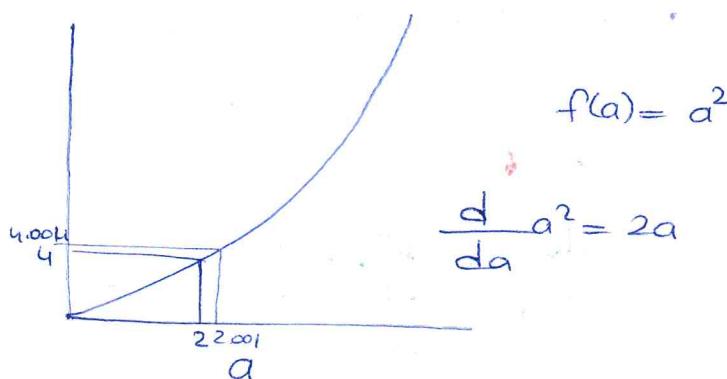
Derivatives



slope (derivative) of f(a) at $a=2$ is 3

slope at $a=5$ is also 3

$$\frac{df(a)}{da} = 3 = \frac{d}{da} f(a)$$



$$f(a) = a^2$$

$$\frac{d}{da} a^2 = 2a$$

$$a=2$$

$$f(a) = 4 \quad a=2,001 \quad f(a) \approx 4.004 \quad (4.004001)$$

slope (derivative) of $f(a)$ at $a=2$ is 4

$$\frac{d}{da} f(a) = 4 \text{ when } a=2$$

$$a=5 \quad f(a)=25$$

$$a=5.001 \quad f(a) \approx 25.010$$

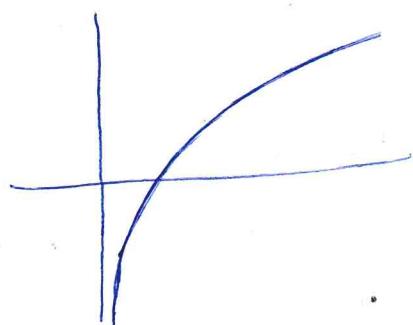
$$\frac{d}{da} f(a) = 10 \text{ when } a=5$$

$$\frac{d}{da} f(a) = \frac{d}{da} a^2 = 2a$$

(6)

$$f(a) = \log_e(a) \\ \ln(a)$$

$$\frac{d}{da} f(a) = \frac{1}{a}$$



$$a=2 \quad f(a) \approx 0.69315$$

$$a=2.001 \quad f(a) \approx 0.69365$$

0.01 0.0005

Computation Graph

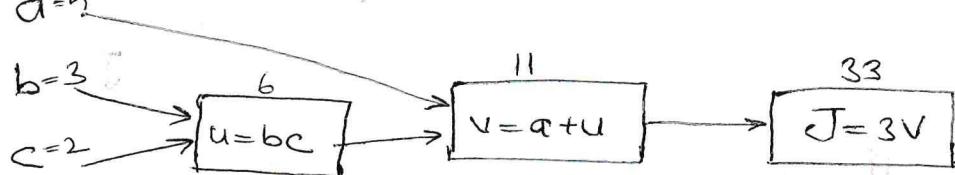
$$J(a, b, c) = 3(a + bc) \\ = 3 \cdot \underbrace{(a + bc)}_{v} = 33$$

$\underbrace{}_u \quad \underbrace{}_v$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



The computation graph organizes a computation with left-to-right computation.

$$\frac{dJ}{dv} = 3$$

$$\frac{dJ}{da} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{da} \\ 3 \cdot 1$$

$$\frac{dv}{da} = 1$$

$$a \rightarrow v \rightarrow J$$

$$\frac{dJ}{db} = \underbrace{\frac{dJ}{du}}_3 \cdot \underbrace{\frac{du}{db}}_{\frac{1}{2}} = 6$$

$$\frac{dJ}{du} = \frac{dJ}{dv} \cdot \frac{dv}{du} = 3 \cdot 1 = 3$$

$$\frac{dJ}{da} = \frac{dJ}{du} \cdot \frac{du}{da} = 3 \cdot 1 = 3$$

$$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db} = 3 \cdot \frac{1}{2} = 1.5$$

$$b = 3 \rightarrow 3.001$$

$$u = b \cdot c = 6 \rightarrow 6.006$$

$$J = 33.006$$

$\frac{d \text{Final Output Var}}{d \text{Var}}$

$\frac{dJ}{dV \text{Var}}$

"dvar" in python

$$\frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} = 9$$

$dVar$ represents the derivative of a final output variable with respect to various intermediate quantities.

Logistic Regression - Gradient Descent

WEEK 2

(7)

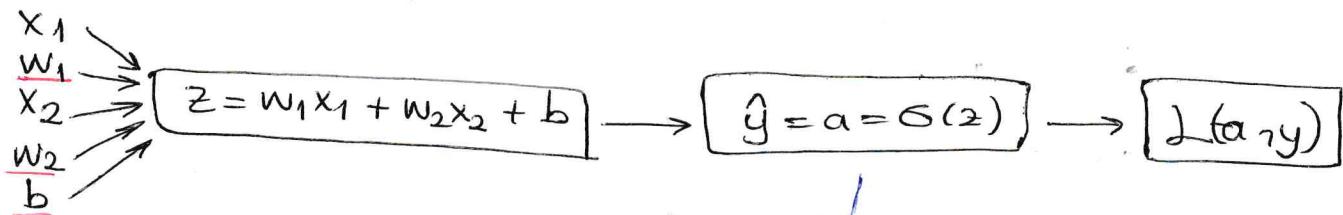
$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$a \rightarrow$ is output of logistic regression

$$L(a, y) = -(y \log(a) + (1-y) \log(1-a))$$

$y \rightarrow$ ground truth label



modify w_1, w_2, b in order to reduce Loss.

$$dz = \frac{dL}{dz} = \frac{dL(a, y)}{da}$$

$$= a - y$$

$$= \frac{dL}{da} \cdot \frac{da}{dz}$$

$$= \left(-\frac{y}{a} + \frac{(1-y)}{(1-a)} \right) \leftarrow a \cdot (1-a)$$

$$"da" = \frac{dL(a, y)}{da}$$

$$= -\frac{y}{a} + \frac{(1-y)}{(1-a)}$$

$$\frac{\partial L}{\partial w_1} = "dw_1" = x_1 \cdot dz \quad dw_2 = x_2 \cdot dz \quad db = dz$$

(in python)

$$\begin{cases} w_1 := w_1 - \alpha dw_1 \\ w_2 := w_2 - \alpha dw_2 \\ b := b - \alpha db \end{cases}$$

Gradient Descent on m examples

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)}) \quad (x^{(i)}, y^{(i)})$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$dw_1^{(i)}, dw_2^{(i)}, db^{(i)}$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial L(a^{(i)}, y^{(i)})}{\partial w_1}}_{dw_1^{(i)}} \quad \text{gradient}$$

$$- (x^{(i)}, y^{(i)})$$

$$\mathbf{J} = 0 ; \mathbf{dw}_1 = 0 ; \mathbf{dw}_2 = 0 ; \mathbf{db} = 0$$

For $i = 1$ to m

$$z^{(i)} = w^T x^{(i)} + b$$

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

$$\mathbf{J} += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$$

$$d_z^{(i)} = a^{(i)} - y^{(i)}$$

$$\left[\begin{array}{l} dw_1 += x_1^{(i)} d_z^{(i)} \\ dw_2 += x_2^{(i)} d_z^{(i)} \\ db += d_z^{(i)} \end{array} \right] \quad \uparrow n=2$$

may need another d_{w_3} for loop

$$\mathbf{J} /= m$$

$$dw_1 /= m ; dw_2 /= m ; db /= m$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

Q: Why is there only one dw variable?
(i.e. no i superscripts in the for loop)

A: The value of dw in the code is cumulative.

Vectorization

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

Non-vectorized

$$z = 0$$

for i in range($n-x$):

$$z += w[i] * x[i]$$

$$z += b$$

- it is going to be slow



Vectorized

$$z = \underbrace{\text{np.dot}(w, x)}_{w^T x} + b$$

$$w \in \mathbb{R}^{n_x}$$

$$x \in \mathbb{R}^{n_x}$$

CPU
GPU

} have parallelization instructions called SIMD instruction.
(Single Instruction Multiple Data)

More Vectorization Examples

WEEK 2.

9

Whenever possible \rightarrow avoid explicit for-loops.

$$u = Av$$

$$u_i = \sum_j A_{ij} v_j$$

$$u = np.zeros((n, 1))$$

for i ...

for j ...

$$u[i] += A[i][j] * v[j]$$

$$u = np.dot(A, v)$$

NOTE: $a+b$ \rightarrow is an element-wise computation

$np.dot \rightarrow$ matrix multiplication

Say you need to apply the exponential operation on every element of a matrix / vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

$$u = np.zeros((n, 1))$$

for i in range(n) :

$$u[i] = \text{math.exp}(v[i])$$

import numpy as np

$$u = np.exp(v)$$

$$np.log(v)$$

$$np.abs(v)$$

$$np.maximum(v, 0)$$

$$v**2$$

instead of $\boxed{dw_1=0 \ dw_2=0}$ we prefer $dw = np.zeros((n-k, 1))$

$$dw = x^{(i)} d_z^{(i)}$$

$$dw / = m$$

Vectorizing Logistic Regression

$$X = \left[\begin{array}{c|c|c|c} 1 & x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \hline 1 & | & | & \dots & | \end{array} \right] \in \mathbb{R}^{(n+1) \times m}$$

$$z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = w^T X + [b \ b \ \dots \ b] = [w^T x^{(1)} b \ w^T x^{(2)} b \ \dots \ w^T x^{(m)} b]$$

$$z = np.dot(w.T, X) + b \xrightarrow{(1,1)} \text{Broadcasting}$$

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(z)$$

Vectorizing Logistic Regression's Gradient Computation

$$d_2^{(1)} = a^{(1)} - y^{(1)} \quad d_2^{(2)} = a^{(2)} - y^{(2)}$$

$$d\hat{z} = [d_2^{(1)} \quad d_2^{(2)} \quad \dots \quad d_2^{(m)}]$$

$$A = [a^{(1)} \dots a^{(m)}] \quad Y = [y^{(1)} \dots y^{(m)}]$$

$$d\hat{z} = A - Y = [a^{(1)} - y^{(1)} \quad a^{(2)} - y^{(2)} \quad \dots]$$

$$db = \frac{1}{m} \sum_{i=1}^m d_2(i)$$

$$= \frac{1}{m} np.sum(d\hat{z})$$

$$dw = \frac{1}{m} \cdot X \cdot d\hat{z}^T$$

$$= \frac{1}{m} \left[\begin{array}{cccc} 1 & x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & 1 & \dots & 1 \end{array} \right] \left[\begin{array}{c} d_2^{(1)} \\ \vdots \\ d_2^{(m)} \end{array} \right]$$

$$\frac{1}{m} \left[x^{(1)} d_2^{(1)} + \dots + x^{(m)} d_2^{(m)} \right]$$

$n \times 1$

$$Z = w^T X + b$$

$$= np.dot(w.T, X) + b$$

→ gradient descent for logistic regression

$$A = g(Z)$$

→ for one iteration

$$d\hat{z} = A - Y$$

$$dw = \frac{1}{m} \cdot X \cdot d\hat{z}^T$$

$$db = \frac{1}{m} \cdot np.sum(d\hat{z})$$

$$W = W - \alpha dw$$

$$b_i = b_i - \alpha db$$

Broadcasting in Python

WEEK 2

(11)

cal = A.sum (axis=0)

→ vertically (0) , horizontally (axis=1)

percentage = 100 * A/cal.reshape(1,4)

↑
(3,4)

↳ it can be used to be sure.

General Principle :

(m,n) + (1,n) → (m,n)

matrix * (m,1) → (m,n)

(m,1) + R

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

$$[1 \ 2 \ 3] + 100 = [101 \ 102 \ 103]$$

Matlab or Octave bsxfun function

A note on python/numpy vectors

a = np.random.randn(5)

a.shape=(5,1)

"rank 1 array"

} don't use

a = np.random.randn(5,1)

→ a.shape = (5,1) column vector

a = np.random.randn(1,5)

→ a.shape = (1,5) row vector

assert(a.shape == (5,1))

→ assertion statement

is really inexpensive

checks the shape.
if is not ^{correct} consistent,
it stops.

! it also help to serve as documentation
for your code.

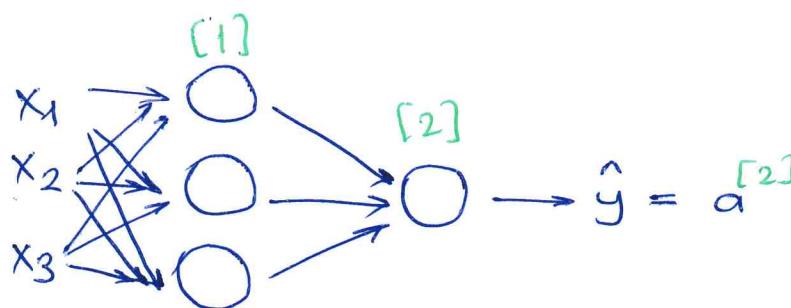
! don't be shy about calling the reshape
operation

Neural Networks Overview

WEEK 3

(1)

Elif Cenur
YILDIZ



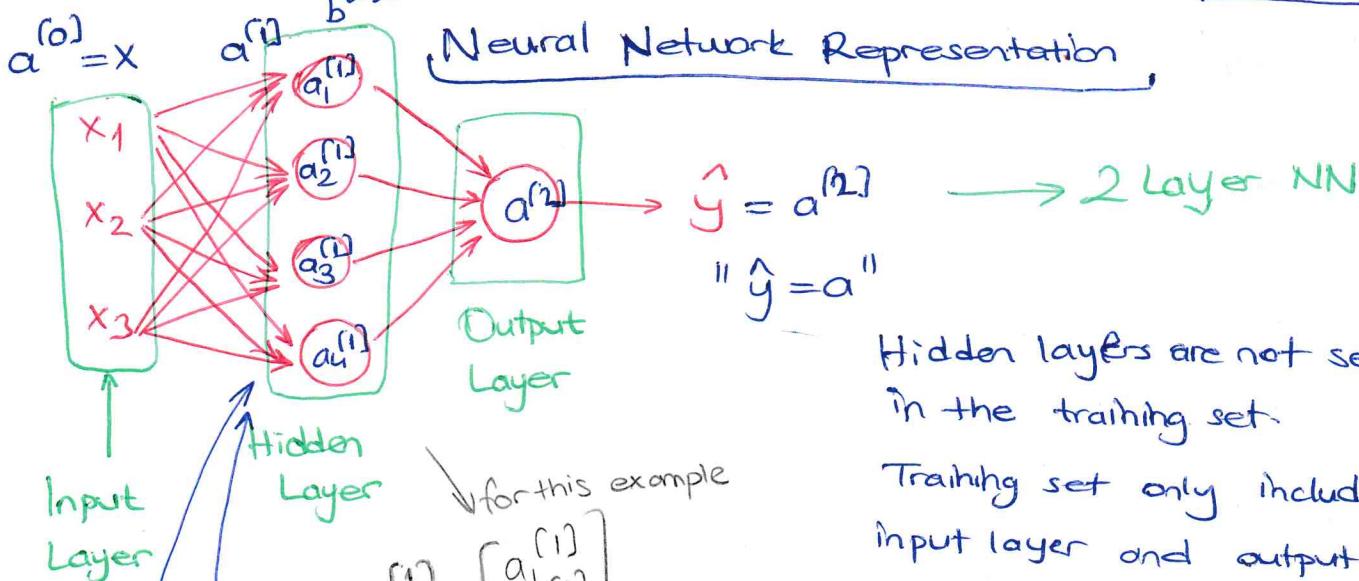
$$w^{[1]} \rightarrow z^{[1]} = w^{[1]}x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]})$$

Below this, another set of equations shows the calculation for the second layer:

$$w^{[2]} \rightarrow z^{[2]} = w^{[2]}a^{[1]} + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]})$$

Neural Network Representation

$$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow L(a^{[2]}, y)$$



Hidden layers are not seen in the training set.

Training set only includes the input layer and output layer.

In conventional usage,
~~hidden layers and output layers are~~ counted.

for this example

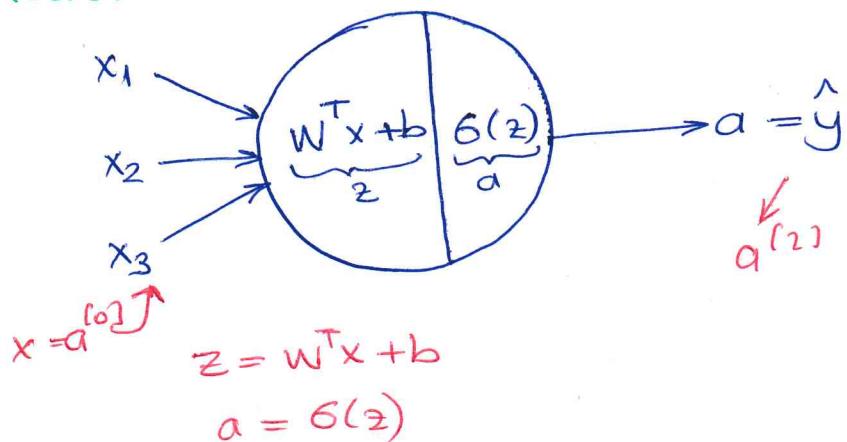
$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}$$

$W^{[1]}, b^{[1]}$
 $(4, 3)$
 ↓ hidden inputs
 hidden layer dimension

Computing a Neural Network's Output

(2)

Neural Network Representation:



$a_i^{[l]}$ → layer
 $a_i^{[l]}$ → node in layer

$$\begin{aligned} z_1^{(1)} &= w_1^{(1)T} x + b_1^{(1)}, & a_1^{(1)} &= g(z_1^{(1)}) \\ z_2^{(1)} &= w_2^{(1)T} x + b_2^{(1)}, & a_2^{(1)} &= g(z_2^{(1)}) \\ z_3^{(1)} &= w_3^{(1)T} x + b_3^{(1)}, & a_3^{(1)} &= g(z_3^{(1)}) \\ z_4^{(1)} &= w_4^{(1)T} x + b_4^{(1)}, & a_4^{(1)} &= g(z_4^{(1)}) \end{aligned}$$

Given input x :

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

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

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

$$a^{(2)} = g(z^{(2)})$$

Vectorizing Across Multiple Examples, WEEK 3

(3)

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

$a^{[2](i)}$ ↗ training example
↗ layer 2

for $i = 1$ to m :

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

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

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

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

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

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

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

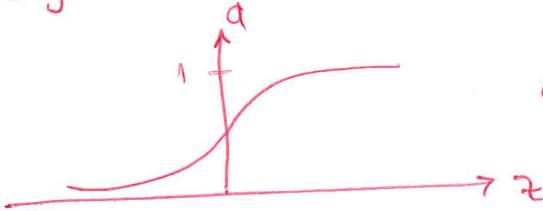
$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

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

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

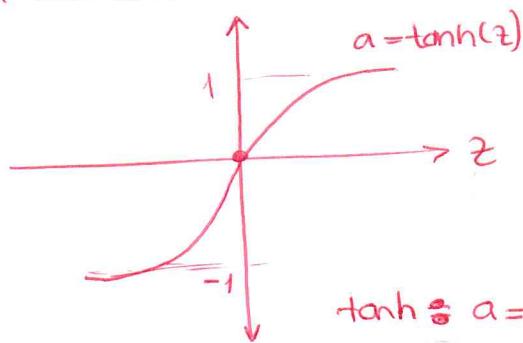
$$A^{[2]} = \sigma(Z^{[2]})$$

4

Sigmoid:

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

never use this,
except for the output
layer, if you are doing
binary classification. or
maybe almost never use this.

tanh function:

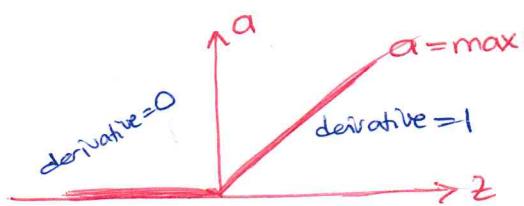
Superior than sigmoid.

When z is very large or very small, the slope of the function would be very close to 0.
This can slow down gradient descent

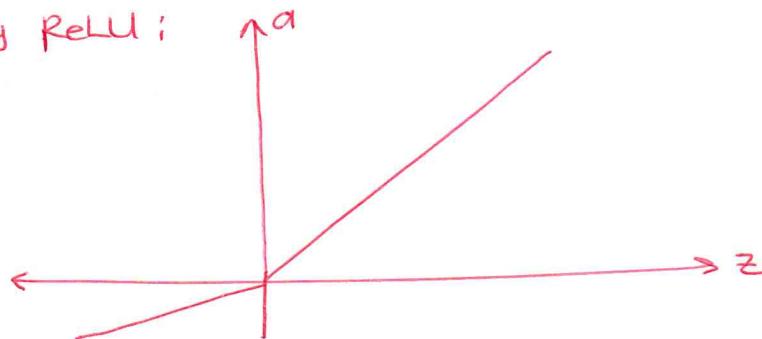
more popular

ReLU function:

(Rectified Linear Unit)



$$\text{ReLU} \Rightarrow a = \max(0, z)$$

Leaky ReLU:

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

Why do we need non-linear activation functions?

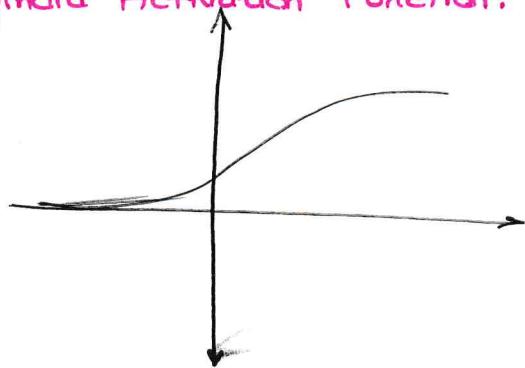
linear activation function = identity activation function

There is just one place where you might use a linear activation function: if you are doing machine learning on the regression problem.

So the one place you might use a linear activation function is usually in the output layer.

Composition of two linear functions is another linear functions.

Sigmoid Activation Function:



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

$$g'(z)$$

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

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

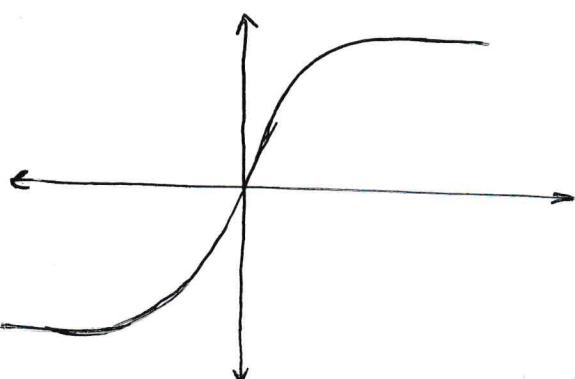
$$= g(z) (1-g(z))$$

$$z=10 \quad g(z) \approx 1 \quad \frac{d}{dz} g(z) \approx 1 \cdot (1-1) \approx 0$$

$$z=-10 \quad g(z) \approx 0 \quad \frac{d}{dz} g(z) \approx 0 \cdot (1-0) \approx 0$$

$$z=0 \quad g(z)=\frac{1}{2} \quad \frac{d}{dz} g(z) = \frac{1}{2} \left(1 - \frac{1}{2} \right) = \frac{1}{4}$$

Tanh Activation Function:



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

$$= \frac{e^z - e^{-z}}{e^z + e^{-z}} = a \quad g'(z) = 1 - a^2$$

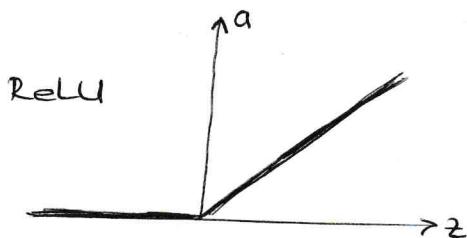
$$g'(z) = \frac{d}{dz} g(z) = 1 - (\tanh(z))^2$$

$$z=10 \quad \tanh(z) \approx 1 \quad g'(z) = 0$$

$$z=-10 \quad \tanh(z) \approx -1 \quad g'(z) = 0$$

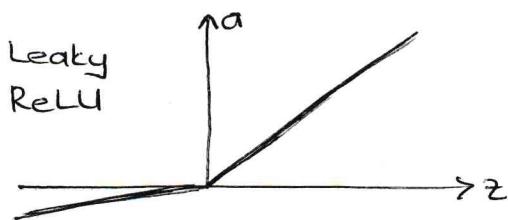
$$z=0 \quad \tanh(z) = 0 \quad g'(z) = 1$$

ReLU and Leaky ReLU:



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

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



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

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

Gradient Descent For Neural Networks

(6)

Parameters: $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}$
 $(n^{(1)}, n^{(0)})$ $(n^{(1)}, 1)$ $(n^{(2)}, n^{(1)})$ $(n^{(2)}, 1)$

$n_x = n^{(0)}$, $n^{(1)}$, $n^{(2)} = 1$
 ↓
 input feature
 hidden units
 output units

Cost Function: $J(W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i)$

When training a neural network, it is important to initialize the parameters randomly rather than to all zeros.

Gradient Descent :

Repeat {

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

$dW^{(1)}, dB^{(1)}, dW^{(2)}, dB^{(2)}$

$W^{(1)} := W^{(1)} - \alpha \cdot dW^{(1)}$

$b^{(1)} := b^{(1)} - \alpha \cdot dB^{(1)}$

 ⋮

}

Forward Propagation:

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

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

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

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

if this is binary classification

Backpropagation:

$$dZ^{(2)} = A^{(2)} - Y$$

$$dW^{(2)} = \frac{1}{m} \cdot dZ^{(2)} \cdot A^{(1)T}$$

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

$$dZ^{(1)} = \underbrace{W^{(2)T} dZ^{(2)}}_{(n^{(2)}, m)} * \underbrace{g^{(1)'}(z^{(1)})}_{(n^{(1)}, m)}$$

it prevents Python from outputting one of those funny rank 1 arrays

$$dW^{(1)} = \frac{1}{m} dZ^{(1)} X^T$$

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

Backpropagation Intuition

WEEK 3

7

$$dz^{(2)} = a^{(2)} - y$$

$$dw^{(2)} = dz^{(2)} a^{(1)T}$$

$$db^{(2)} = dz^{(2)}$$

$$dz^{(1)} = w^{(2)T} dz^{(2)} * g'(z^{(1)})$$

$$dw^{(1)} = dz^{(1)} x^T$$

$$db^{(1)} = dz^{(1)}$$

$$w^{(2)} \quad \frac{\text{dimensions}}{(n^{(2)}, n^{(1)})}$$

$$z^{(2)}, dz^{(2)} \quad (n^{(2)}, 1) \rightarrow (1, 1)$$

$$z^{(1)}, dz^{(1)} \quad (n^{(1)}, 1) \quad \xrightarrow{\text{element wise}}$$

$$dz^{(1)} = w^{(2)T} dz^{(2)} * g'(z^{(1)})$$

$$(n^{(1)}, 1) = (n^{(1)}, n^{(2)}) \quad \downarrow \quad (n^{(2)}, 1) * (n^{(1)}, 1)$$

$$dw^{(1)} = dz^{(1)} \cdot x^T$$

$$db^{(1)} = dz^{(1)}$$

$$dz^{(2)} = A^{(2)} - y$$

$$dw^{(2)} = \frac{1}{m} dz^{(2)} \cdot A^{(1)T}$$

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

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

$$dw^{(1)} = \frac{1}{m} dz^{(1)} X^T$$

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

Vectorized Implementation

$$z^{(1)} = w^{(1)} X + b^{(1)}$$

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

was mentioned in previous page.

$$dz^{(2)} = a^{(2)} - y$$

$$dw^{(2)} = dz^{(2)} \cdot a^{(1)T}$$

$$db^{(2)} = dz^{(2)}$$

why?

" $dw \neq dz \cdot x$ "

$w_2^{(1)} = [\quad]$
for one output



Random Initialization

if W ~~and b~~ ^{is} initialized with 0, all units in hidden layer will produce the same value and will be symmetric.

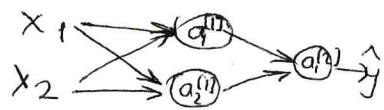
As a solution : you initialize these parameters (W) with ^{very small} random values and you prevent the symmetry ^{breaking} problem.

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

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

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

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



if we give ~~very~~ big numbers to W (initially), the derivatives would be very small ^(close to 0) and learning would be very slow.

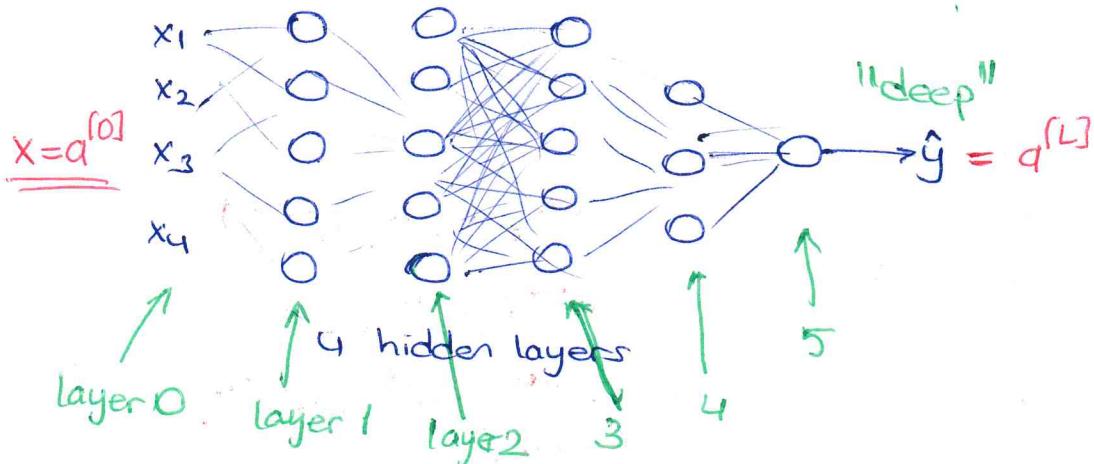
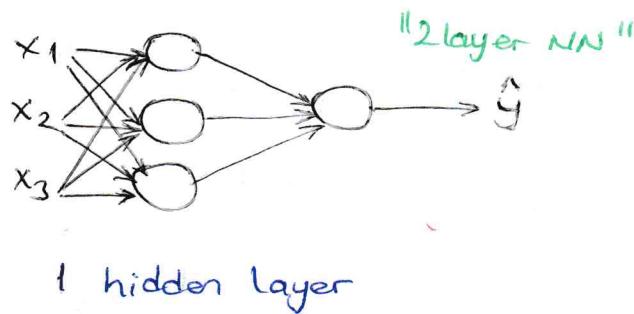
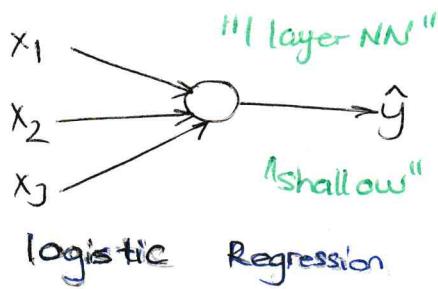
Notes :

1. The tanh activation usually works better than sigmoid activation function for hidden ~~units~~ units because the mean of its output is closer to zero, so it centers the data better for the next layer. The output of the tanh is between -1 and 1, it thus centers the data which makes the learning simpler for the next layer.

Deep L-layer Neural Network

WEEK 4

(1)



$$l = 4 \ (\# \text{layers})$$

$$n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 5, n^{[4]} = 3, n^{[5]} = 1$$

$n^{[l]}$ = # units in layer l

$$n^{[0]} = n_x = 3$$

$$n^{[L]} = 1$$

$a^{[l]}$ = activation in layer l

$$a^{[l]} = g^{[l]}(z^{[l]}), \quad w^{[l]} = \text{weights for } z^{[l]}$$

$$b^{[l]}$$

Forward Propagation in a Deep Network

$$Z^{[1]} = W^{[1]} A^{[0]} + b^{[1]} \quad X = A^{[0]}$$

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

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

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

$$Z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(Z^{[l]})$$

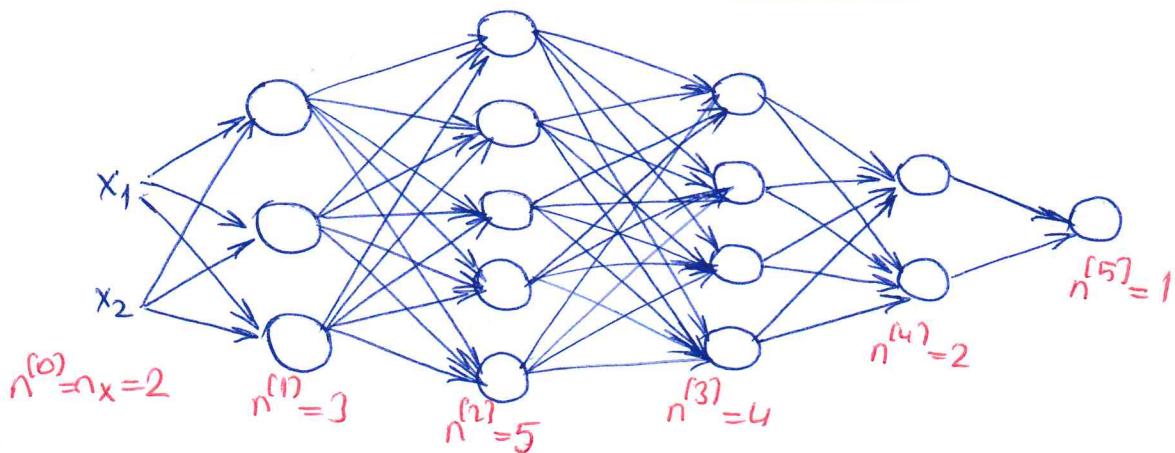
for loop

for $l = 1 \dots L$

There is no way to implement this without an explicit For loop.

In this place, it's perfectly okay to have an explicit For loop.

Getting Your Matrix Dimensions Right



$$z^{(l)} = w^{(l)} \cdot x + b^{(l)}$$

$$\begin{matrix} (3,1) & (3,2) & (2,1) & (3,1) \\ (n^{(1)}, 1) & (n^{(1)}, n^{(0)}) & (n^{(2)}, 1) & (n^{(2)}, 1) \end{matrix}$$

$$* \quad w^{(1)} = (n^{(1)}, n^{(0)})$$

$$* \quad w^{(2)} = (5, 3) \quad (n^{(2)}, n^{(1)})$$

$$* \quad z^{(2)} = w^{(2)} \cdot a^{(1)} + b^{(2)}$$

$$* \quad (5,1) \quad (5,3) \quad (3,1) \quad (5,1) \quad (n^{(2)}, 1)$$

$$w^{(3)} = (4, 5)$$

$$w^{(4)} : (3, 4) , \quad w^{(5)} : (1, 2)$$

$w^{(l)} : (n^{(l)}, n^{(l-1)})$
$b^{(l)} : (n^{(l)}, 1)$
$dw^{(l)} : (n^{(l)}, n^{(l-1)})$
$db^{(l)} : (n^{(l)}, 1)$

Vectorized Implementation

$$z^{(l)}, a^{(l)} : (n^{(l)}, 1)$$

$$z^{(l)}, A^{(l)} : (n^{(l)}, m)$$

$$dz^{(l)}, dA^{(l)} : (n^{(l)}, m)$$

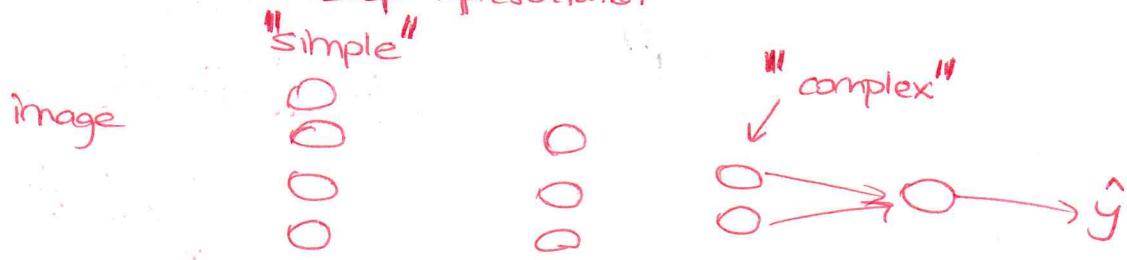
$$z^{(l)} = w^{(l)} \cdot X + b^{(l)}$$

$(n^{(1)}, m) \quad (n^{(1)}, n^{(0)}) \quad (n^{(0)}, m) \quad (n^{(l)}, 1)$

↓
with broadcasting
it acts like $(n^{(1)}, m)$

Why deep representation?

Intuition about deep representation

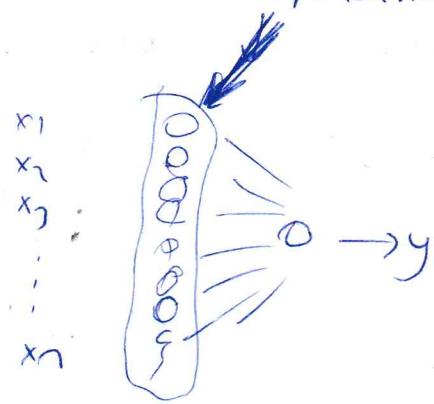


Audio → low level audio → phoneme → words → sentence
CAT → phrases

CIRCUIT THEORY and DEEP LEARNING

Informally: There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden ~~units~~ to compute.

it necessitate exponentially large units



deep learning is a branding. It is actually neural network with many hidden layers.

Forward and Backward Functions

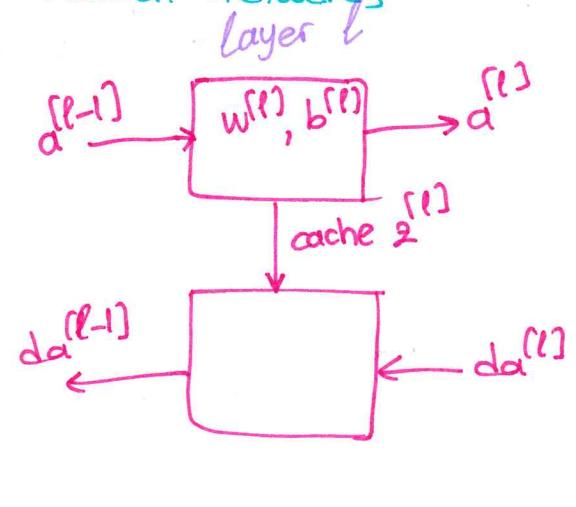
Layer l : $w^{[l]}, b^{[l]}$

Forward: Input $a^{[l-1]}$, output $a^{[l]}$
cache $z^{[l]}$

$$z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

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

Building blocks of deep neural networks



Backward: Input $da^{[l]}$, output $da^{[l-1]}$
cache $z^{[l]}$
 $dw^{[l]}$
 $db^{[l]}$

Forward Propagation for layer l

$$z^{[l]} = w^{[l]} \cdot A^{[l-1]} + b^{[l]}$$

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

Backward Propagation for layer l

$$\begin{aligned} dz^{[l]} &= da^{[l]} * g^{[l]}'(z^{[l]}) \\ dw^{[l]} &= dz^{[l]} * a^{[l-1]} \\ db^{[l]} &= dz^{[l]} \\ da^{[l-1]} &= w^{[l]T} \cdot dz^{[l]} \\ dz^{[l]} &= w^{[l+1]T} \cdot dz^{[l+1]} * g^{[l+1]}'(z^{[l]}) \end{aligned}$$

$$\begin{aligned} dz^{[l]} &= dA^{[l]} * g^{[l]}'(z^{[l]}) \\ dw^{[l]} &= \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T} \\ db^{[l]} &= \frac{1}{m} \text{np.sum}(dz^{[l]}, \text{axis}=1, \text{keepdims}=\text{True}) \\ dA^{[l-1]} &= W^{[l]T} \cdot dz^{[l]} \end{aligned}$$

Parameters vs Hyperparameters

WEEK 4

(5)

Parameters: $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, w^{(3)}, b^{(3)}, \dots$

Hyperparameters : Learning rate

iterations

hidden layers L

hidden units $n^{(1)}, n^{(2)}, \dots$

choice of activation function

} determine
the final
parameters (w, b)

Later: momentum, mini-batch size, regularization parameters

Applied deep learning is a very empirical process.