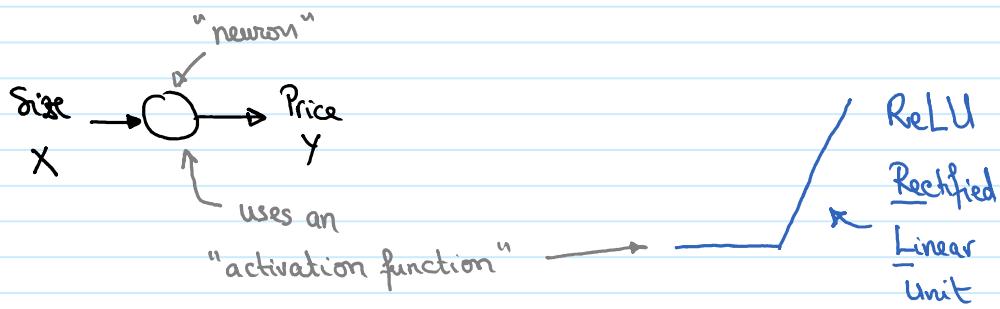
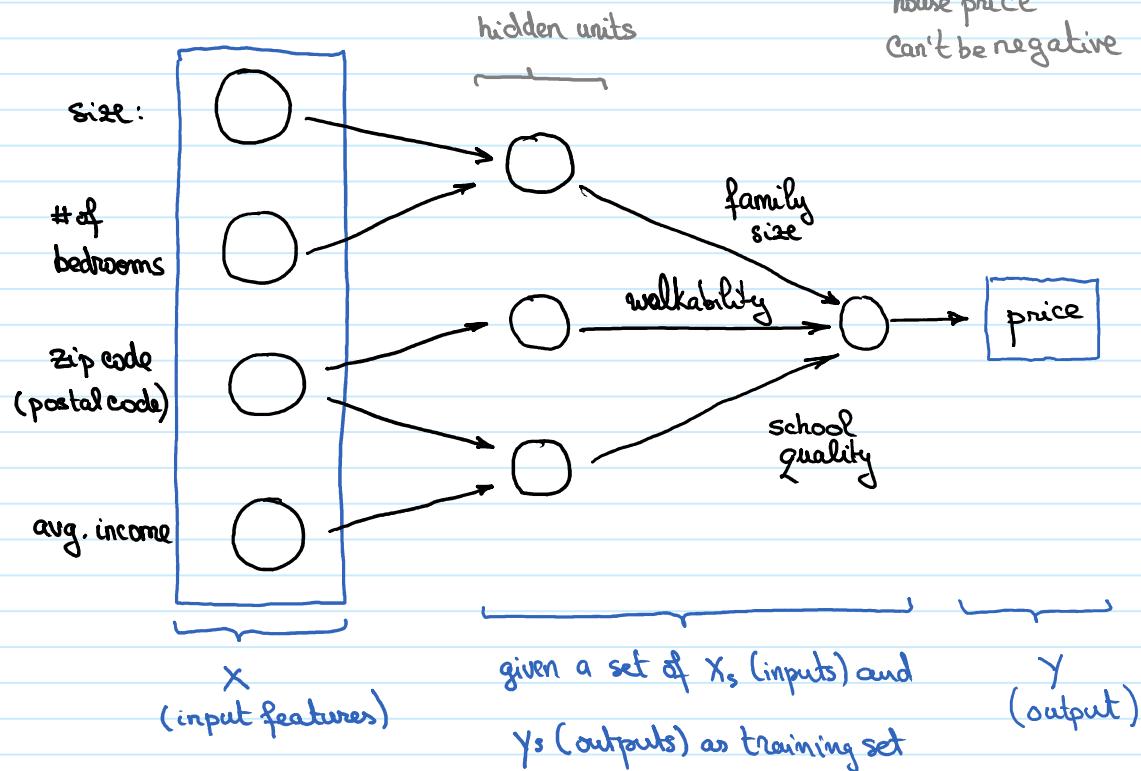
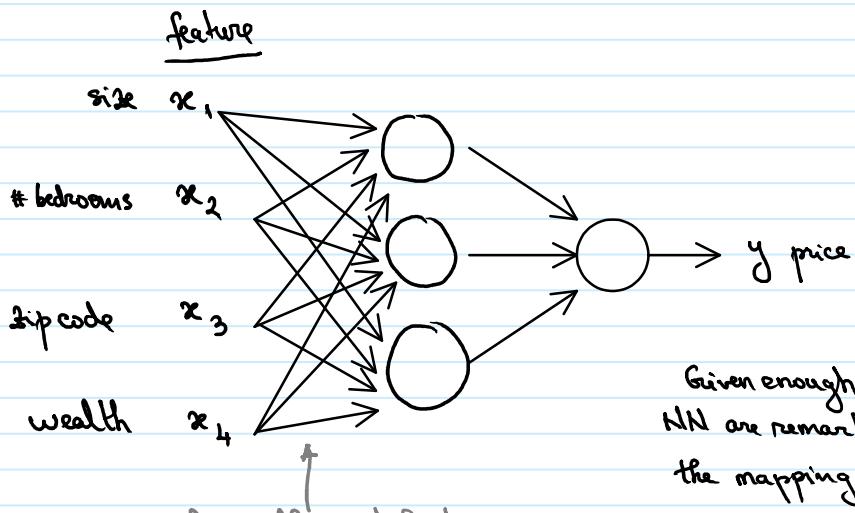


Week 1 Intro to Neural Networks

Saturday, 23 June, 2018 20:30

Simple model - housing price prediction





Given enough training data (x, y)
 NN are remarkably good @ finding
 the mapping function of $x \rightarrow y$

when all input features
 are connected to all nodes
 the network is called
 densely connected (or fully connected)

SUPERVISED) LEARNING EXAMPLES

Input (x)	Output (y)	Application
Home features	Price	Real Estate
Ad, user info	Click on ad? (0/1)	Online Advertising
Image	Object (1 ... 1,000)	Photo tagging
Audio	Text transcript	Speech recognition
English	Chinese	Machine translation
Image, Radar info	Position of other cars	Autonomous driving

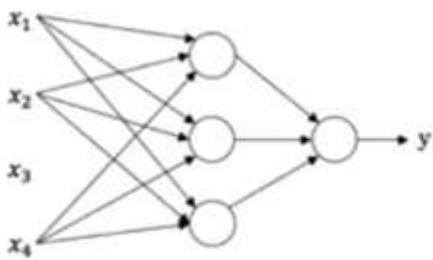
} standard
NN

} convolution
NN (CNN)

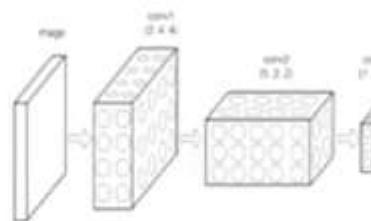
} Recurrent
NN
(RNN)

} advance
NN

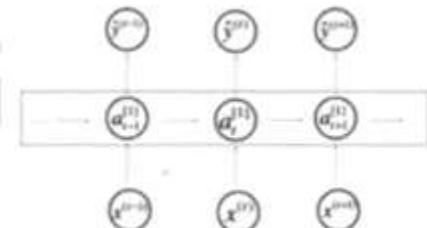
TYPES OF NEURAL NETS



Standard NN



Convolutional NN



Recurrent NN

Screen clipping taken: 2018-06-23 21:38

where there are spatial relationships between input features
(image processing)

where the input data is a sequence that has information encoded in its temporal order
(language is commonly represented as an RNN)

Structured Data

Size	#bedrooms	...	Price (1000\$)
2100	3		400
1600	3		330
2400	3		369
:	:		:
3000	4		540

User Age	Ad Id	...	Click
41	93242		1
80	93287		0
18	87312		1
:	:		:
27	71244		1

Screen clipping taken: 2018-06-23 21:43

- Each of the input features & output labels has a well defined value / meaning
→ more economic value
(better ads system, product recommendations)

Unstructured Data



Audio



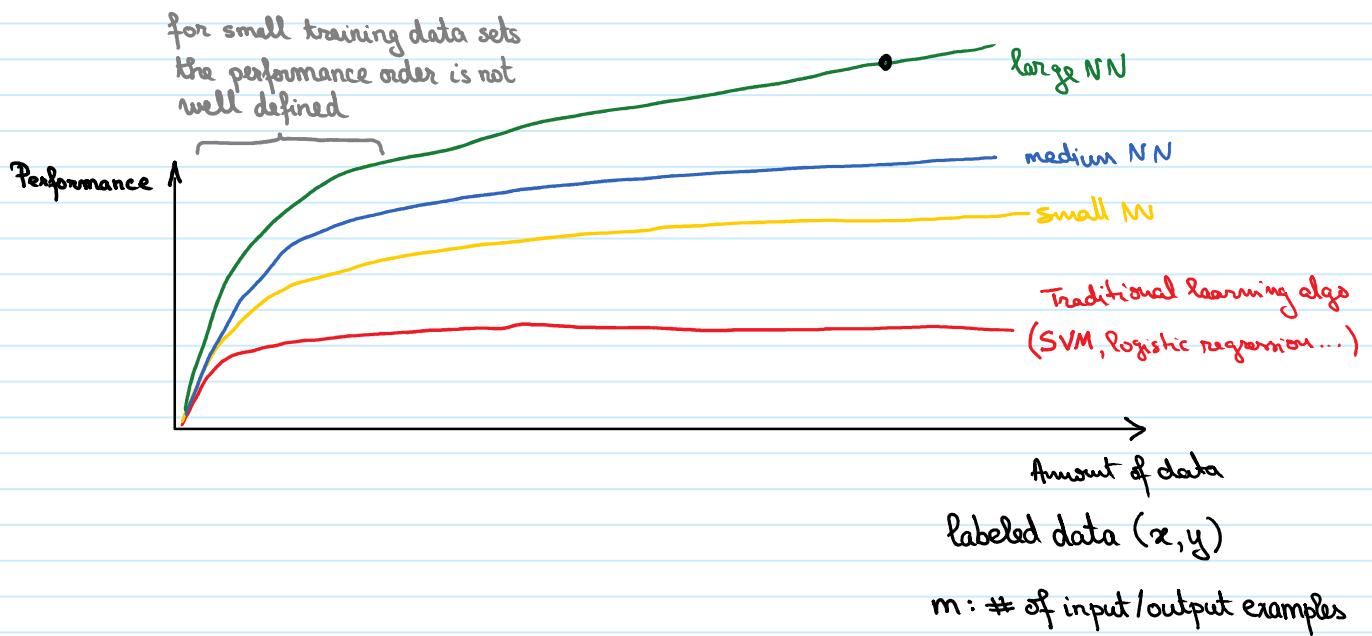
Image

Four scores and seven years ago...

Text

Pictures or raw audio or text

- The input features can be
→ pixel values in an image
→ individual words in a text
- Computers has more trouble interpreting unstructured data.
→ more news coverage



Scale of NN (many layers, many nodes, many connection)

+

Scale of input data (large digitized input sets)

+

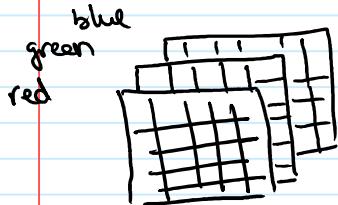
fast hardware

High performance outputs of neural nets.

Week 2 Neural Networks Basics

Sunday, 24 June, 2018 22:56

Input (x)
an RGB 64×64 px img. \rightarrow Output $y = \{1; 0\}$ cat or no cat.



} unroll each
channel &
concatenate
them into one
long arr vector

$$x = \begin{bmatrix} 255 \\ 251 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix} \quad \left. \begin{array}{l} \text{red} \\ \vdots \\ \text{grn} \\ \vdots \\ \text{blue} \end{array} \right\} \quad \begin{array}{l} \text{size of } x = n_x = \\ 64 \times 64 \times 3 = \\ 12,288 \text{ values.} \\ \text{↳ this } x \text{ is} \\ \text{the input feature} \\ \text{vector} \end{array}$$

Notation for the course

Thursday, 05 July, 2018 23:02

Binary Classification:

Single training example is represented by a pair:

$$(x, y)$$

where x is an n_x dimensional feature vector

$$\text{i.e. } x \in \mathbb{R}^{n_x}$$

and y (the label) is $\{0, 1\}$

The training set is composed of m training examples:

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

$m = m_{\text{train}}$ vs $m_{\text{test}} = \# \text{ of test examples.}$

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \quad \begin{array}{c} \uparrow \\ n_x \text{ features} \\ \downarrow \end{array} \quad Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

$\xleftarrow{\hspace{1cm}} m \xrightarrow{\hspace{1cm}}$ examples

$$Y \in \mathbb{R}^{1 \times m}$$
$$Y \text{.shape} \rightarrow (1, m)$$

$$X \in \mathbb{R}^{n_x \times m} \rightarrow X \text{.shape} \rightarrow (n_x, m)$$

Logistic Regression

Monday, 25 June, 2018 21:00

Given an input feature vector \mathbf{x} , we want to determine \hat{y} (an estimate of y).

\hat{y} is the probability that $y = 1$ given input feature \mathbf{x} .

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

i.e. if \mathbf{x} is a picture, \hat{y} tells us what is the chance that \mathbf{x} is a cat picture

Inputs:

$\mathbf{x} \in \mathbb{R}^{n_x}$ (input feature vector)

$\mathbf{w} \in \mathbb{R}^{n_x}$ and $b \in \mathbb{R}^1$ (parameters)

↳ corresponds to y intercept

Outputs:

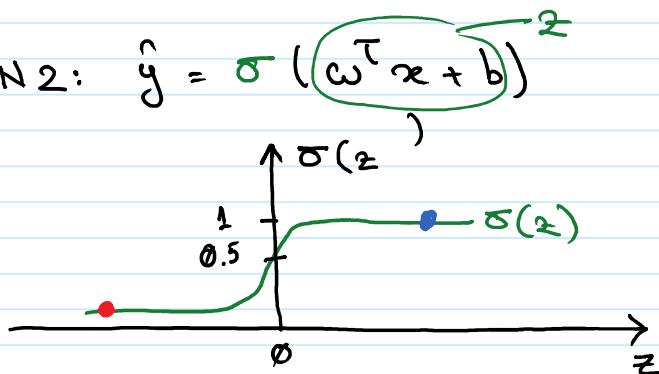
$$\hat{y}$$

- OPTION 1: $\hat{y} = \mathbf{w}^\top \mathbf{x} + b$ (this is the approach used for linear regression)

This doesn't work because \hat{y} is supposed to be a probability

i.e. $\hat{y} \in [0; 1]$

- OPTION 2: $\hat{y} = \sigma(\mathbf{w}^\top \mathbf{x} + b)$



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

- If $z \rightarrow \infty \Rightarrow \sigma(z) \rightarrow \frac{1}{1+0} = 1$

- If $z \rightarrow -\infty \Rightarrow \sigma(z) \rightarrow \frac{1}{1+\infty} = 0$

ASIDE:

In some courses w & b are not kept separate.

→ a new feature $x_0 := 1$ is defined such that $x \in \mathbb{R}^{n_x+1}$

$$\rightarrow \text{in this case } \hat{y} = \sigma(\theta^T x), \text{ where } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \left[\begin{array}{c} b \\ w \end{array} \right]$$

Logistic Regression Cost Function

Thursday, 05 July, 2018 23:06

Given:

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

→ a training set of m examples $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

we want to find $\hat{y}^{(i)}$ such that $\hat{y}^{(i)} \approx y^{(i)}$, where $i \in \{1 \dots m\}$

predictions we made using training example inputs (\hat{y})

match the "ground truth" (training example output y)

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

$x^{(i)}, y^{(i)}, z^{(i)}$ are parameters/variables/etc associated with the i^{th} training example

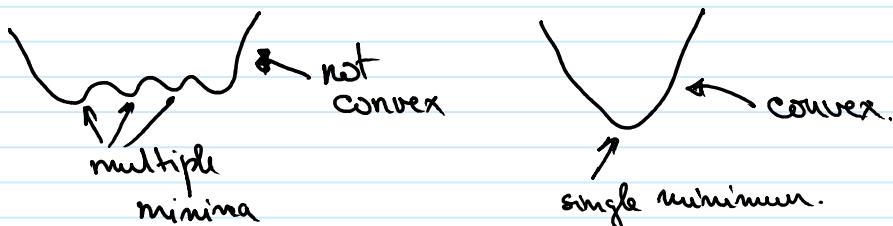
WHAT LOSS (OR ERROR) FUNCTION CAN WE USE?

A loss function measures how good the neural network predicted output \hat{y} is when the training value output (ground truth) is y

OPTION 1 AVERAGE SQUARE OF THE ERROR

$$\bullet L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

This loss function is not used much because the optimization problem is not convex
i.e. there is no global minimum towards which Gradient Descent can optimize



OPTION 2

LOGISTIC REGRESSION Loss (AKA CROSS ENTROPY LOSS)

$$\bullet \mathcal{L}(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

Intuition: we want $\mathcal{L}(\hat{y}, y)$ to be as small as possible when \hat{y} is close to y .

$$\rightarrow \text{If } y = 1 \Rightarrow \mathcal{L}(\hat{y}, y) = -\log \hat{y} \quad \begin{matrix} \leftarrow \text{want this to be small} \\ \rightarrow \text{want } \log \hat{y} \text{ large} \rightarrow \text{want } \hat{y} \text{ large} \end{matrix}$$

but $\hat{y} = \sigma(z) \Rightarrow \hat{y} = 1$ at maximum

$$\rightarrow \text{If } y = 0 \Rightarrow \mathcal{L}(\hat{y}, y) = -\log (1-\hat{y}) \quad \begin{matrix} \leftarrow \text{want this to be small} \\ \rightarrow \text{want } \log (1-\hat{y}) \text{ large} \rightarrow \text{want } 1-\hat{y} \text{ large} \\ \text{so want } \hat{y} \text{ small} \end{matrix}$$

but $\hat{y} = \sigma(z) \Rightarrow \hat{y} = 0$ at minimum

So by using the Logistic Regression loss function \hat{y} will be pushed towards y at the extremes.

COST FUNCTION:

Cost function averages the loss function over all training examples.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right]$$

The COST FUNCTION applies to all examples and it is used to determine the values of $w \notin b$.

The loss function $\mathcal{L}(\hat{y}, y)$ measures how well we are doing on a single example

Gradient descent

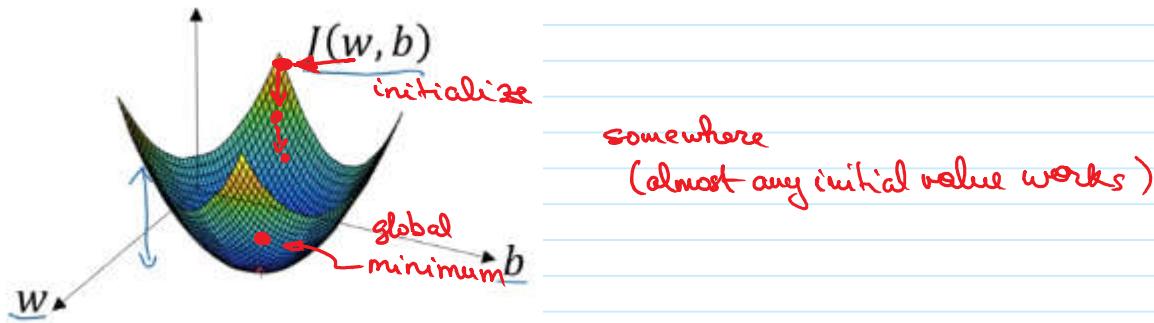
Thursday, 05 July, 2018 23:06

Given a cost function $J(w, b)$ over a given training set

$\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ we want to find parameters

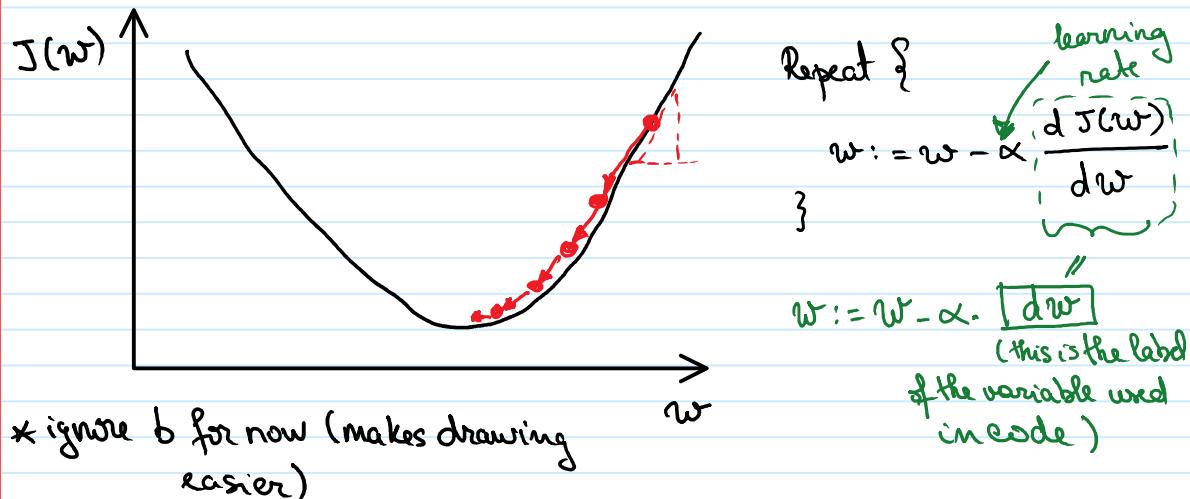
$w \notin b$ such that $J(w, b)$ is as small as possible.

i.e. find $w \notin b$ that minimize $J(w, b)$



Screen clipping taken: 2018-06-25 22:14

It turns out our $J(w, b)$ is a convex function (like a bowl)



When taking into account optimizing both $w \notin b$ we have

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

$$\left. \begin{array}{l} \frac{\partial J(w, b)}{\partial w} \\ \frac{\partial J(w, b)}{\partial b} \end{array} \right\}$$

"partial derivative"

$dJ(w, b)$ when b is eliminated

$$b := b - \alpha \frac{d J(w, b)}{d b}$$

$$\frac{\partial J(w, b)}{\partial b}$$

$\rightarrow d w \rightarrow$
*when implemented
in code*
 $\rightarrow d b \rightarrow$

Reminder:

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

Computation Graph

Thursday, 05 July, 2018 23:09

In a Neural Net:

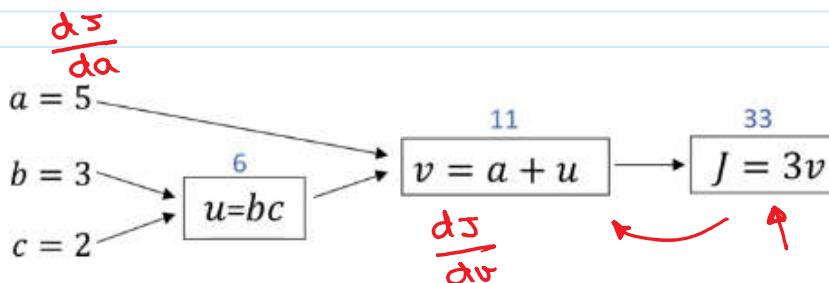
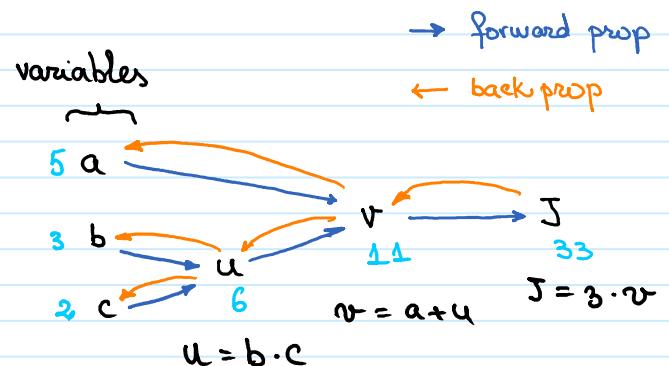
- Forward pass (**Forward propagation**) computes the output of the NN
- Backward pass (**Back propagation**) computes gradients (derivatives)

COMPUTATION GRAPH:

$$J(a, b, c) = 3(a + bc)$$

$$\begin{array}{c} \underbrace{a}_{\text{variables}} \\ \underbrace{b}_{\text{variables}} \\ \underbrace{c}_{\text{variables}} \end{array}$$

$$\underbrace{\underbrace{a + bc}_{u} \underbrace{v}_{\text{variables}}}_{J}$$



Screen clipping taken: 2018-06-25 23:01

$\frac{dJ}{dv} = ?$ if we change the value of v a little bit how does the value of J change?

$$J = 3v$$

$$\left. \begin{array}{l} v = 11 \rightarrow 11.001 \\ J = 33 \rightarrow 33.003 \end{array} \right\} \Rightarrow \boxed{\frac{dJ}{dv} = 3}$$

$$\left. \begin{array}{l} a = 5 \rightarrow 5.001 \\ v = 11 \rightarrow 11.001 \end{array} \right\} \boxed{\frac{dv}{da} = 1}$$

$$J = 33 \rightarrow 33.001$$

$$\frac{dJ}{da} = ?$$

$$= \boxed{\frac{dJ}{dv}} \times \boxed{\frac{dv}{da}}$$

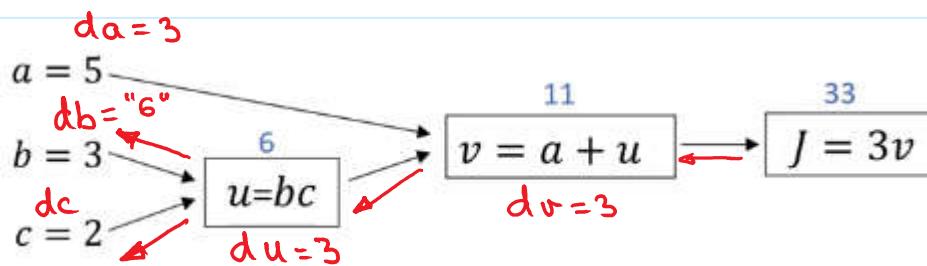
$$= 3 \times 1 = 3$$

aka "chain rule"

$$a \rightarrow v \rightarrow j \Rightarrow \frac{dJ}{da} = \frac{dJ}{dv} \times \frac{dv}{da}$$

$\frac{dJ}{da} \rightarrow "da" = 3$
 $\frac{dJ}{dv} \rightarrow "dv" = 3$

d Final Output Var
 d_{var}
 where d_{var} can be $\{a, b, c, u, v\}$



Screen clipping taken: 2018-06-25 23:14

$$\frac{dJ}{du} = \underbrace{\frac{dJ}{dv}}_3 \cdot \underbrace{\frac{dv}{du}}_1 = 3 \cdot 1 = 3$$

$$\frac{dJ}{db} = \underbrace{\frac{dJ}{dv}}_3 \cdot \underbrace{\frac{dv}{du}}_1 \cdot \underbrace{\frac{du}{db}}_{c=2} = 3c = 6$$

$$\frac{dJ}{dc} = 3 \cdot b = 9$$

Logistic Regression Gradient Descent

Thursday, 05 July, 2018 23:22

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

For a given example a

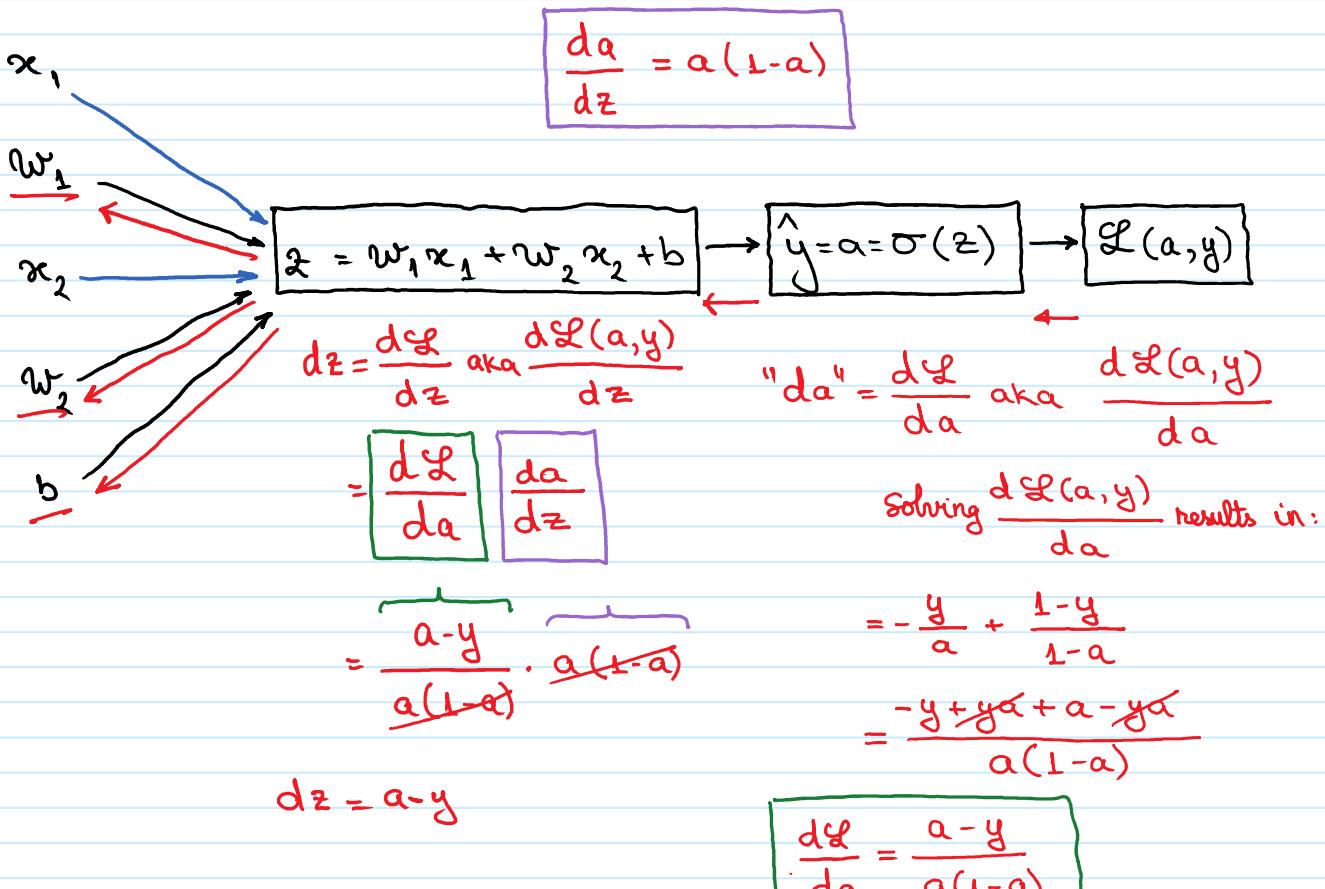
$$\hat{y} = a = \sigma(z)$$

$$\mathcal{L}(\hat{y}, y) = - (y \log(a) + (1-y) \log(1-a))$$

Let's say we only have two features x_1 & x_2

We also have to figure out w and b (want to modify w_1, w_2 & b to minimize \mathcal{L})

Remember: $a = \sigma(z) = \frac{1}{1+e^{-z}} \Rightarrow \frac{da}{dz} = \frac{d}{dz} \frac{1}{1+e^{-z}}$



$$\frac{d\mathcal{L}}{dw_i} = "d_{w_i}" = x_i \cdot d\mathcal{L} = x_i(a-y)$$

$$dw_1$$

$$dw_2 = x_2 \cdot dz = x_2(a-y)$$

$$db = dz = a-y$$

Repeat {

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

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

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

}

This is for a SINGLE example

For multiple examples we use $J(w, b)$ the cost function.

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

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

So we know how to compute $dw_1^{(i)}, dw_2^{(i)}, db^{(i)}$

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

$dw_1^{(i)}$ for the training example
 $(x^{(i)}, y^{(i)})$

dw_1 = average dw_i over all training examples.

Non-vectorized Gradient Descent

Tuesday, 26 June, 2018 19:30

GRADIENT DESCENT ON M EXAMPLE NON-VECTORIZED

Initial values: $J = \emptyset$; $\underline{dw_1} = \emptyset$; $\underline{dw_2} = \emptyset$; $\underline{db} = \emptyset$

for loop { For $i := 0$ to m # iterate over all training examples.

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

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

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

$$d_2^{(i)} := a^{(i)} - y^{(i)}$$

for loop
for multiple
features

$$\left. \begin{aligned} dw_1 &:= dw_1 + x_1^{(i)} dz^{(i)} \\ dw_2 &:= dw_2 + x_2^{(i)} dz^{(i)} \end{aligned} \right\}$$

assume we only have 2 features otherwise calculate $d w_3, d w^4 \dots$

$$db := db + dg^{(i)}$$

$$\mathfrak{I} := \frac{\mathfrak{I}}{m}$$

$$d\omega_L := \frac{d\omega_i}{m}$$

$$d\omega_2 := \frac{d\omega_2}{m}$$

$$db := \frac{db}{m}$$

dW , is used as an accumulator

$$d\omega_1 = \frac{\partial \mathcal{J}}{\partial \omega_1}$$

$$w_1 := w_1 - \alpha d w_1$$

$$w_2 := w_2 - \alpha d w_2$$

$$b := b - \alpha d_b$$

All of the above implements a single step of Gradient Descent .

→ we need to repeat this for every step of Gradient Descent

The Problem with this algorithm is that it uses two for-loops .

→ for-loops are not computationally efficient and slow execution

Vectorization is used to implement this algorithm without using any for-loops

Vectorization

Friday, 06 July, 2018 20:44

The art of removing "for loops" in code.

$$z = w^T x + b \quad w = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad x = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad w \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x}$$

Non-vectorized implementation:

$$z = \emptyset$$

for i in range (n_x):

$$z += w[i] \times x[i]$$

$$z += b$$

Vectorized:

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

$$\underbrace{w^T x}_{w^T x}$$

np stands for numpy

GPU
CPU

} parallelization using SIMD (single instruction multiple data)

Whenever possible, avoid explicit for-loops:

$$U = A \cdot v$$

$$U_i = \sum_j A_{i,j} \cdot v_j$$

Non-vectorized

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

for $i = \dots$

for $j = \dots$

$$u[i] = A[i][j] \times v[j]$$

Vectorized

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

If we need to apply the exponential operation on every element of a matrix/vector.

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

Non-vectorized

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

for i in range(n):

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

Vectorized

$u = np.exp(v)$

Some Numpy vectorized functions are:

$np.exp(v)$

$np.abs(v)$

$v**2$

$np.log(v)$

$np.maximum(v, 0)$

$1/v$

$$\begin{aligned}
 J &= 0, \quad \boxed{\cancel{dw_1 = 0}, \cancel{dw_2 = 0}}, \quad db = 0 \quad dw = np.zeros((n_x, 1)) \\
 \text{for } i = 1 \text{ to } m: \\
 z^{(i)} &= w^T x^{(i)} + b \\
 a^{(i)} &= \sigma(z^{(i)}) \\
 J &+= -[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})] \\
 dz^{(i)} &= a^{(i)}(1 - a^{(i)}) \\
 \boxed{\left. \begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{array} \right\} n_x = 2} \quad dw += x^{(i)} dz^{(i)} \\
 J &= J/m, \quad \boxed{dw_1 = dw_1/m, \quad dw_2 = dw_2/m, \quad db = db/m} \quad dw /= m
 \end{aligned}$$

Andrew Ng

Screen clipping taken: 2018-06-26 20:29

VECTORIZING Logistic REGRESSION FORWARD PASS

Normally we would have to compute the following for each training example

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

To vectorize it we will use matrix multiplication:

$$Z = [z^{(1)}, z^{(2)}, \dots, z^{(m)}]$$

$w^T = [w_1, w_2, w_3, \dots, w_{n_x}]$, where n_x is the number of features

$$X = \left[\begin{array}{c|c|c|c}
 & & & \\
 x^{(1)} & x^{(2)} & \dots & x^{(m)} \\
 \hline
 \end{array} \right] \quad \left. \begin{array}{l} n_x \\ \text{features} \end{array} \right\} \quad \text{where } X \in \mathbb{R}^{n_x \times m}$$

m examples

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

$$Z := np.dot(w.T, x) + b$$

b will automatically be cast as a $[1 \times m]$ vector for this operation (see "Broadcasting" in Python)

vectorized

To calculate the a (aka \hat{y})

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

VECTORIZING LOGISTIC REGRESSION BACKWARD PASS

$$dz^{(1)} = a^{(1)} - y^{(1)}$$

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

$$dz^{(3)} = a^{(3)} - y^{(3)}$$

:

Let's define:

$$dZ = \underbrace{[dz^{(1)}, dz^{(2)}, \dots, dz^{(m)}]}_{m \text{ examples}}$$

then:

$$A = [a^{(1)}, \dots, a^{(m)}]$$

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

$$dZ = A - Y = [a^{(1)} - y^{(1)}, a^{(2)} - y^{(2)}, \dots]$$

To calculate db

$$d\mathbf{b} = \frac{1}{m} \sum_{i=1}^m d\mathbf{z}^{(i)}$$

$$= \frac{1}{m} * \text{np.sum}(d\mathbf{Z})$$

vectorized

$$d\mathbf{w} = \frac{1}{m} \mathbf{X} \cdot d\mathbf{Z}^\top$$

$$= \frac{1}{m} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \cdot \begin{bmatrix} d\mathbf{z}^{(1)} \\ d\mathbf{z}^{(2)} \\ \vdots \\ d\mathbf{z}^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} \left[x^{(1)} d\mathbf{z}^{(1)} + \dots + x^{(m)} d\mathbf{z}^{(m)} \right] \text{ an } m \times 1 \text{ vector}$$

Vectorized Gradient Descent

Friday, 06 July, 2018 20:49

GRADIENT DESCENT ON M-EXAMPLES VECTORIZED

for step in range (1000) : # of steps to iterate

$$\left[\begin{array}{l} Z := W^T X + b \\ \quad := \text{np.dot}(W.T, X) + b \\ \\ A := \sigma(Z) \\ \\ dZ := A - Y \\ \\ dW := \frac{1}{m} X dZ^T \\ \\ db := \frac{1}{m} \text{np.sum}(dZ) \\ \\ W := W - \alpha dW \\ \\ b := b - \alpha db \end{array} \right]$$

Why Sigma as cost function?

Friday, 06 July, 2018 20:53

WHY σ AS COST FNC.?

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

$\hat{y} = p(y=1 | x)$ probability of $y=1$ given a set of input features x

$$\begin{aligned} \text{If } y=1 : \quad p(y|x) &= \hat{y} \\ y=0 : \quad p(y|x) &= 1-\hat{y} \end{aligned} \quad \left. \right\}$$

These two equations above can be represented by

$$\begin{aligned} p(y|x) &= \hat{y}^y \cdot (1-\hat{y})^{1-y} \rightarrow y=1 \Rightarrow p(y|x) = \hat{y}^1 \cdot (1-\hat{y})^0 \\ &\qquad\qquad\qquad = \hat{y} \\ &\rightarrow y=0 \Rightarrow p(y|x) = \hat{y}^0 \cdot (1-\hat{y})^1 \\ &\qquad\qquad\qquad = 1-\hat{y} \end{aligned}$$

$$\log p(y|x) = \log(\hat{y}^y (1-\hat{y})^{1-y})$$

$$= y \log(\hat{y}) + (1-y) \log(1-\hat{y})$$

$$= -\mathcal{L}(\hat{y}, y) \quad \text{Negative of the loss function.}$$

Negative because when training an algorithm we want to make probabilities large, but the loss function should be minimized in logistic regression. So minimizing the loss function corresponds to maximizing the log of the probability.

Cost Function over all training set

Friday, 06 July, 2018 20:54

ENTIRE TRAINING SET

HOW DOES THE COST FUNCTION OVER ALL M EXAMPLES
Look Like?

$$p(\text{labels in training set}) = \prod_{i=1}^m p(y^{(i)} | x^{(i)})$$

assuming the training examples are drawn independently (i.e. are drawn identically, are drawn independently distributed)

So if we would like to carry out maximum likelihood estimation then we want to find the parameters that maximizes the chance the observation is in the training set.

Maximizing $p(\text{labels in training set})$ = maximizing the
 $\log p(\text{labels in training set})$

$$\begin{aligned} \log [p(\text{labels in training set})] &= \log \left(\prod_{i=1}^m p(y^{(i)} | x^{(i)}) \right) \\ &= \sum_{i=1}^m \underbrace{\log [p(y^{(i)} | x^{(i)})]}_{-\mathcal{L}(\hat{y}^{(i)}, y^{(i)})} \end{aligned}$$

The principle of maximum likelihood estimation which means to choose the parameters of maximizing $\log(p(\cdot))$

$$\Rightarrow \text{maximize } - \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

because we want to minimize the cost we get rid of the "-"

because we want to minimize the cost we get rid of the "-" sign and we also divide by m to better scale our quantities (for convenience) and we get the cost function

$$\text{Cost } J(w, b) = \frac{1}{m} \sum \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

So by minimizing $J(w, b)$ we are carrying out maximum likelihood estimation w/ the logistic model under the assumption that our training examples were iid (identically independently distributed).

! Type of Vector Multiplication

Thursday, 28 June, 2018 23:00

Given:

$$w = \begin{bmatrix} | & | \\ x_{n,1} & \dots & x_{n,m} \\ | & | \end{bmatrix}; y = [y_1, \dots, y_m]; k = [k_1, \dots, k_m]$$

"MATRIX MULTIPLICATION" (DOT PRODUCT) Inner Product

$$\text{np.dot}(w, y^T) = \begin{bmatrix} x_{1,1} \cdot y_1 + \dots + x_{1,m} \cdot y_m \\ \vdots \\ x_{n,1} \cdot y_1 + \dots + x_{n,m} \cdot y_m \end{bmatrix}$$

- dimension wise: $A[n, m] \times B[p, r] = C[n, r]$

OUTER PRODUCT

$$\text{np.outer}(y, k) = \begin{bmatrix} y_1 \cdot k_1 & y_1 \cdot k_2 & \dots & y_1 \cdot k_m \\ \vdots & & & \\ y_m \cdot k_1 & \dots & & y_m \cdot k_m \end{bmatrix}$$

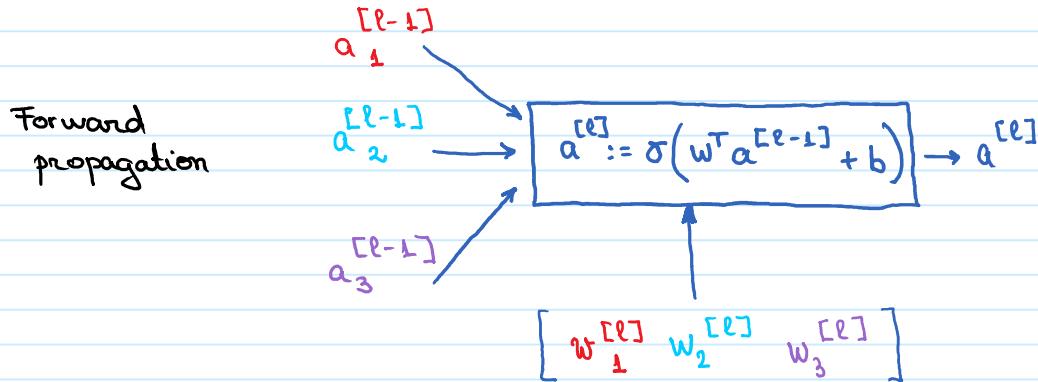
ELEMENT-WISE MULTIPLICATION

$$y * k = [y_1 \cdot k_1 \quad y_2 \cdot k_2 \quad \dots \quad y_m \cdot k_m]$$

! MII Summary

Monday, 20 August, 2018 23:36

Given a Neural Network with L total layers on layer l:



For the 1st layer ($l=1$) $a^{[l-1]} = x$ (input features from training example)

To the last layer ($l=L$) $a^{[l]} = \hat{y}$ (predicted output)

Given a training set of m examples (x, y) we use Gradient Descent to find the vector w and scalar b such that $J(w, b, x, y)$ is minimized.

$$g \left(m \left\{ \begin{bmatrix} x_1^{(1)} & | & | \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(m)} \\ \vdots & \vdots & & \vdots \\ x_n^{(1)} & | & | \end{bmatrix} \cdot \underbrace{\begin{bmatrix} w_1 & w_2 & \dots & w_n \end{bmatrix}^T}_{w} + b \right\} \right) = \hat{y}$$

g is known as the activation function and it must be non-linear.

i.e σ , tanh, ReLU

$$J() = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}, y) \quad (\text{average loss across all examples, for a given definition of loss } \mathcal{L})$$

OBSERVATIONS:

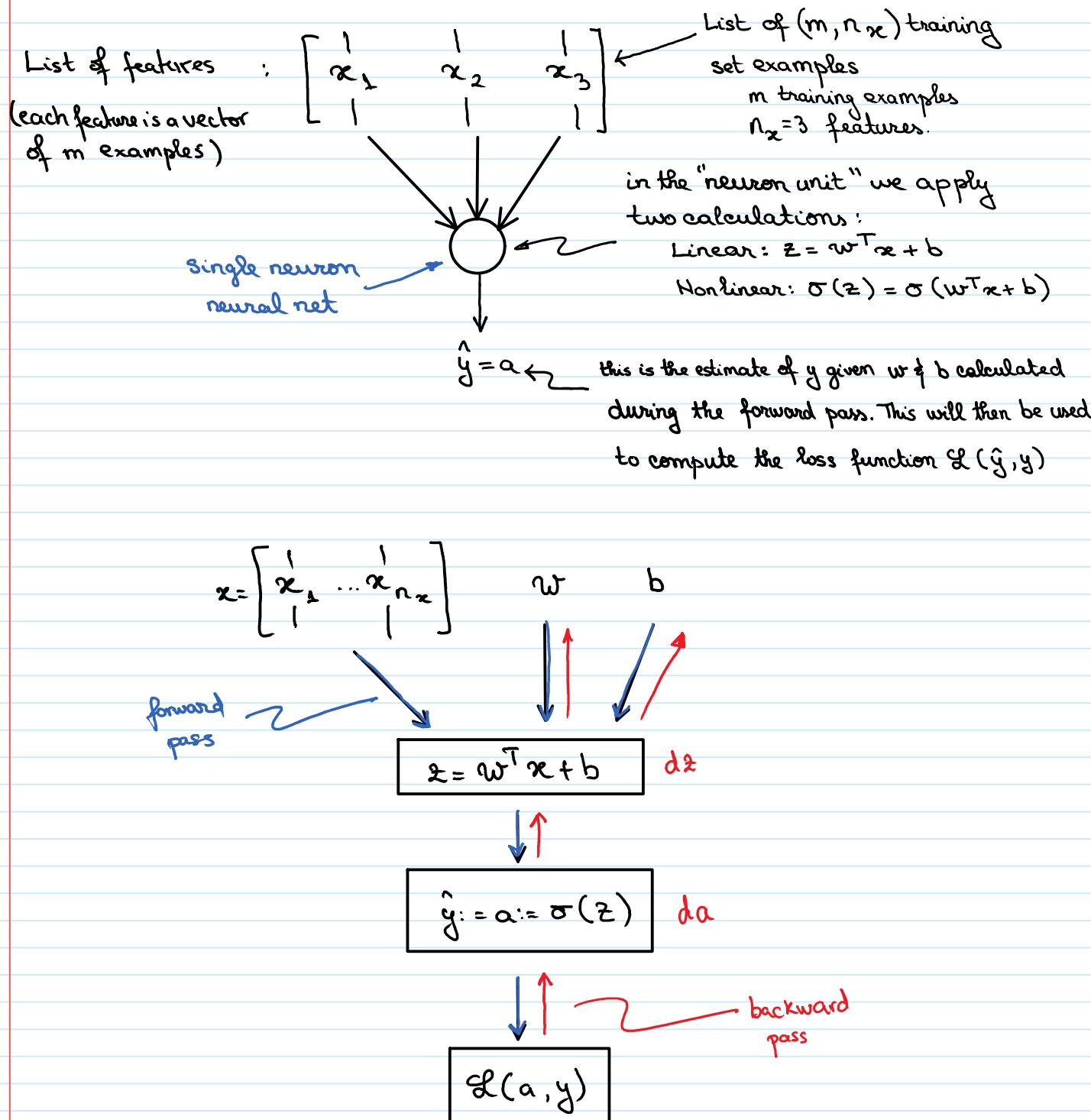
- Each value in vector w is associated with each feature in the training data set
- w_1 is derived by running Gradient Descent on all m examples of feature 1
(similarly for the other weights)
- The same scalar b applies to all m training examples and all features.

Week 3 Shallow NNs

Sunday, 01 July, 2018 13:05

Overview of NN implementation

Friday, 06 July, 2018 00:02



One forward pass outputs $m \times \hat{y}$ predictions which are then used to calculate the cost function:

$$J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$J = \frac{1}{m} \sum_i \text{se}(\hat{y}^{(i)}, y^{(i)})$$

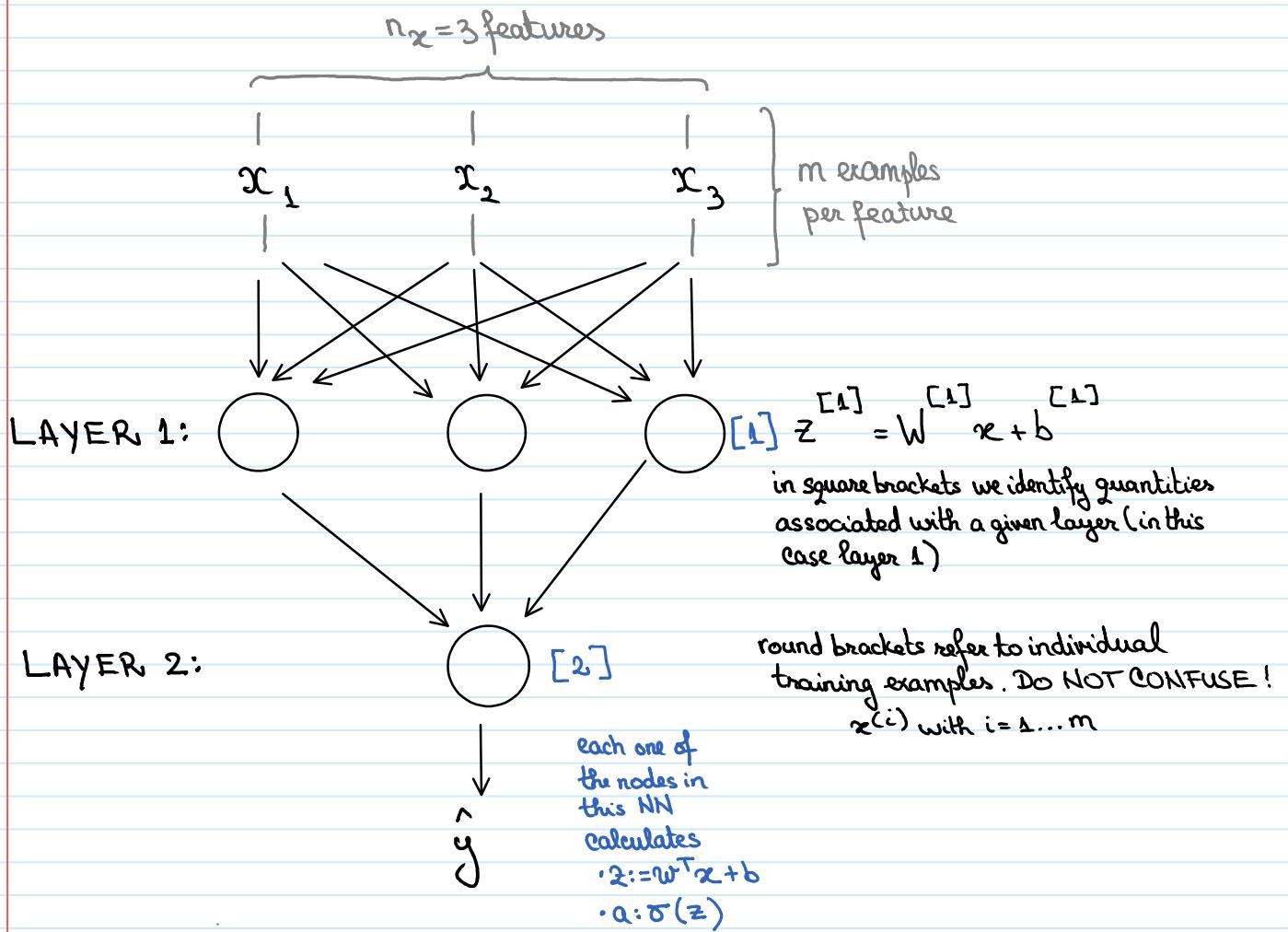
Then in the backward pass the cost function J is used to determine da .

da is a vector of gradients (one value per input feature) which is used to determine $d w$ (the gradient by which we reduce the weights to minimize J).

We then repeat the forward & backward pass until we converge on a minimum cost J

! A more general neural net

Friday, 06 July, 2018 00:02



DETAILED DESCRIPTION

$n_x = 3 \text{ features}$

$x = \begin{bmatrix} | & | & | \\ x_1 & x_2 & x_3 \\ | & | & | \end{bmatrix} \quad m \text{ examples}$

$w^{[1]}, b^{[1]}$

$d w^{[1]}, d b^{[1]}$

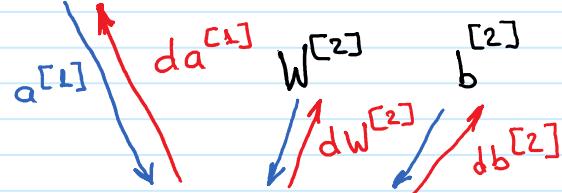
$z^{[1]} = W^{[1]} \cdot x + b^{[1]}$

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

$$\downarrow \uparrow dz^{[1]}$$

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

Q: Are $w^{[1]}$ & $b^{[1]}$ containing parameters for all neurons in layer 1?



A: Yes. In this case:

$$w.\text{shape} = [3, 3]$$

$$b.\text{shape} = [3, 1]$$

(more details below)

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

$$\downarrow \uparrow dz^{[2]}$$

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

$$\hat{y} = a^{[2]} \quad \uparrow \quad da^{[2]}$$

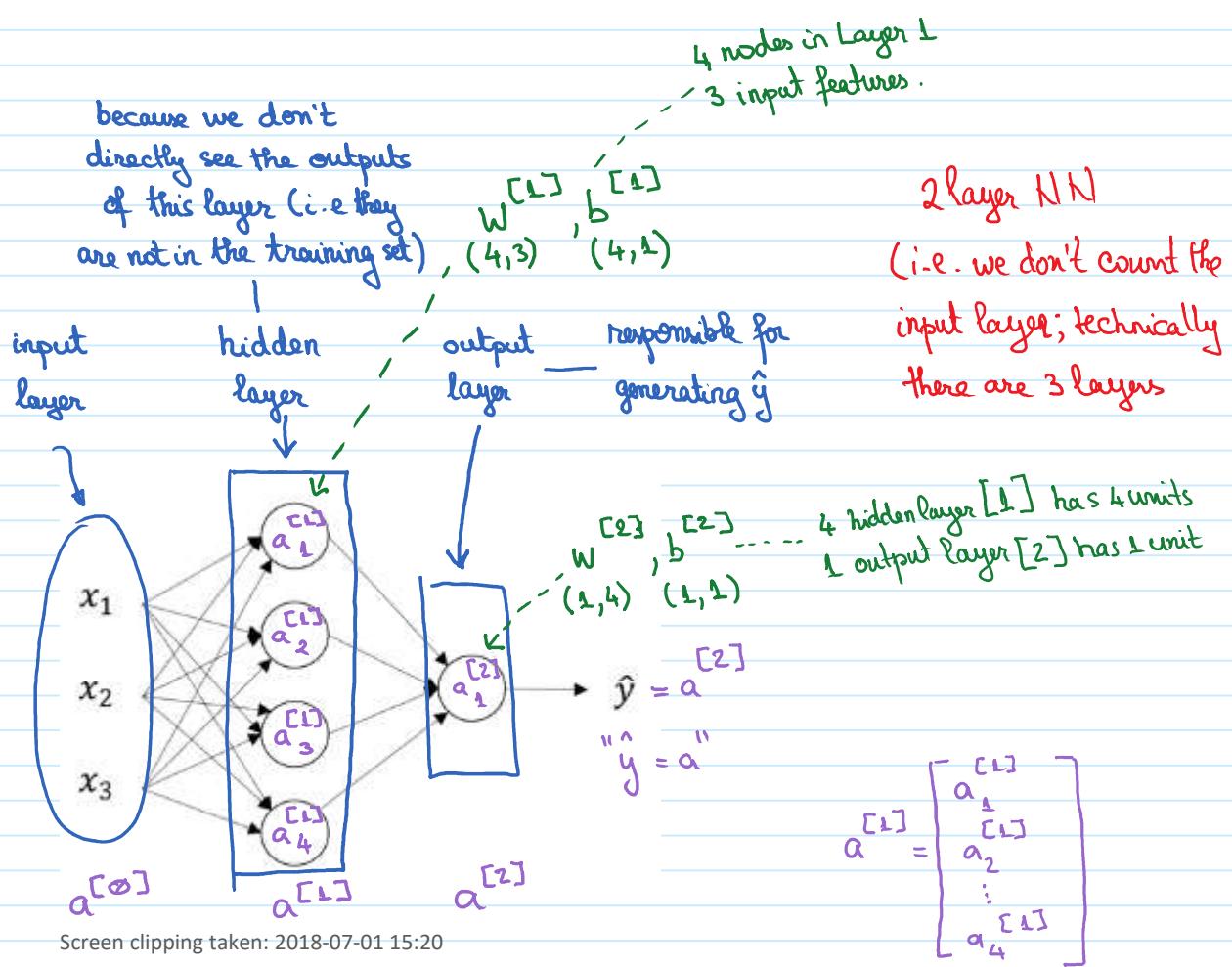
$$\mathcal{L}(a^{[2]}, y)$$

Summary:

A neural net w/ 2 layers takes a logistic regression and repeats it twice.

Neural Network Representation

Friday, 06 July, 2018 00:02

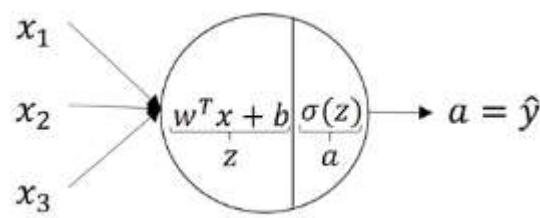


$a^{[\emptyset]} := X$ (the input features)

where " a " stands for activations \rightarrow it refers to the values different layers of the NN are passing on to the next layer

Computing NN output for a single example

Monday, 02 July, 2018 15:00



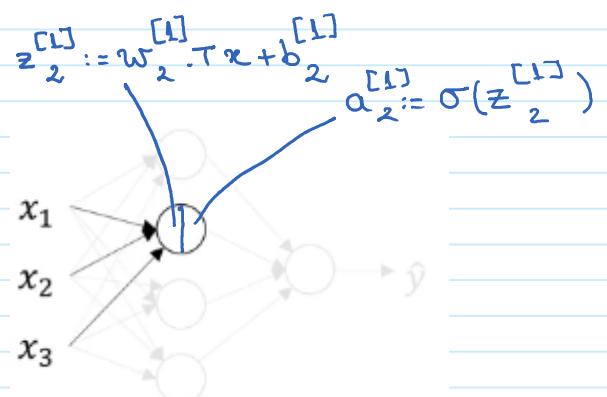
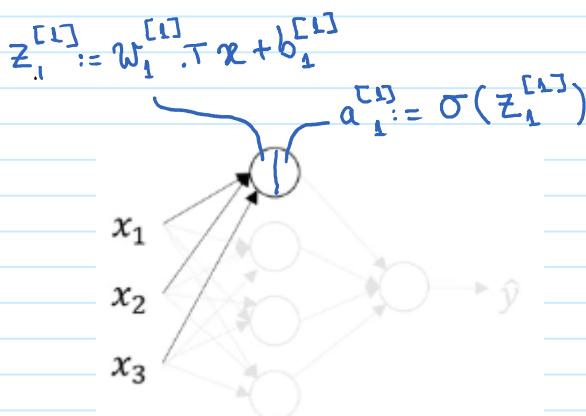
A node in a NN computes two values

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

$$a := \sigma(z)$$

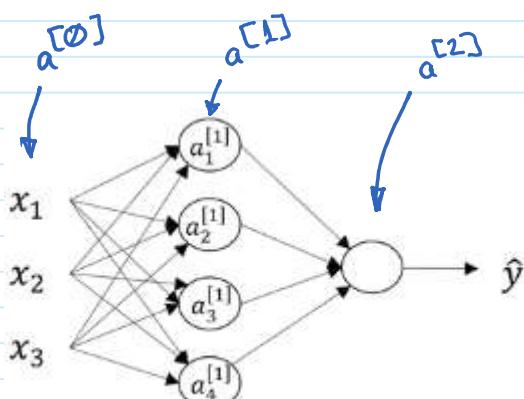
Screen clipping taken: 2018-07-02 14:51

$[l]$ ↪ layer # in NN
 $a^{[l]}$ ↪ node # in layer l



Screen clipping taken: 2018-07-02 15:23

Screen
clipping
taken:
2018-07-02
15:16



Screen clipping taken: 2018-07-02 15:26

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, a_4^{[1]} = \sigma(z_4^{[1]})$$

Screen
clipping
taken:
2018-07-02
15:27

We want to compute \mathbf{z} in vectorized form.

$$\begin{array}{c}
 W^{[1]} \\
 \left[\begin{array}{c} w_1^{[1] \cdot T} \\ w_2^{[1] \cdot T} \\ w_3^{[1] \cdot T} \\ w_4^{[1] \cdot T} \end{array} \right] \\
 (4, 3)
 \end{array}
 \times
 \underbrace{\left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right]}_{(3, 1)} +
 \underbrace{\left[\begin{array}{c} b_1 \\ b_2 \\ b_3 \\ b_4 \end{array} \right]}_{(4, 1)} =
 \left[\begin{array}{c} w_1^{[1] \cdot T} x + b_1^{[1]} \\ w_2^{[1] \cdot T} x + b_2^{[1]} \\ w_3^{[1] \cdot T} x + b_3^{[1]} \\ w_4^{[1] \cdot T} x + b_4^{[1]} \end{array} \right] =
 \underbrace{\left[\begin{array}{c} z_1 \\ z_2 \\ z_3 \\ z_4 \end{array} \right]}_{(4, 1)}$$

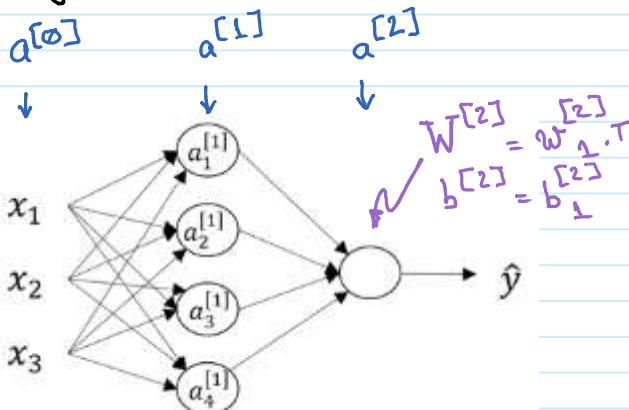
Think about it as having 4 logistic regressions (one per node) for each of the 3 inputs and we stack them together
 → we end up with a $(4, 3)$ matrix

Rule of thumb:
 The vertical dimension corresponds to a layer in a network.

$$a^{[1]} := \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} := \sigma(\mathbf{z}^{[1]})$$

Now if we replace the input layer (i.e. $[x_1, x_2, x_3]$) w/ $a^{[0]}$ and the output layer (i.e. \hat{y}) w/ $a^{[2]}$

we get:



$$\mathbf{z}^{[1]} := W^{[1]} \cdot a^{[0]} + b^{[1]}$$

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

$$\mathbf{z}^{[2]} := W^{[2]} \cdot a^{[1]} + b^{[2]}$$

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

Screen clipping taken: 2018-07-02 15:53

(1, 1) (1, 1)

So to calculate \hat{y} (predicted output) in vectorized form for a single example (i.e. $a^{[0]} = [x_1 \dots x_m]$) we need to calculate the following equations:

Logistic Regression

1 layer \nmid 1 node

$$z := W^T x + b$$

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

Neural Net

2 layers w/ 3 \nmid 1 nodes respectively

$$z^{[1]} := W^{[1]} a^{[0]} + b^{[1]}$$

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

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

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

Computing a NN output for multiple examples

Thursday, 05 July, 2018 23:58

For one example we computed the predicted output as:

$$x \longrightarrow a^{[2]} = \hat{y}$$

For multiple examples we'll calculate the predicted outputs as:

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

in this case
 $a^{[2](i)}$
 \uparrow
layer 2 \nwarrow
example set i

where m is the # of examples used.

A non-vectorized form of this (i.e. using a for loop) is:

for $i = 1$ to m :

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

Note: the same weights
 $W^{[1]}, b^{[1]}$ and
 $W^{[2]}, b^{[2]}$ are
applied to all examples
during a forward pass

THIS IS NOT WHAT
HAPPENS IN A LIVING
ORGANISM!

$$\sum^{[1]} := W^{[1]} \begin{bmatrix} | & | \\ x^{(1)} & \dots & x^{(m)} \\ | & \dots & | \\ & (n_x, m) \end{bmatrix}^T + b^{[1]}$$

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

$$\text{---} \Gamma_2 \Gamma_1 \Gamma_2 \Gamma_1 \Gamma_2 \Gamma_1 \Gamma_2$$

$$H := \cup (\sqcup)$$

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

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

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

value of the activation of the first training example of the first node of layer 1

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

index into each node on layer 1
(scan through hidden units)

Scan through the training examples

Justification for vectorized implementation

$$\begin{aligned} z^{1} &= w^{(1)} x^{(1)} + b^{(1)}, & z^{[1](2)} &= w^{(1)} x^{(2)} + b^{(1)}, & z^{[1](3)} &= w^{(1)} x^{(3)} + b^{(1)} \\ w^{(1)} &= \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix}, & w^{(1)} x^{(1)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, & w^{(1)} x^{(2)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, & w^{(1)} x^{(3)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \\ w^{(1)} \begin{bmatrix} 1 & x^{(1)} & 1 \\ 1 & x^{(2)} & 1 \\ 1 & x^{(3)} & 1 \end{bmatrix} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} z^{1} \\ z^{[1](2)} \\ z^{[1](3)} \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} = Z^{[1]} \end{aligned}$$

$$\hat{y}^{(i)} = \omega^{(i)T} X + b^{(i)}$$
$$\omega^{(i)T} X + b^{(i)} = \hat{y}^{(i)}$$
$$\hat{y}^{(i)} = \begin{bmatrix} \text{green} & \text{blue} & \text{orange} \end{bmatrix} \begin{bmatrix} \text{green} \\ \text{blue} \\ \text{orange} \end{bmatrix} + b^{(i)} = \text{green}^2 + \text{blue}^2 + \text{orange}^2 + b^{(i)}$$

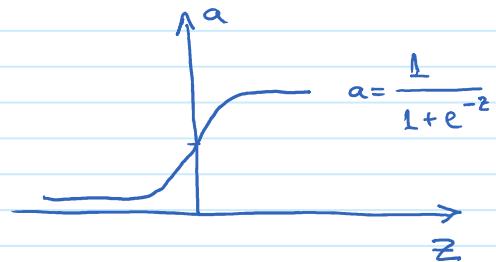
Andrew Ng

Screen clipping taken: 2018-07-02 22:08

Activation Functions

Thursday, 05 July, 2018 23:57

So far we have been using the σ function $\sigma(z) = \frac{1}{1+e^{-z}}$ as the activation function for each node.



Given a more general form of the forward pass is:

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

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

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

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

Where g can be a non-linear function that is not the σ .

If we use $g(z) = \tanh(z)$ as the activation function it almost always works better than the σ

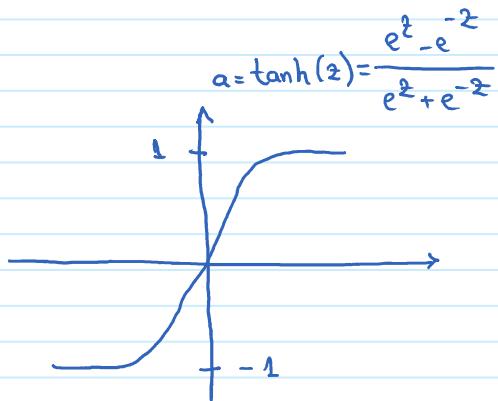
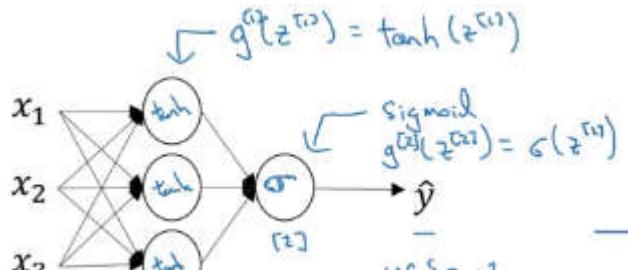
→ that is because the mean activation is closer to 0 as the function takes values between $(-1; 1)$ so it helps "center" your data around 0 .

→ centering makes the learning process for the next layer easier.

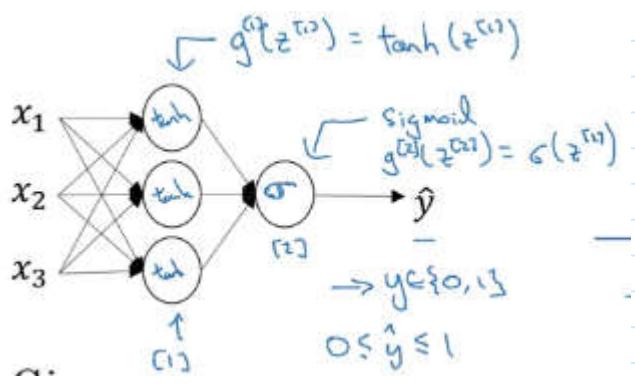
EXCEPTION: On the output layer of a binary classification problem

the σ function is better as the output \hat{y} is between $0 \leq 1$

We can use different activation functions for the different layers of the NN



tanh is a shifted and scaled version of the σ function



Screen clipping taken: 2018-07-02 22:37

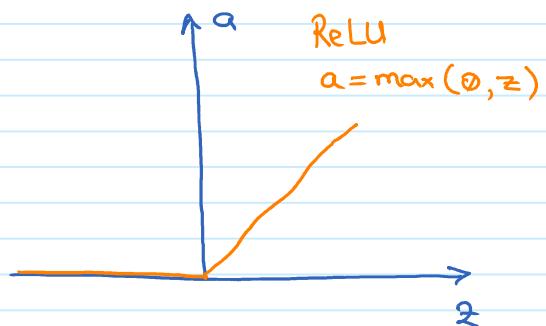
Flat Activation Function Problem

Thursday, 05 July, 2018 23:55

When the derivative of the activation function is close to 0 learning occurs very slowly (slows down gradient descent)

ReLU to the Rescue!

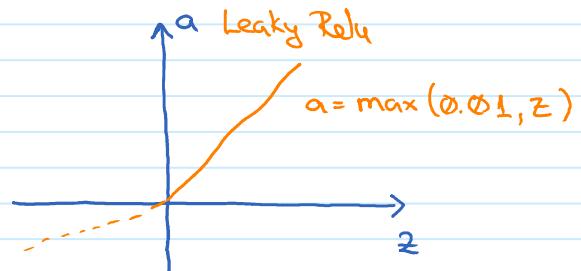
Rectified Linear Unit



Note: ReLU is not differentiable when $z = 0$; but in practice we'll almost never encounter $z = 0.00000...$.

One disadvantage of ReLU is that the derivative is 0 for $z < 0$.

Some use the leaky ReLU, though in practice quite few nodes will have positive z so learning will still occur quite fast.



use 0.01 b/c it empirically seems to do well (but this can be modified)

σ : only use it for the output layer

tanh: try it some times

ReLU: use this most of the time

Evaluating the performance of different activation function in practice instead of always choosing the same "best" one ensures that data idiosyncrasies are adequately

dealt with and also hardens the architecture against evolutions of the algorithm.

Also evaluate: network topology, initialization values, etc

Why non-linear activation functions?

Thursday, 05 July, 2018 23:20

Why not set $g(z) = z$ (the linear activation function)? ?

In this case:

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

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

$$= w^{[2]}(w^{[1]}x + b^{[1]}) + b^{[2]}$$

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

$$= w'x + b'$$

So it turns out we are computing a linear function of the inputs

→ there's a theorem showing that if we are computing a linear func. of the inputs a network with many layers can be replaced by a network w/ only a single layer so the hidden layers don't do anything (they are not more expressive than a linear function of the inputs).

ALL B/C the composition of a linear func is itself a linear func.

EXCEPTION: If we do machine learning on a regression problem (i.e. housing prices) so $y \in \mathbb{R}$ i.e. $y = \$0 \dots \$1,000,000$ the output layer should use a linear activation function (so that the predicted output

can go from ϕ to ∞). But the hidden layers should still NOT USE
a linear activation function.

Derivatives of Activation Functions

Wednesday, 04 July, 2018 22:51

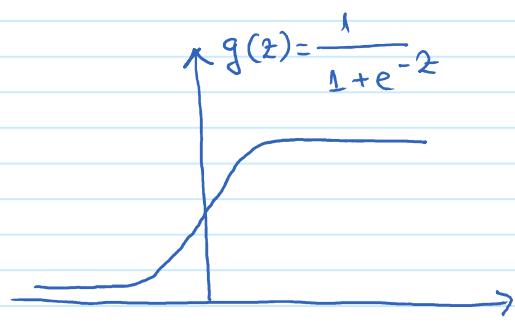
- Sigmoid activation function

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

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

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

$$\boxed{g'(z) = a (1 - a)}$$



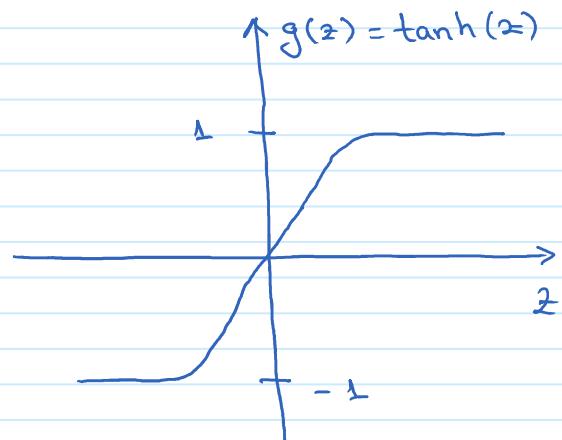
- Tanh activation function

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

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

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

$$\boxed{g'(z) = 1 - a^2}$$



- ReLU & Leaky ReLU

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

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases}$$

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

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0.01, & z \leq 0 \end{cases}$$

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

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

We can set $g'(z) = 1$ for $z = 0$

↳ g' becomes a subgradient of the optimization func.

Gradient descent NN

Thursday, 05 July, 2018 22:41

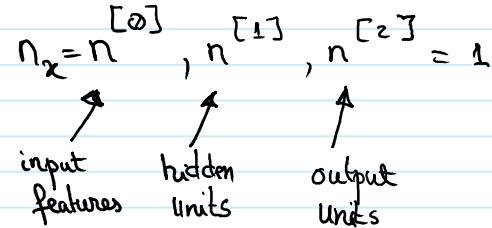
Goal of this section is to determine what parameter values to use when running gradient descent

Gradient descent consists of repeating 3 stages:

①) INITIALIZE PARAMETERS

- REPEAT
- 1) FORWARD PROPAGATION (compute predictions)
 - 2) BACKWARD PROPAGATION (compute derivatives)
 - 3) UPDATE WEIGHTS

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



Cost function (assume binary classification):

$$J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum \mathcal{L}(\hat{y}, y)$$

$\uparrow a^{[2]}$

To train the NN we need to do Gradient Descent:

Important to initialize parameters randomly (instead of all 0's)

After we initialize each loop of Gradient Descent

Repeat {

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

* Compute derivatives

$$\frac{dW^{[1]}}{dW^{[1]}} = \frac{dJ}{dW^{[1]}}, \frac{db^{[1]}}{db^{[1]}} = \frac{dJ}{db^{[1]}}, \dots$$

* Update weights:

$$W^{[L]} = W^{[L]} - \alpha dW^{[L]}$$

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

} until it looks like parameters are converging.

FORWARD PROPAGATION (compute predictions):

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

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

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

$$A^{[2]} = g(Z^{[2]}) \text{ for binary classification } A^{[2]} = \sigma(Z^{[2]})$$

BACK PROPAGATION (compute derivatives) :

$$dZ^{[2]} = A^{[2]} - Y \quad \text{this is all vectorized across examples}$$

↑ assume we are doing
binary classification so
 $Y = [y^{(1)}, y^{(2)} \dots y^{(m)}]$

the activation function for the
output layer is $\sigma(z) \rightarrow g'(z) = \sigma(z) - Y$

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

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

↑ sum horizontally across
matrix

↑ prevents Python from
outputting $(n^{[2]},)$ instead
it outputs $(n^{[2]}, 1)$

$$dZ^{[L]} = \underbrace{W^{[L] \top} dZ^{[L+1]}}_{(n^{[L]}, m)} * \underbrace{g^{[L]'}(Z^{[L]})}_{(n^{[L]}, m)}$$

element wise product

↑ the derivative of whatever activation
function we used for the hidden layer

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

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

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

$(n^{[L]}, 1)$

↑
we don't want
a $(n^{[L]},)$ output

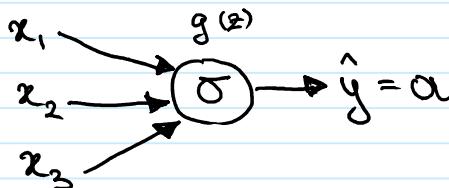
Backpropagation Intuition

Thursday, 05 July, 2018 22:55

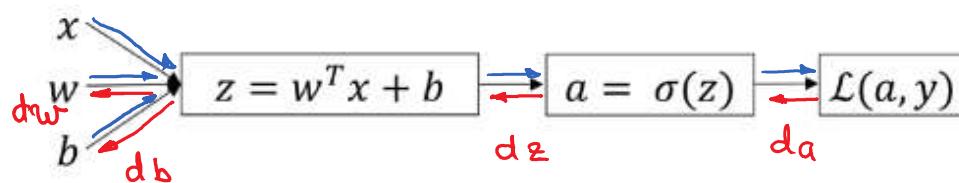
Goal: Develop an intuition of how the parameters used in backpropagation are derived.

We'll start from the (single node, single layer) Logistic Regression and progress to the (multi node, multilayer) Neural Net.

• Logistic REGRESSION



Logistic regression



→ Forward Prop

← Back Prop

Screen clipping taken: 2018-07-05 22:38

$$da = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$d2 = a - y$$

$$dw = d2 \cdot x$$

$$db = d2$$

See below for derivations.

NOTE: We are applying a lot of chain rules (see [COMPUTATION GRAPH](#) in Week 2)

$$\mathcal{L}(a, y) = -y \log(a) - (1-y) \log(1-a)$$

↓ take derivative of
 $\mathcal{L}(a, y)$ wrt a

$$da = \frac{d}{da} \mathcal{L}(a, y)$$

$$|| da = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\parallel \text{da} = -\frac{y}{a} + \frac{1-y}{1-a}$$

Then:

$$dz = \frac{d}{dz} \mathcal{L}(a, y)$$

(apply chain rule)

$$= \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{d}{dz} a = \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{d}{dz} g(z) = \frac{\partial \mathcal{L}}{\partial a} \cdot g'(z) = da \cdot g'(z)$$

for logistic regression

$$a = g(z) = \sigma(z) \Rightarrow g'(z) = a(1-a)$$

as per [DERIVATIVES OF ACTIVATION FUNCTIONS](#)

So:

$$dz = da \cdot g'(z)$$

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

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

$$= -y(1-a) + a(1-y)$$

$$= -y + ya + a - ay$$

$$\parallel dz = a - y$$

Then we compute

$$dw = \frac{d}{dw} \mathcal{L}(a, y)$$

$$= \underbrace{\frac{d}{da} \mathcal{L}(a, y)}_{a-y} \cdot \frac{d}{dz} a \cdot \frac{d}{dw} z$$

$$= (a - y) \cdot \frac{d}{dw} z$$

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

$$\Rightarrow \frac{d}{dw} z = x \quad \left\{ \begin{array}{l} \frac{d}{db} z = 1 \end{array} \right.$$

$$\| \quad d\mathbf{w} = (\mathbf{a} - \mathbf{y}) \cdot \mathbf{x} \text{ or } d\mathbf{z} \cdot \mathbf{x}$$

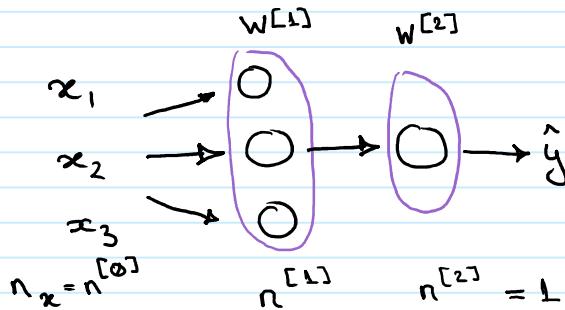
Similarly to $d\mathbf{w}$

$$\| \quad d\mathbf{b} = d\mathbf{z}$$

2018-07-07 (Sat) @ 10:30

- NEURAL NETWORK 2 LAYER:

Equations for a single example:

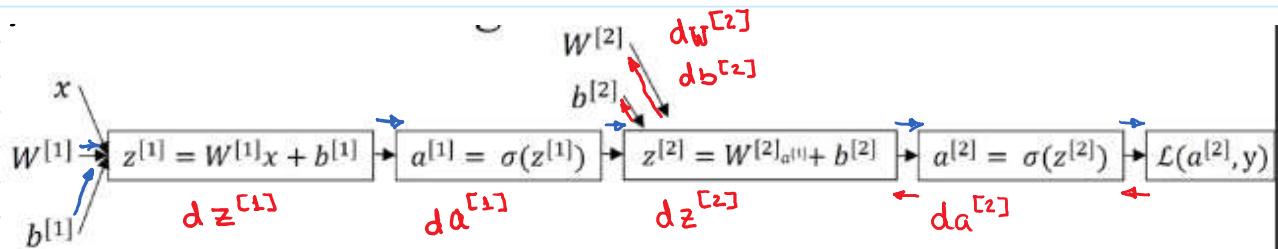


$$W^{[2]} : (n^{[2]}, n^{[1]})$$

$$z^{[2]}, dz^{[2]} : (n^{[2]}, 1)$$

$$z^{[1]}, dz^{[1]} : (n^{[1]}, 1)$$

Instead of having a single computational step (for a single layer) we have **two steps** one for each layer.



Screen clipping taken: 2018-07-07 10:55

Note: We don't need to take derivatives for the input \mathbf{x} b/c for supervised learning the input \mathbf{x} is fixed so we are not trying to optimize \mathbf{x} .

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

$$\| \quad \begin{cases} d\mathbf{w}^{[2]} = dz^{[2]} a^{[1] \top} \\ \dots \end{cases}$$

In Linear Regression " $d\mathbf{w} = dz \cdot \mathbf{x}$ "

In NN $a^{[1]}$ plays the role of \mathbf{x} , and it's transposed b/c \mathbf{W} is a row vector, whereas $d\mathbf{w}^{[2]}$ is a column vector.

$$W^{[2]} = \boxed{\quad \quad \quad}$$

$$\| \begin{cases} \alpha w = \alpha z - a \\ \alpha b^{[2]} = d z^{[2]} \end{cases}$$

b/c w is a row vector, whereas αw is a column vector.
 $w = [\quad]$
and x was a column vector.

$$\| \underbrace{d z^{[1]}}_{(n^{[1]}, 1)} = \underbrace{(W^{[2]}, T)}_{(n^{[2]}, n^{[2]})} d z^{[2]} * \underbrace{g^{[1]}'(z^{[1]})}_{(n^{[1]}, 1)}$$

element wise multiplication

★ Tip: Ensuring the dimensions match up resolves many bugs in implementing back propagation

$$\| \begin{cases} d w^{[L]} = d z^{[L]} \cdot X^T \\ d b^{[L]} = d z^{[L]} \end{cases}$$

Vectorized equations for multiple examples:

SINGLE EXAMPLE

$$d z^{[2]} = a^{[2]} - y$$

$$d W^{[2]} = d z^{[2]} a^{[1]T}$$

$$d b^{[2]} = d z^{[2]}$$

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

$$d W^{[1]} = d z^{[1]} x^T$$

$$d b^{[1]} = d z^{[1]}$$

VECTORIZED IMPLEMENTATION

$$d z^{[2]} = A^{[2]} - Y$$

$$d W^{[2]} = \frac{1}{m} d z^{[2]} A^{[1]T}$$

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

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

element wise multiplication

$$d W^{[1]} = \frac{1}{m} d z^{[1]} X^T$$

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

Screen clipping taken: 2018-07-07 11:24

For single example we have :

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

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

To vectorize this for multiple examples we stack the z in columns (one col. per example)

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

and then we have

$$\mathbf{Z}^{[l]} = \mathbf{W}^{[l]} \mathbf{X} + \mathbf{b}$$

$$A^{[l]} = g^{[l]}(\mathbf{Z}^{[l]})$$

Random initialization of weights

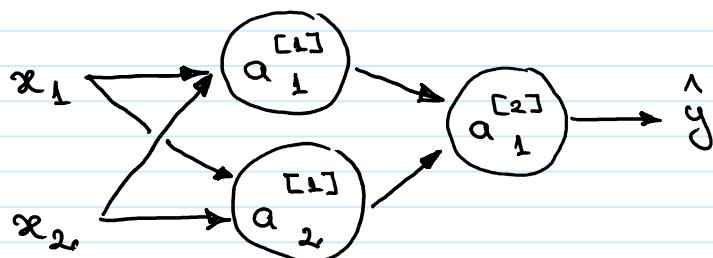
Saturday, 07 July, 2018 11:41

For NN it's important to initialize the weights randomly.

→ for logistic regression it is ok to initialize the weights to \emptyset

→ for NN this won't work w/ gradient descent.

WHAT HAPPENS WHEN INITIALIZING WEIGHTS TO ZERO?



Note: initializing b to \emptyset is ok
but $W = \emptyset$ is a problem

$$n^{[0]} = 2 \quad n^{[1]} = 2 \quad n^{[2]} = 1$$

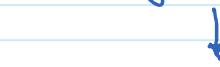
$$W^{[1]} = \begin{bmatrix} \emptyset & \emptyset \\ \emptyset & \emptyset \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} \emptyset \\ \emptyset \end{bmatrix} \quad W^{[2]} = \begin{bmatrix} \emptyset & \emptyset \end{bmatrix}$$

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



these units are identical

(aka symmetric → they compute
exactly the same function)



by induction: after every iteration →
of training they still compute
the exact same function

(ie both hidden units compute the
same function → so there is no point
of having more than one hidden

Proof:

$$\text{Given } dW = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

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

$$W^{[1]} = \begin{bmatrix} \underline{\underline{u}} \\ \underline{\underline{v}} \end{bmatrix}$$

(∴ at 1st ... and all remaining)

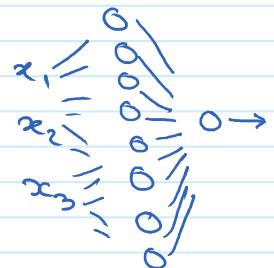
same function - or there is no power

of having more than one hidden
unit)

$$W = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

(i.e. the 1st row = 2nd row so still symmetric)

This argument works the same for networks w/
many hidden units \rightarrow if they all are initialized to the
same value (i.e. 0) they all compute the same function
and can be replaced by a single unit.



RANDOM INITIALIZATION

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

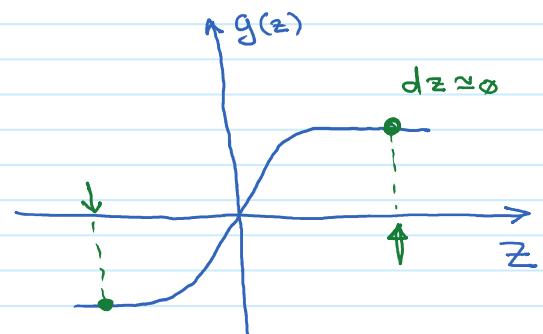
$$b^{[1]} = \text{np.zeros}(2, 1)$$

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

$$b^{[2]} = \emptyset$$

Why 0.01? Why do we want to initialize the weights to small values?
Why not 100 or 1000?

When we use a sigmoid output function
(tanh or σ) if the weights are
too large z is very large and we
end up in the regions where the
function is flat so the derivative ≈ 0
so gradient descent is slow and the
learning is very slow



If we don't use sigmoids this is not as big of an issue.

Sometimes $\phi(\cdot)$ is not ideal → for shallow networks it is OK
(not many hidden units)
→ for deep networks there is better

Week 4 Deep NNs

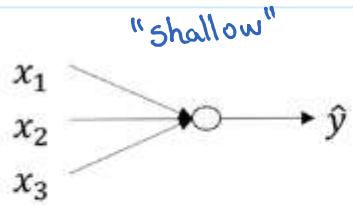
Monday, 09 July, 2018 20:55

DEEP L-LAYER NEURAL NETWORKS

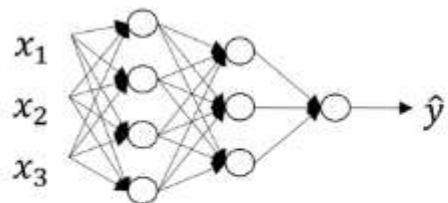
Deep NN nomenclature

Monday, 09 July, 2018 22:21

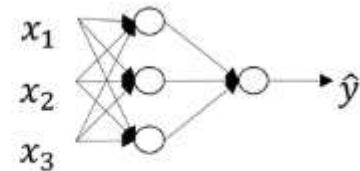
Nomenclature of Deep neural network



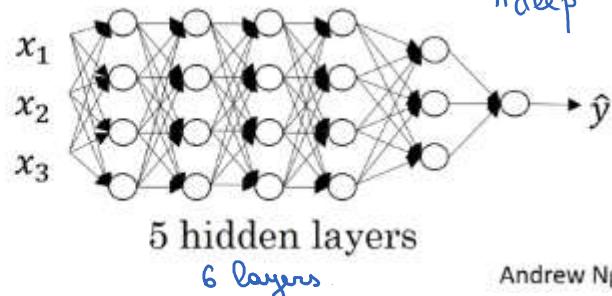
logistic regression
1 layer



2 hidden layers
3 layers

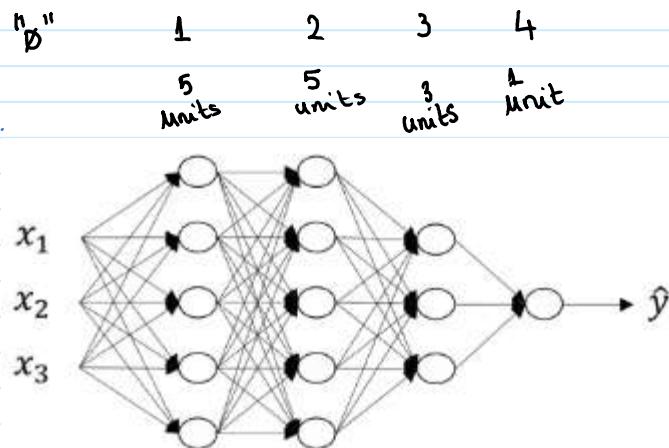


1 hidden layer
2 layers



Andrew Ng

Screen clipping taken: 2018-07-09 22:08



4 layer NN $\Rightarrow L = 4$

Screen
clipping
taken:
2018-07-09
22:11

L : # of layers in the NN

$a^{[l]}$: activations in layer l
 $a^{[l]} = g^{[l]}(z^{[l]})$

$n^{[l]}$: # of units (nodes) in a layer l

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

$$n^{[1]} = 5$$

$$n^{[2]} = 5$$

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

$$n^{[4]} = 1 = n^{[L]}$$

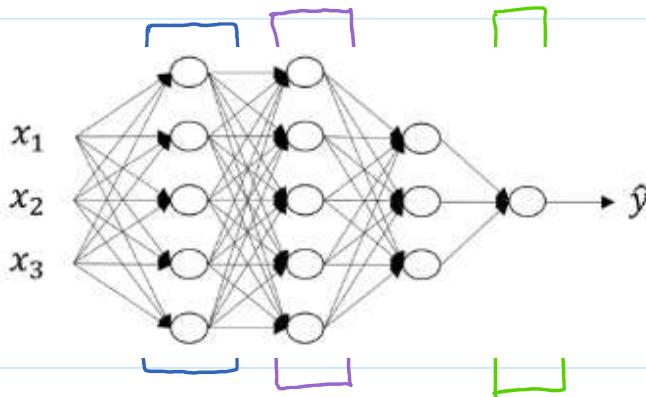
$w^{[l]}$ } weights used to compute
 $b^{[l]}$ } $z^{[l]}$ in layer l

$X = a^{[\phi]}$ (the input features are
the activations for
layer ϕ)

$$a^{[L]} = a^{[4]} = \hat{y}$$

Forward Prop in deep NNs

Monday, 09 July, 2018 22:23



Screen
clipping
taken:
2018-07-09
22:24

Example of FWD Prop on a single
example x

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

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

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

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

• • •

$$z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}$$

$$a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$$

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

General Forward Propagation Equation:

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

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

Vectorized form (for multiple
training examples):

$$\underline{z}^{[1]} := W^{[1]} A^{[0]} + b^{[1]}$$

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

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

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

$$A^{[4]} := g^{[4]}(\underline{z}^{[4]}) = \hat{Y}$$

Layer L

for l=1 to 4

there is no
way to get
rid of
(vectorize
away)

the for loop
to iterate
through
layers to
calculate
activations

Layer 4

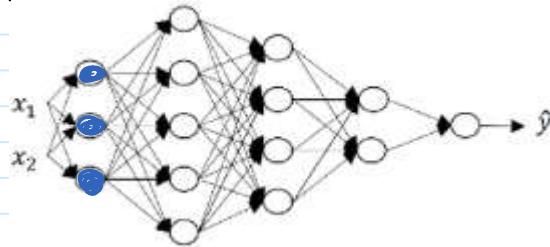
prediction for all
training examples stacked
horizontally

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

Getting matrix dimensions right

Monday, 09 July, 2018 22:49

$$L = 5$$



Screen
clipping
taken:
2018-07-09
22:52

$$n^{[0]} = 2 = n_x$$

$$n^{[1]} = 3$$

$$n^{[2]} = 5$$

$$n^{[3]} = 4$$

$$n^{[4]} = 2$$

$$n^{[5]} = 1$$

for now we've seen
only NN w/ a single
output unit. In later
courses we'll see NN
w/ multiple output
units.

FOR SINGLE EXAMPLE IMPLEMENTATION:

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

↑
vector of
activations
for the 1st
hidden layer

$$(3, 1)$$

$$(n^{[1]}, 1)$$

↑
vector of
input features

$$(2, 1)$$

$$(n^{[0]}, 1)$$

In General:

$$W^{[l]} : (n^{[l]}, n^{[l-1]})$$

$$\text{so for } W^{[2]} : (n^{[2]}, n^{[1]}) = (5, 3)$$

so for W 's dimensions:

$$\begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \ddots & \cdot \\ \cdot & \cdot & \ddots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$$\underbrace{(3, 2)}$$

$$|W| = (n^{[1]}, n^{[0]})$$

b/c we are going to compute

$$\underbrace{z^{[2]} = W^{[2]} \underbrace{a^{[1]}}_{(5, 1)} + b^{[2]}}$$

$$\underbrace{\quad}_{(5, 1)} \quad \underbrace{\quad}_{(3, 1)} \quad \underbrace{\quad}_{(5, 3)}$$

$$W^{[3]} : (4, 5)$$

$$W^{[4]} : (2, 4)$$

$$W^{[5]} : (1, 2)$$

For the dimensions of b remember b is the same dimension as z

In General

$$b^{[e]} : (n^{[e]}, 1)$$

$$b^{[1]} : (3, 1)$$

$$b^{[2]} : (5, 1)$$

$$|dw^{[e]}| = |W^{[e]}| = (n^{[e]}, n^{[e-1]})$$

$$|db^{[e]}| = |b^{[e]}| = (n^{[e]}, 1)$$

$$|a^{[e]}| = |z^{[e]}| = (n^{[e]}, 1)$$

FOR MULTIPLE EXAMPLE IMPLEMENTATION (VECTORIZED)

The dimensions of $\begin{cases} W, dw \\ b, db \end{cases}$ stay the same as for a single example

The dimensions of $z, a \nparallel x$ will change a little.

$$|z^{[e]}| : (n^{[e]}, 1)$$

$$|\Sigma^{[e]}| : (n^{[e]}, m) = |d\Sigma^{[e]}|$$

$$|x| = |a^{[0]}| : (n_x, 1)$$

$$|\underline{X}| : (n_x, m) \leftrightarrow \begin{matrix} \text{entire} \\ \text{training set} \end{matrix}$$

$$|a^{[e]}| : (n^{[e]}, 1)$$

special case $e = 0$

$$|a^{[l]}| : (n^{[l]}, 1)$$

$$|A^{[l]}| : (n^{[l]}, m) = |dA^{[l]}|$$

†

special case $l=0$

$$|A^{[0]}| = |X| = (n^{[0]}, m) \\ = (n_x, m)$$

Note for $|b^{[l]}|$

$$\underbrace{Z^{[l]}}_{(n^{[l]}, m)} = \underbrace{W^{[l]} \cdot A^{[l-1]}}_{(n^{[l]}, n^{[l-1]})} + \underbrace{b^{[l]}}_{(n^{[l]}, 1)}$$

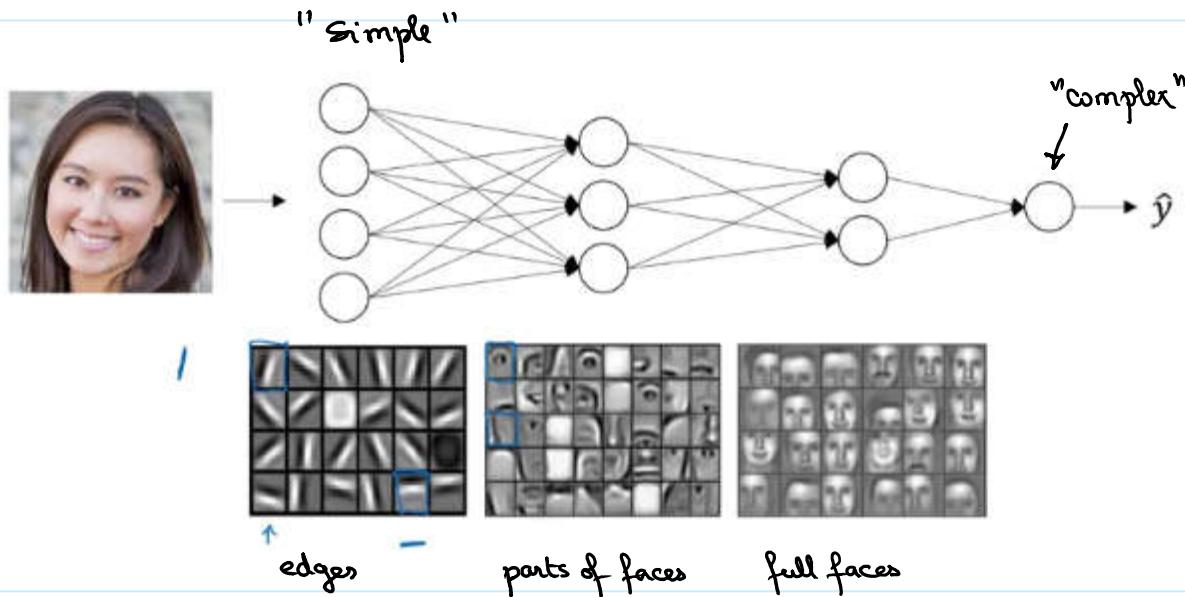
but when the addition occurs $b^{[l]}$ is duplicated
m times by python broadcasting
and added to the $W^{[l]} A^{[l-1]}$
product.

$$\text{Also } |dZ^{[l]}| = |Z^{[l]}| = |A^{[l]}| = |dA^{[l]}| = (n^{[l]}, m)$$

Why deep representations?

Tuesday, 10 July, 2018 22:09

Why do deep NN work well at solving many problems?



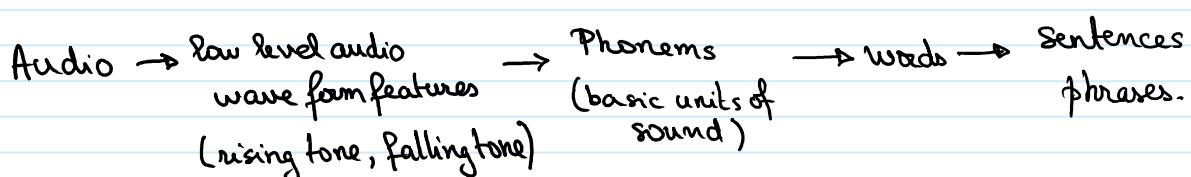
Screen clipping taken: 2018-07-10 22:12

First Layer detects simple patterns (edges)

Second Layer detects slightly more complex patterns (parts of faces)

Third Layer detects highly complex patterns (faces)

A NN composes simple patterns and builds them up in a hierarchical representation.



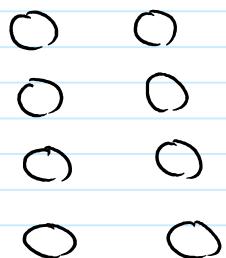
CIRCUIT THEORY & DEEP LEARNING

Circuit theory: describes what kind of functions we can compute using logic gates (AND, NOT, OR ...)

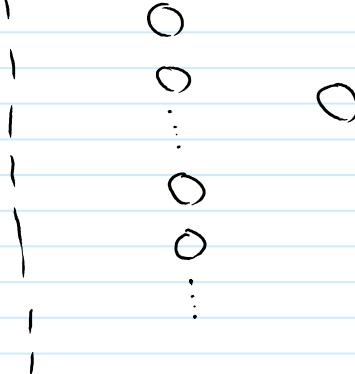
An informal argument as to why "deep" neural nets are good @ solving problems:

There are functions that can be computed with a "small" L-layer deep neural net that shallower networks require exponentially more hidden nodes to compute.

DEEP NN (14 nodes on 5 layers)



SHALLOW NN (200 nodes on 2 layers)



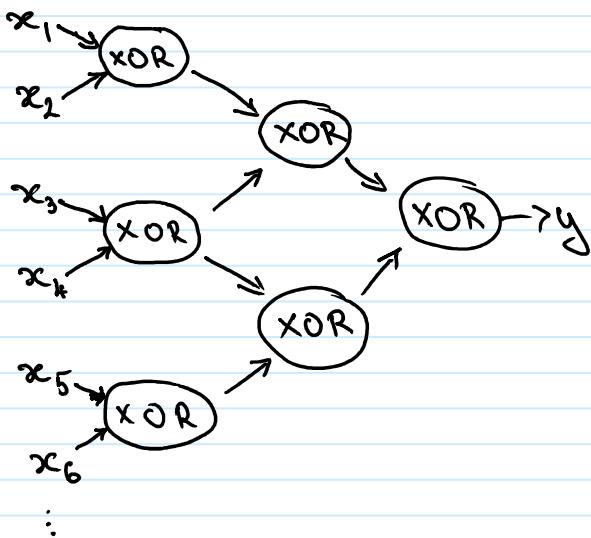
Example:

Compute the XOR of many inputs

WHEN MANY HIDDEN LAYERS

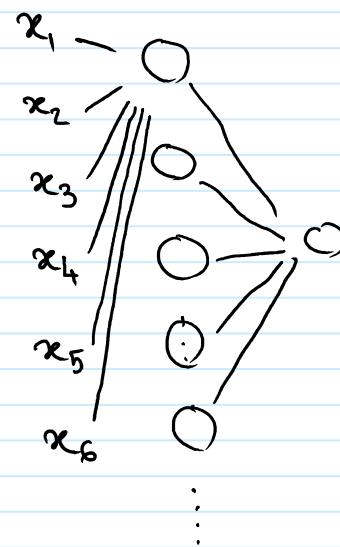
ARE ALLOWED

← O(log n) layers deep →



WHEN FEW HIDDEN UNITS ARE

ALLOWED



The hidden layer needs to be exhaustively large

The hidden layer needs to be exhaustively large to accommodate all the possible combinations $O(2^n) \rightarrow$ exponentially large.

Very deep NN (with maybe dozens of hidden layers) seem to work on some hard problems.

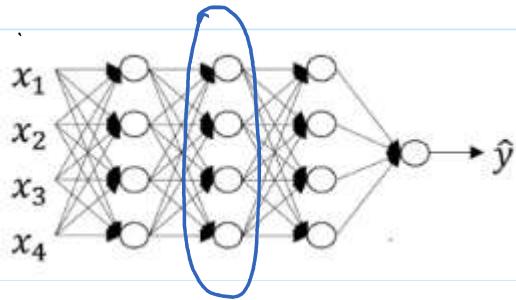
Building blocks of deep NN

Tuesday, 10 July, 2018 22:53

Key Components

Forward Propagation

Backward Propagation



Screen clipping taken: 2018-07-10 22:57

For a given layer l : $W^{[l]}, b^{[l]}$

FWD PROP:

Input: $a^{[l-1]}$

Output: $a^{[l]}$, cache $z^{[l]}$ (we'll use this during Back Prop)

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

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

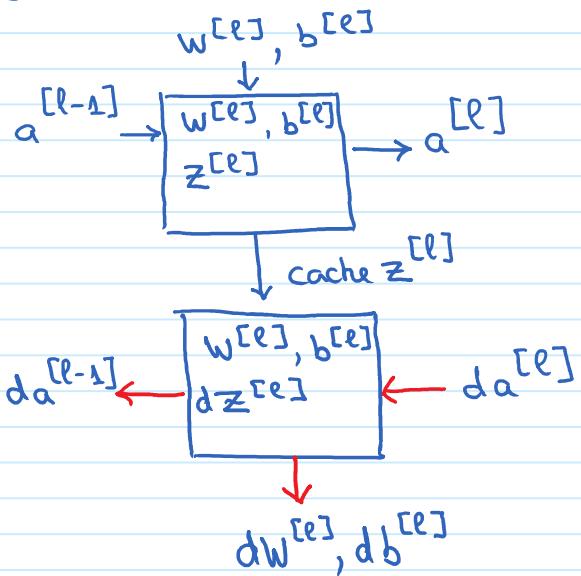
BACK PROP:

Input: $da^{[l]}$, cache $z^{[l]}$

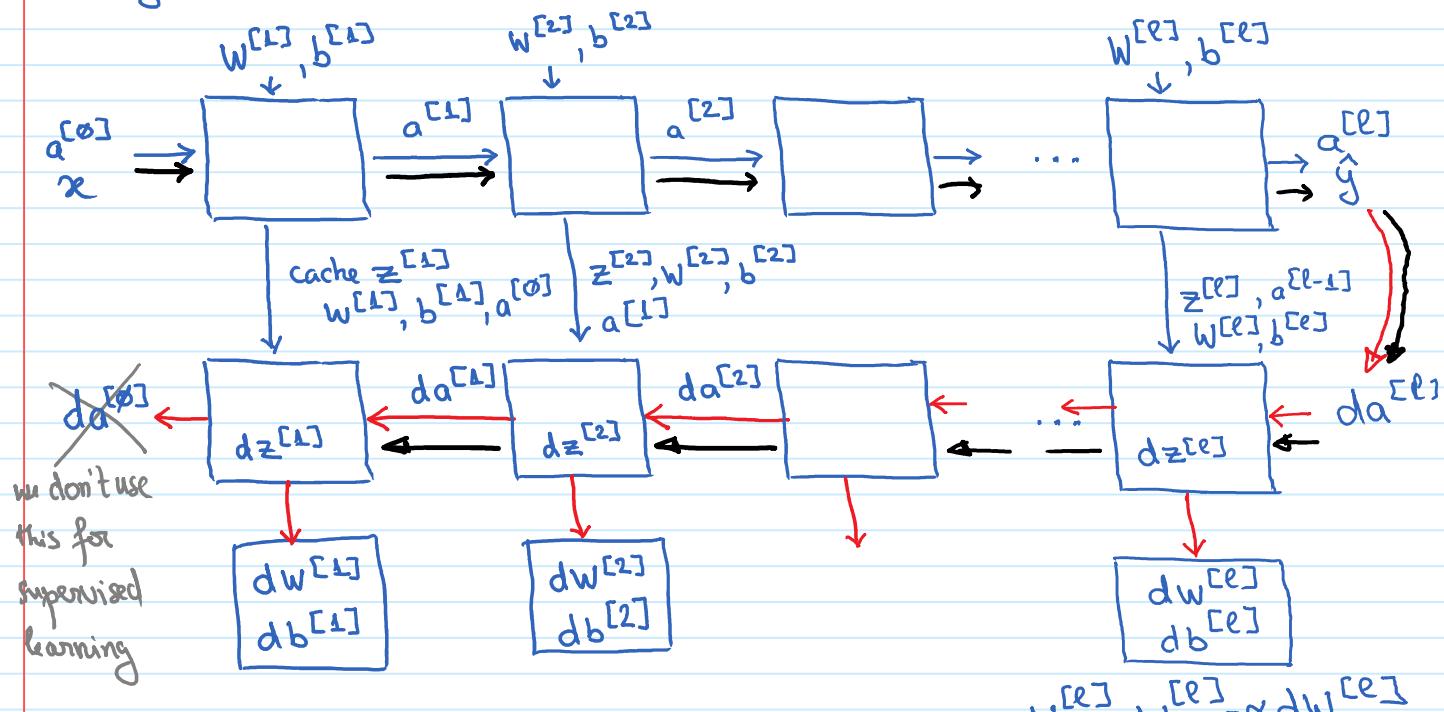
Output: $da^{[l-1]}$,
 $dW^{[l]}, db^{[l]}$

For layer l :

$$\dots \quad \downarrow \quad w^{[l]}, b^{[l]}$$



For a general NN:



We cache $z^e, w^e, b^e, a^{[l-1]}$

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

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

FWD & BACK prop

Tuesday, 10 July, 2018 23:21

FORWARD PROPAGATION (for layer l):

- **Inputs:** $a^{[l-1]}$ ($w^{[l]}, b^{[l]}$ saved in cache)
- **Outputs:** $a^{[l]}$ ($z^{[l]}$ saved in cache)

SINGLE EXAMPLE	VECTORIZED
$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$ $a^{[l]} = g^{[l]}(z^{[l]})$	$\underline{z}^{[l]} = \underline{w}^{[l]} \underline{A}^{[l-1]} + \underline{b}^{[l]}$ $A^{[l]} = g^{[l]}(\underline{z}^{[l]})$

$$X = A^{[\emptyset]} \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow A^{[L]} = \hat{y}$$

BACWARD PROPAGATION (for layer l)

- **Inputs:** $da^{[l]}$
- **Outputs:** $da^{[l-1]}, dw^{[l]}, db^{[l]}$

SINGLE EXAMPLE

$$dz^{[e]} = da^{[e]} * g^{[e]}'(z^{[e]})$$

element wise multiplication

$$dw^{[e]} = dz^{[e]} \cdot a^{[e-1]}$$

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

$$da^{[e-1]} = W^{[e]T} \cdot dz^{[e]}$$

If we plug the definition of $da^{[e-1]}$ into $da^{[e]}$ and compute $dz^{[e]}$ we get:

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

See [Backprop Intuition](#) in Week 3 for how to understand the derivation of these formulae.

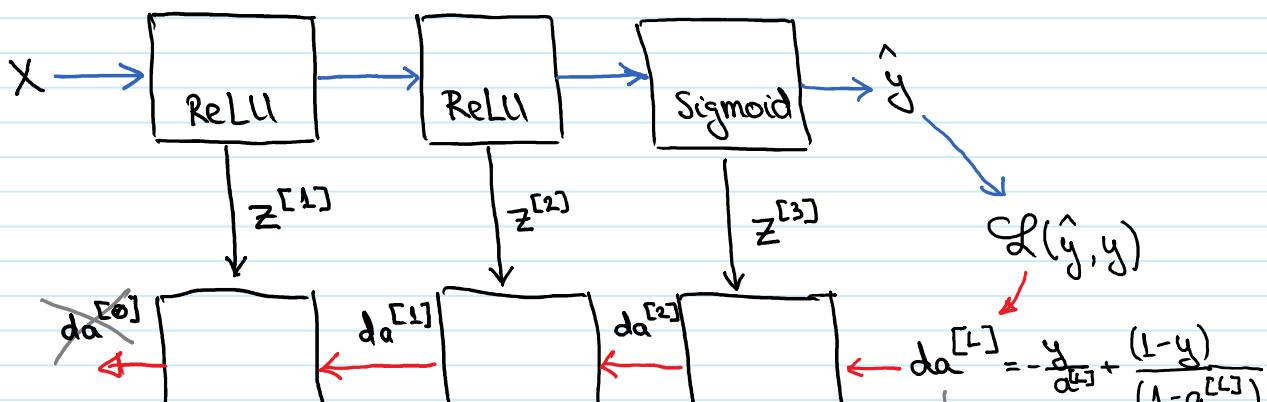
VECTORIZED

$$dZ^{[e]} = dA^{[e]} * g^{[e]}'(Z^{[e]})$$

$$dW^{[e]} = \frac{1}{m} dZ^{[e]} \cdot A^{[e-1]T}$$

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

$$dA^{[e-1]} = W^{[e]T} dZ^{[e]}$$



$$da^{[L]} = -\frac{y}{a^{[L]}} + \frac{(1-y)}{(1-a^{[L]})}$$

When we do binary classification using logistic regression it can be shown:

$$da^{[L]} = \frac{d}{d\hat{y}} L(\hat{y}, y) = \frac{d}{da^{[L]}} L(a^{[L]}, y)$$

$$= -\frac{y}{a^{[L]}} + \frac{1-y}{1-a^{[L]}}$$

$$= -\frac{y}{a^{[L]}} + \frac{1-y}{1-a^{[L]}}$$

When we VECTORIZE the scheme above we get :

$$dA^{[L]} = \left[-\frac{y^{(1)}}{a^{[L]}} + \frac{1-y^{(1)}}{1-a^{[L]}} \quad \dots \quad -\frac{y^{(m)}}{a^{[L]}} + \frac{1-y^{(m)}}{1-a^{[L]}} \right]$$

\uparrow
the activations for all training examples

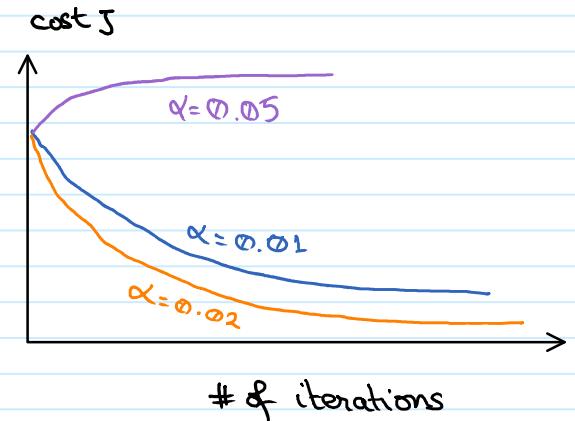
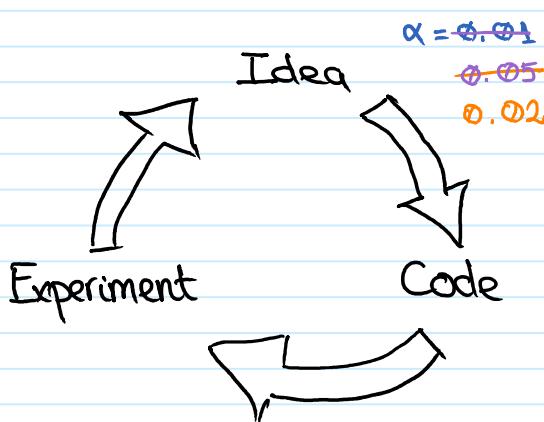
Parameters vs Hyperparameters

Wednesday, 11 July, 2018 21:38

Parameters : $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]} \dots$

Hyperparameters : {
Learning rate α
of iterations
of hidden layers L
of hidden units $n^{[1]}, n^{[2]} \dots$
choice of activation function
:
Momentum term
Minibatch size
Regularization parameters
LATER COURSES}

APPLYING DEEP LEARNING IS A VERY EMPIRICAL PROCESS :
(aka need to try different things and see what works)



Deep Learning is applied to many different domains :

Computer vision, speech, natural language processing, online advertising,
(NLP)

web search, product recommendations etc.

Intuitions from one domain sometimes carry over to other domains sometimes they don't.

→ Just try out a range of values and see what works.

Even when working in one domain hyperparameters best values change as the hardware improves or has changed, data sets have changed.

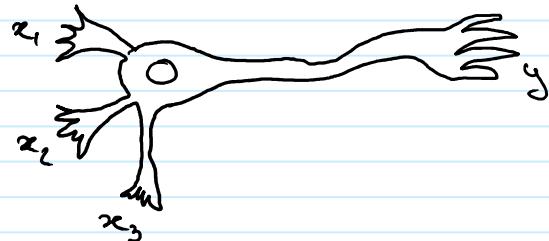
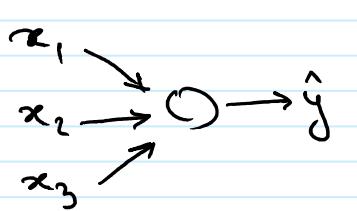
HEURISTIC: Every now and then try different hyperparameters.

Deep Learning and the human brain

Wednesday, 11 July, 2018 21:56

Punch Line: Not very much in common!

There is a very loose analogy between a logistic regression unit and a biological neuron.



We don't know how a neuron works!