



- ✓ • Geoffrey Hinton (backprop, NN)
- ✓ • Pieter Abbeel (reinforcement learning)
- ✓ • Ian Goodfellow (GANs)

- Yoshua Bengio
- Yuxiang Lin

- Andrej Karpathy
- Ruslan Salakhutdinov

- Yann LeCun

# Deep Learning : Week 1

## Intro to Deep Learning

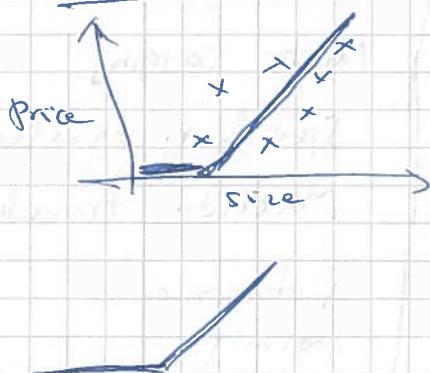
Intro:

- ① AI is new Electricity
- ② end-to-end learning
- ③ train / dev / test
- ④ CNN
- ⑤ RNN / LSTM

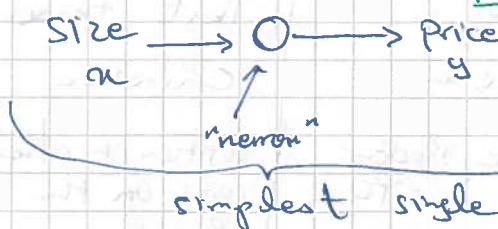
1. Neuronal Networks and Deep Learning
2. Improving Deep Neuronal Networks: Hyperparameter tuning, Regularization, Optimization
3. Structuring your ML project
4. CNN
5. Natural Language Processing Building sequence models



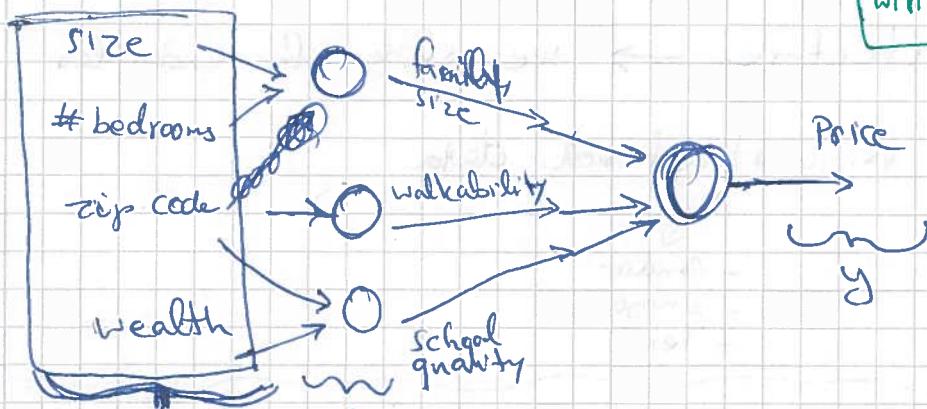
What is a neuronal Network?



$$\text{Price} = f(\text{size})$$



ReLU : Rectified Linear Unit



take all four input features (Densely connected)

### Learning Objectives

- ① Understand the major trends driving the rise of Deep learning.
- ② Be able to explain how deep learning is applied to supervised learning
- ③ Understand what are the major categories of models (such as CNN and RNN) and where they should be applied
- ④ Be able to recognize when Deep learning will and will not work well.

- NN are best for the supervised learning settings
  - finding a function that maps an input  $x$  to an output  $y$ .  $(x, y)$   
 $y = f(x)$

## Supervised Learning with NNs

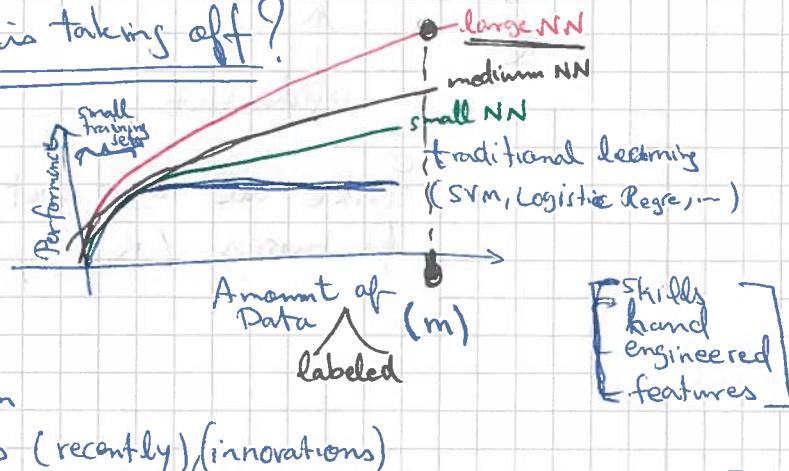
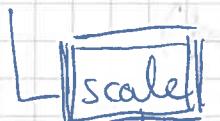
	Input ( $x$ )	Output ( $y$ )	Application
Standard NN	Home features Ad, user info	Price click on ad? (0/1)	Real estate Online advertisement
CNN	Image	Object (1, ..., 1000)	Photo tagging
RNN	Audio	Text transcript	Speech recognition
Custom/ Hybrid	English Image, Radar INFO	Chinese Position of other Cars on the Road	Machine translation Autonomous Driving

- Cleverly selecting what should be  $x$  and what should be  $y$
- NN architecture → see slides for examples
- Structured vs. unstructured data

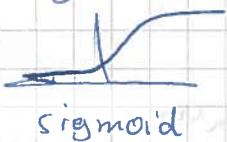
- ↓                                  ↓
- Home info
  - Ad info
  - Audio
  - Image
  - Text

## Why deep learning is taking off?

- Drivers of DL :

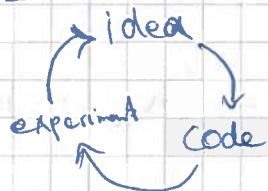


- Algorithm improvements  
e.g.



made computing very much faster

- faster computation



10 min  
1 day  
1 month

} try a lot more ideas

## Week 2

## Neural Networks Basics

### Binary Classification:

- forward pass / backward pass
- logistic regression

An algorithm for binary classification

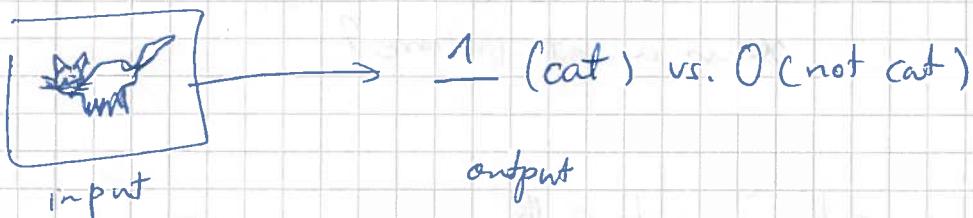


Image is  $64 \times 64$  matrices of RGB

→ "unroll" pixel values into a vector  $\mathbf{x} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$  width of  
size  $64 \times 64 \times 3 = 12288$

$$n = n_{\text{in}} = 12288$$

Notation:

$(\mathbf{x}, y)$  single training example

$$\mathbf{x} \in \mathbb{R}^{n_{\text{in}}}, y \in \{0, 1\}$$

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

### Learning Objectives:

- ① Build a logistic Regression mode as a shallow NN
- ② Implement the main steps of an ML algorithm, including making predictions, derivative computation, and Gradient descent.
- ③ Implement computationally efficient, highly vectorized version of a model
- ④ Understand how to compute derivatives for logistic regression using Backprop.
- ⑤ Become familiar with Python & numpy
- ⑥ work with iPython notebooks
- ⑦ Implement vectorization across multiple training examples

←

$m_{\text{train}}$        $m_{\text{test}} = \# \text{ test examples}$

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \\ E & & & m \end{bmatrix} \quad X \in \mathbb{R}^{n_a \times m}$$

~~X~~  
X. shape =  $(n_a, m)$

$$Y = [y^{(1)} \ y^{(2)} \ \cdots \ y^{(m)}] \quad Y \in \mathbb{R}^{1 \times m}$$

Y. shape =  $(1, m)$

---

## Logistic Regression

in a supervised learning problem

A learning algorithm when output labels  $Y$  are either 0 or 1

Given  $x$ , want  $\hat{y} = P(y=1|x)$ ,  $0 \leq \hat{y} \leq 1$

Intuitively

what is the chance that  
 $x$  is a cat picture?

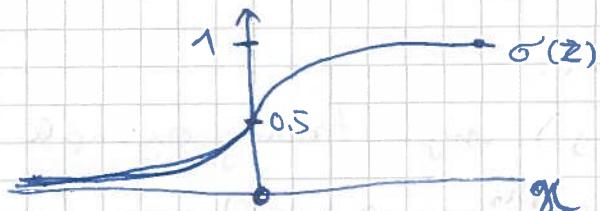
$$x \in \mathbb{R}^{n_a}$$

Parameters :  $w \in \mathbb{R}^{n_a}$ ,  $b \in \mathbb{R}$

Output :  $\hat{y} = \sigma(\underbrace{w^T x + b}_{\text{linear regression would not work}})$

$[\hat{y} = w^T x + b] \rightarrow$  linear regression would not work because  
 $\hat{y}$  should be  $\in [0, 1]$

Sigmoid function :  $\sigma$



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

④ if  $z$  large  $\Rightarrow \sigma(z) \approx \frac{1}{1+0} = 1$  / if  $z$  large negative  
 $\Rightarrow \sigma(z) \approx 0$

## ← Logistic Regression:

your job is to learn the parameters  $w, b$

so that  $\hat{y}$  becomes a good estimate of  $P(y=1|n)$

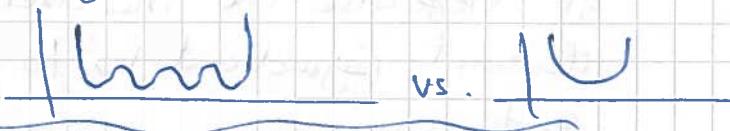
## Logistic Regression Cost Function:

- ground truth

We could define the ~~Cost function~~ / "Loss function" /  
or "error function" as ~~Cost~~

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2 \quad (\text{square error})$$

but in logistic regression people don't do it because  
the optimization ~~problem~~ problem becomes non-convex  
when learning the parameters. It will have multiple  
local optima and gradient descent won't find a  
global optima



Loss Function: defines how good the estimate  $\hat{y}$  is  
of  $y$  is.

works on  
one single  
training  
example

~~Cost function~~  $L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$

Intuition:

If  $y=1$ :  $L(\hat{y}, y) = -\log \hat{y}$  we want it to be as small as possible

$\Rightarrow$  let's make  $\log \hat{y}$  as big as possible  $\Rightarrow$  make  $\hat{y}=1$

If  $y=0$ :  $L(\hat{y}, y) = -\log (1-\hat{y})$  we want it to be as small as possible  $\Rightarrow$  make  $\log (1-\hat{y})$  as large as possible  
 $\Rightarrow$  make  $\hat{y}=0$

Cost Function: for the entire training set

average of the loss function  $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$

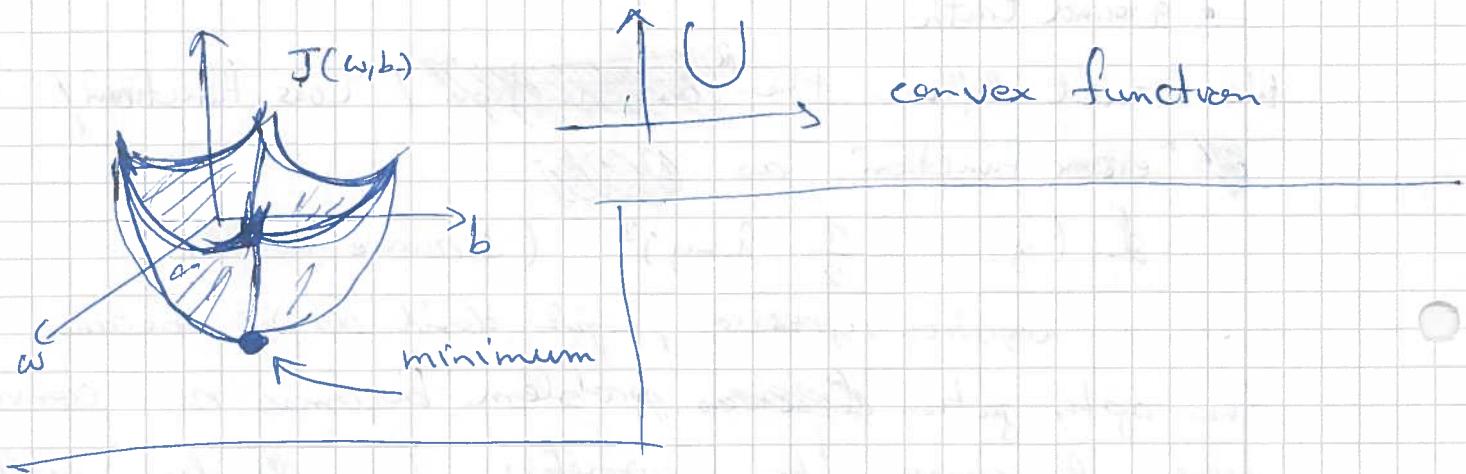
5

## Gradient Descent

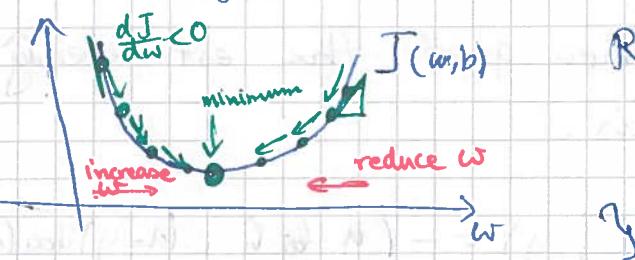
Recap:  $\hat{y} = \sigma(w^T + b)$ ,  $\sigma(z) = \frac{1}{1 + e^{-z}}$

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

- We want to find  $w, b$  that makes our cost function  $J(w, b)$  minimum.



So: initialize  $w$  and  $b$ , to zero or random,  
go step wise at the direction of the steepest  
descent (gradient), to converge to the global  
optima.



Repeatedly:

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

*"dw"*

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

$\frac{dJ}{dw}$  slope of the function

so that we know in which direction  
we need to step in order to go down hill.

$J(w, b)$

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

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

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

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

*"Partial derivative"*  
*"slope of the Function in the w direction"*

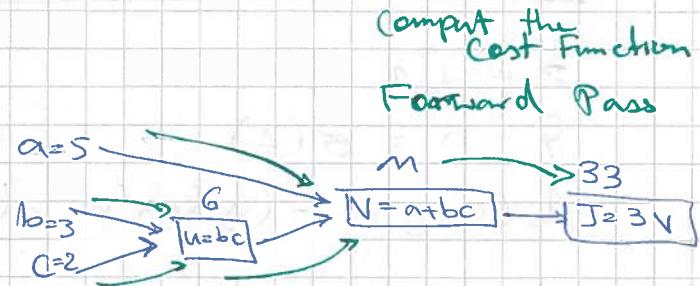
## Computation Graph

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

$$u = b * c$$

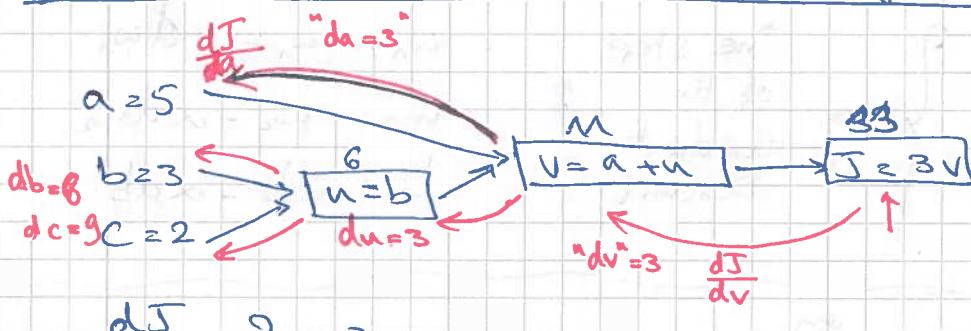
$$v = a + u$$

$$J = 3v$$



- Comes in handy when there is an output variable that you want to optimize.
- In case of logistic regression  $J$  is the cost function, that we want to minimize
- left to right pass computes the value of  $J$
- A right to left pass would be the most natural to computation of the derivatives

## Derivatives with the Computation Graph



Backward Pass  
Compute the derivative  
of the Cost Function

In code following convention:

use variable name  
"dvar" for  
 $\frac{d \text{FinalOutputVar}}{d \text{var}}$

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

$$\frac{dJ}{da} = 3 = \underbrace{\frac{dJ}{dv}}_{\substack{\text{chain rule} \\ a \rightarrow v \rightarrow J}} \cdot \underbrace{\frac{dv}{da}}_{=1} = 3 \times 1$$

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

$$\frac{dJ}{db} = \underbrace{\frac{dJ}{dv}}_{=3} \cdot \underbrace{\frac{dv}{db}}_{=2} = 3 \times 2 = 6$$

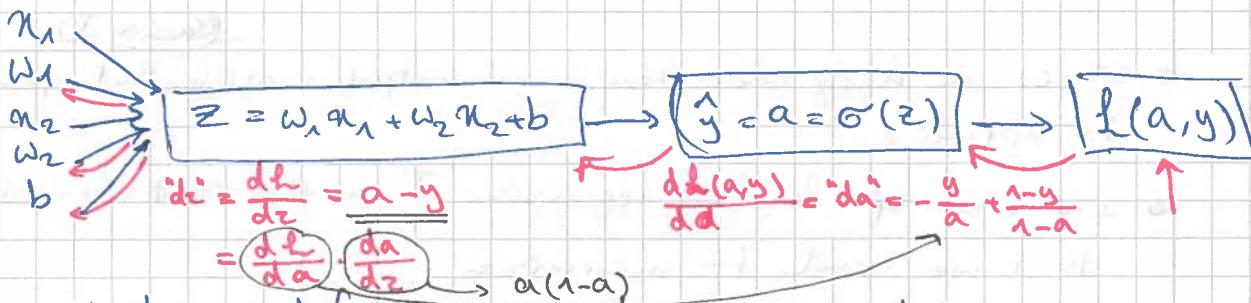
$$\frac{dJ}{dc} = \underbrace{\frac{dJ}{dv}}_{=3} \cdot \underbrace{\frac{dv}{dc}}_{=3} = 3 \times 3 = 9$$

## Logistic Regression Gradient Descent

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

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

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



- Task: modify the parameters  $w_1, w_2$ , and  $b$  in order to reduce the loss  $L(a, y)$

• Forward propagation: compute the loss

• Backward propagation: compute the derivatives

$$\left. \begin{aligned} \frac{\partial h}{\partial w_1} &= "d w_1" = a_1 dz \\ \frac{\partial h}{\partial w_2} &= a_2 dz \\ \frac{\partial h}{\partial b} &= dz \end{aligned} \right\} \Rightarrow \text{One step of the Gradient Descent}$$

$$\begin{aligned} w_1 &:= w_1 - \alpha d w_1 \\ w_2 &:= w_2 - \alpha d w_2 \\ b &:= b - \alpha d b \end{aligned}$$

## Gradient Descent on m Examples

- Cost Functions

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

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

Previous lecture  
for 1 training example  
 $(x^{(i)}, y^{(i)})$   
 $d w_1^{(i)}, d w_2^{(i)}, d b^{(i)}$

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

## Algorithm : Logistic Regression on m examples

→ ~~multiple~~, 2 features

$$J = 0, dw_1 = 0, dw_2 = 0, db = 0$$

For  $i=1$  to  $m$ :

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

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

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

$$dw_1 += a_n^{(i)} dz^{(i)}$$

$$dw_2 += a_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J / m$$

$$dw_1 / m$$

$$dw_2 / m$$

$$db / m$$

- $dw_1, dw_2, db$  don't have a superscript. They are accumulators.

One step of Grad.Dsc. will then be:

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

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

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

- Repeat the above code till termination criteria is reached

- Two weaknesses

- two For loops: 1. over m training examples  
2. over n features

- Explicit for makes algorithm inefficient

→ Vectorization (get rid of for loops)

- In deep learning vectorization is no more a "nice-to-have" but a necessity

## Vectorization

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

### Non-Vectorized

$$Z = 0$$

for  $i$  in range( $n\_a$ ):

$$Z += w[i] * a[i]$$

$$Z *= b$$

$$w = [ \dots ] , a = [ \dots ]$$

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

$$a \in \mathbb{R}^{n_a}$$

### Vectorized

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

300 times faster!

because it uses parallel processing of multicore CPU / GPU

- Vectorization uses parallelization instructions

SIMD  
simd

- Rule of thumb:

Whenever possible avoid using explicit For Loops!

## Vectorizing Logistic Regression

$$J=0, \underline{dw1}=0, \underline{dw2}=0, db=0$$

for  $i = 1$  to  $m$ :

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

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

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

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

$$\left. \begin{array}{l} dw_1 += a_1^{(i)} dz^{(i)} \\ dw_2 += a_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{array} \right\} n_a = 2$$

$$J = J/m, dw_1 = dw_1/m, dw_2 = dw_2/m, db = db/m$$

Vectorized:

$$dw = \text{np.zeros}((n_a, 1))$$



$$Z = \text{np.dot}(w.T, X) + b$$

$$A = \sigma(Z)$$



broad-casting  
→

$$dw += a^{(i)} dz^{(i)}$$

$$X = \begin{bmatrix} | & | & | \\ a^{(1)} & a^{(2)} & \cdots & a^{(m)} \\ | & | & | \end{bmatrix} \quad (n_a, m) \quad X \in \mathbb{R}^{n_a \times m}$$

$$\xrightarrow{w^T} \begin{bmatrix} | & | & | \\ a^{(1)} & a^{(2)} & \cdots & a^{(m)} \\ | & | & | \end{bmatrix}$$

$$Z = [z^{(1)} \ z^{(2)} \ \cdots \ z^{(m)}] = w^T X + [b \ b \ \cdots \ b] = \underbrace{[w^T a^{(1)} + b \ w^T a^{(2)} + b \ \cdots \ w^T a^{(m)} + b]}$$

← • Broadcasting: Expanding values/vectors to the required shape.



$$dz^{(n)} = a^{(1)} - y^{(n)}, dz^{(2)} = a^{(2)} - y^{(2)}, \dots$$

$$dZ = \underbrace{[dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}]}_{n \times m}$$

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

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

$d\omega = 0$	$db = 0$
$d\omega += n^{(1)} dz^{(1)}$	$db += dz^{(1)}$
$d\omega += n^{(2)} dz^{(2)}$	$db += dz^{(2)}$
$\vdots$	$\vdots$
$d\omega / = m$	$db / = m$

m times

vectorized

$$d\omega = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

$$= \frac{1}{m} \text{np.sum}(dZ)$$

$$d\omega = \frac{1}{m} \times dZ^T$$

$$= \frac{1}{m} \left[ \begin{matrix} n^{(1)} & \dots & n^{(m)} \\ | & \dots & | \\ 1 & \dots & 1 \end{matrix} \right] \left[ \begin{matrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{matrix} \right]$$

$$= \frac{1}{m} \left[ \begin{matrix} n_1^{(1)} dz^{(1)} + \dots + n_m^{(1)} dz^{(m)} \\ n_1^{(2)} dz^{(2)} + \dots + n_m^{(2)} dz^{(m)} \\ \vdots \\ n_1^{(m)} dz^{(1)} + \dots + n_m^{(m)} dz^{(m)} \end{matrix} \right]$$

$n \times 1$   
vector

All together:

$$\bar{Z} = \omega^T X + b$$

$$= \text{np.dot}(\omega.T, X) + b$$

$$A = \sigma(\bar{Z})$$

$$dZ = A - Y$$

$$d\omega = \frac{1}{m} \times dZ^T$$

$$db = \frac{1}{m} \text{np.sum}(dZ)$$

$$\omega := \omega - \alpha d\omega$$

$$b := b - \alpha db \quad | \text{GD}$$

forward pass

backward pass

A single iteration of the Gradient Descent

Broadcasting:

$$(m, n) + (1, n) \rightsquigarrow (m, n)$$

$$\text{matrix} * (m, 1) \rightsquigarrow (m, n)$$

$$(m, 1) + \text{IR} \rightsquigarrow (m, 1)$$

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

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



## A Note on Python/numpy Vectors

\* Broadcasting advantage/powerful, but also if not used with care can lead to very strange very hard to find bugs.

$a = np.random.rand(5) \rightsquigarrow [ \cdot \cdot \cdot \cdot \cdot ]$

$a.shape \rightsquigarrow (5, )$  rank 1 array!

$a = np.random.rand(5, 1) \rightsquigarrow [ \cdot \cdot \cdot \cdot \cdot ]$

$a.shape \rightsquigarrow (5, 1)$  column vector

$\text{assert}(a.shape == (5, 1)) \leftarrow !$

$a = a.reshape((5, 1))$

- ① Don't use "rank 1 arrays", always use column/row vectors with defined shape       $n \times 1$ ,  $n \times n$  matrices
- ② Do assertions regularly
- ③ User reshape to make sure your vectors have correct shape

## Explanation of Logistic Regression cost Function

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

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

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

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

$y$  has to be 1 or 0 because it is a binary classification

$$P(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)} \quad \left. \right\} \rightarrow \text{summarizes } P(y|x)$$

- ② log function is strictly monotonic func
-

$$\log P(y|n) = \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ = -\frac{1}{2} (\hat{y}, y)$$

- negative sign because:

- we want to maximize the probability  
 $\Rightarrow$  minimize the loss function

- For cost function on the entire training set:

Assuming all samples are drawn i.i.d. !

identically &  
independently  
distributed.

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

$$\log P(- \quad - \quad -) = \log (- \quad - \quad -)$$

$$\underbrace{\log p(-)}_{\alpha} = \sum_{i=1}^m \underbrace{\log p(y^{(i)} | x^{(i)})}_{-L(\hat{y}^{(i)}, y^{(i)})} = - \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

④ Maximum Likelihood Estimations:

 simply means: Choose parameters that maximizes this thing

⑤ Cost Function:

maximize likelihood

$\Rightarrow$  minimize the cost

$\Rightarrow$  get rid of  $(-)$  sign

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

Add for convenience  
 to make sure the  
 quantities better scale

Programming Exercise:

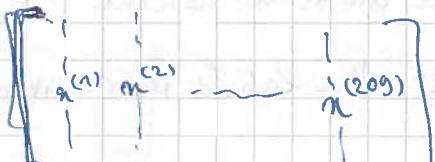
① reshape data

image  
 $(209, 64, 64, 3)$

$$X.reshape(a, b, c, d)$$

size dataset

$$X.reshape(X.shape[0], -1).T \rightsquigarrow (b * cd, 1)$$



④ H5 data format

⑤ Standardize :

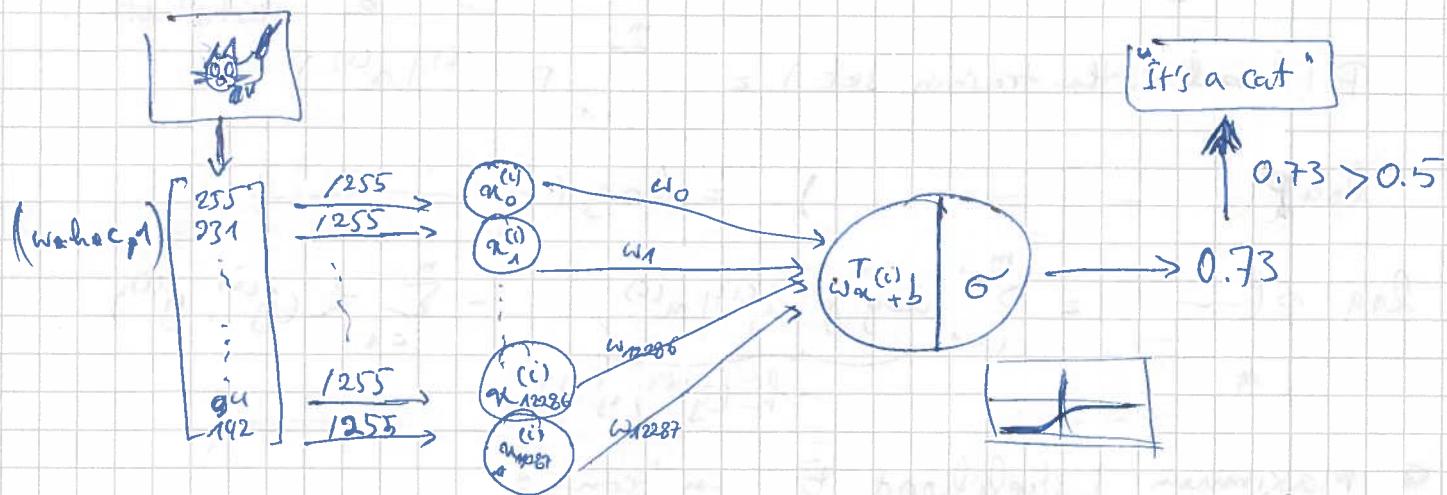
- subtract mean

- divide by standard dev

- for Images:

- just divide by 255

## General Architecture of the Learning Alg.



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

$$\text{loss function} \leftarrow L(a^{(i)}, y^{(i)}) = -y^{(i)} \log a^{(i)} + (1-y) \log(1-a^{(i)})$$

$$\text{Cost Function: } J = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

### Steps

- ① init params
- ② learn the params for this node by minimizing the cost
- ③ Use the learned params to make predictions (on the test set)
- ④ Analyse & conclude

// ① initialize ( $\omega, b$ )

- ② Optimize the loss iteratively to learn parameters ( $\omega, b$ )
  - compute the cost and its gradient
  - update the params using the gradient descent

③ use the learned ( $\omega, b$ ) to predict the labels for a given set of examples.

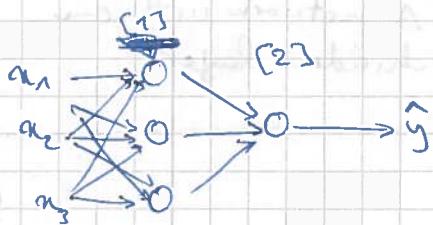
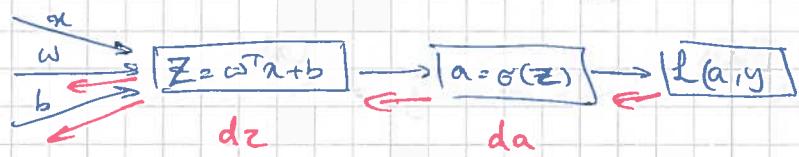
1. preprocessing the dataset is important
2. You implemented each function separately, ~~then you combined them~~  
 - initialize(), propagate(), optimize(),  
 then you built your model() by combining them.
3. Tuning the learning rate (which is an example of a "hyperparameters") can make a big difference to the algorithm.

## Week 3

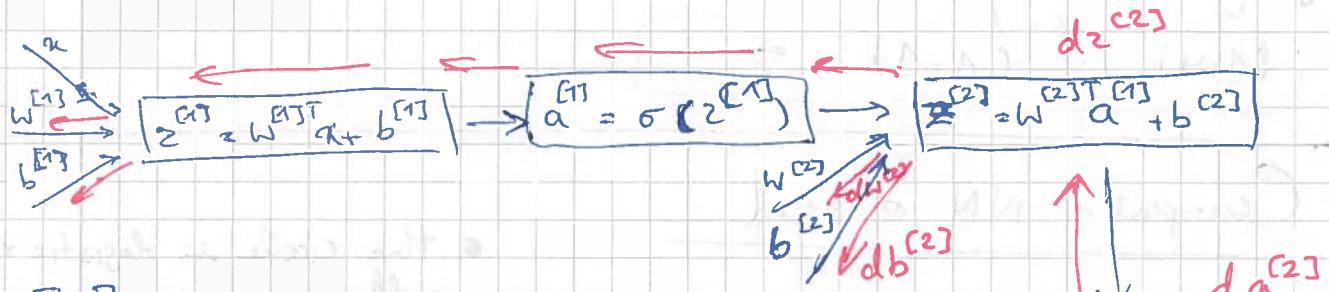
### Shallow neural Networks.

#### Learning Objectives

- Understand hidden units and hidden layers
- Be able to apply a variety of activation functions in neural network.
- Build your first forward & Backward Propagation with a hidden layer.
- Apply random initialization to your neural network.
- Become fluent in deep learning notation on Neural Network Representation
- Build and train neural networks with one hidden layer.

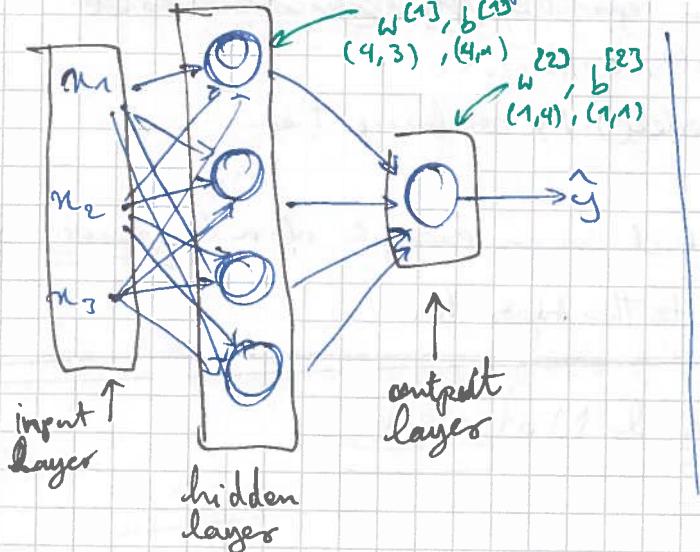


How does a NN look like?



- [1]  $\rightarrow$  layer 1, [2]  $\rightarrow$  layer 2
- NN is like taking the logistic Regression and repeating it twice

## Neural Network Representation



Input features:

$$a^{[0]} = X$$

$a$  stands for activation.

activation: Different values that different layers of NN are passing on.

$a^{[1]}_1$  → value of first node in hidden layer 1

$$a^{[1]} = \begin{bmatrix} a^{[1]}_1 \\ a^{[1]}_2 \\ a^{[1]}_3 \\ a^{[1]}_4 \end{bmatrix} \quad \rightsquigarrow 4 \text{ hidden units/nodes}$$

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

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

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

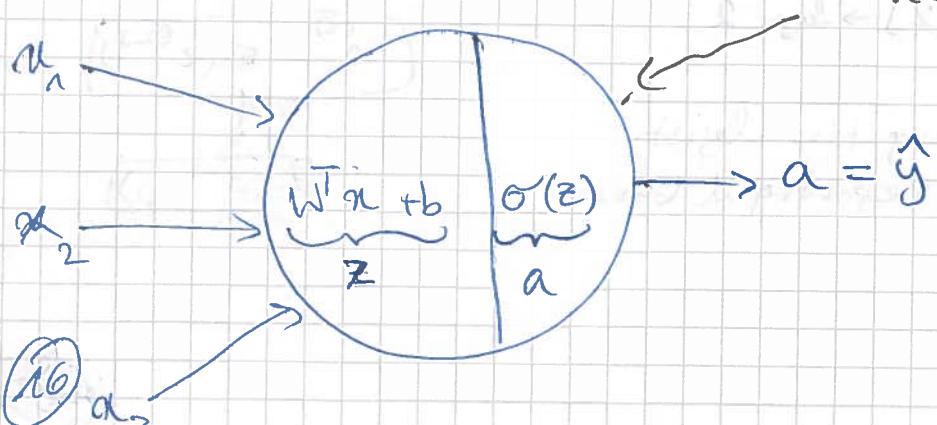
$w^{[2]}, b^{[2]}$

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

A 2 layer NN,

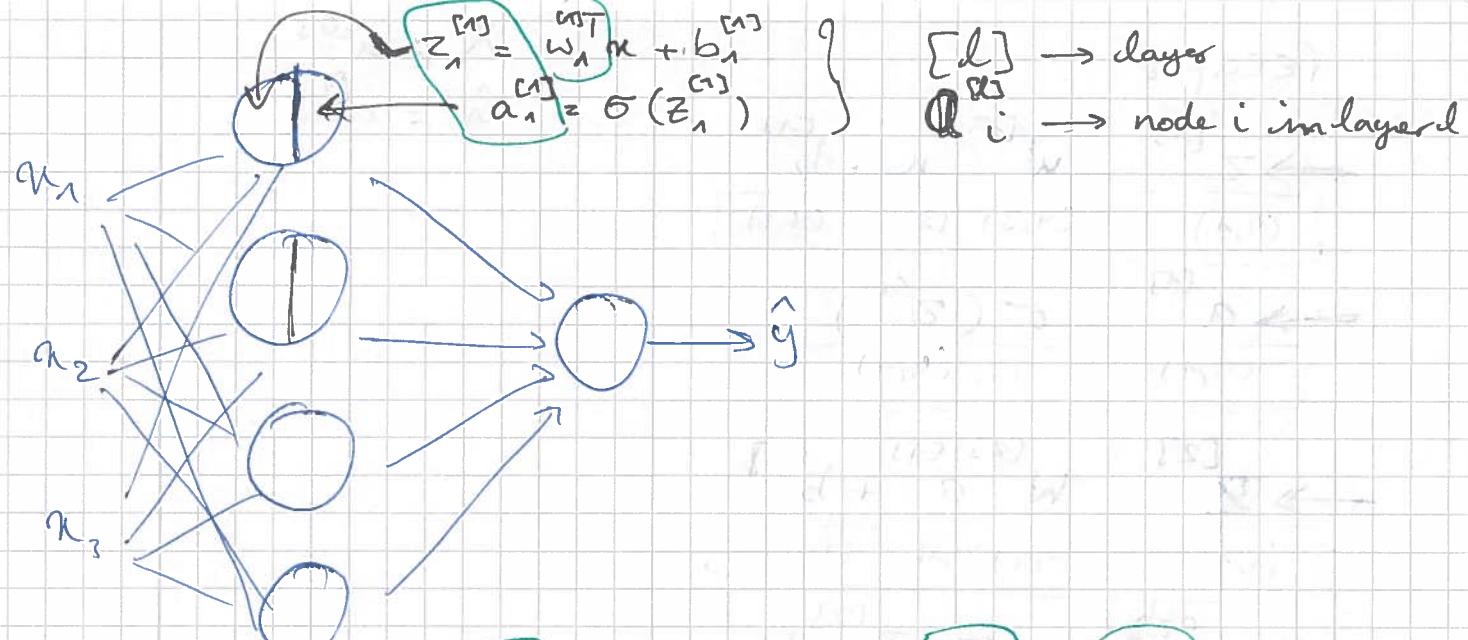
- the input layer is not counted.
- A network with one hidden layer.

## Computing NN's output



- The circle in logistic regression really represents two steps of computation

- A network with hidden layers just does this a lot more times



- accordingly:

$$\begin{aligned} z^{[1]}_2 &= w_2^T x + b_2 \\ a^{[1]}_2 &= \sigma(z^{[1]}_2) \\ z^{[1]}_4 &= w_4^T x + b_4 \\ a^{[1]}_4 &= \sigma(z^{[1]}_4) \end{aligned}$$

$$\begin{aligned} z^{[1]}_3 &= w_3^T x + b_3 \\ a^{[1]}_3 &= \sigma(z^{[1]}_3) \end{aligned}$$

### Vectorization:

$$W^{[1]} = \begin{bmatrix} w_1^T \\ w_2^T \\ w_3^T \\ w_4^T \end{bmatrix}$$

$W^{[1]}$  is a  $4 \times 3$  matrix.

- we can think of it as we have 4 logistic regression units there and each has a parameter vector  $w$  of size 3 and by stacking them together we end up with a  $4 \times 3$  matrix.

$$X \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} = \begin{bmatrix} w_1^T x + b_1 \\ w_2^T x + b_2 \\ w_3^T x + b_3 \\ w_4^T x + b_4 \end{bmatrix}$$

$$z^{[1]} = \begin{bmatrix} z^{[1]}_1 \\ z^{[1]}_2 \\ z^{[1]}_3 \\ z^{[1]}_4 \end{bmatrix} = Z^{[1]}$$

$$a^{[1]} = \begin{bmatrix} a^{[1]}_1 \\ a^{[1]}_2 \\ a^{[1]}_3 \\ a^{[1]}_4 \end{bmatrix} = \sigma(Z^{[1]})$$

recap:

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

$(4,1)$        $(4,3)$        $(3,1)$        $(4,1)$

$$X = a^{[0]}$$

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

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

$(4,1)$        $(4,1)$

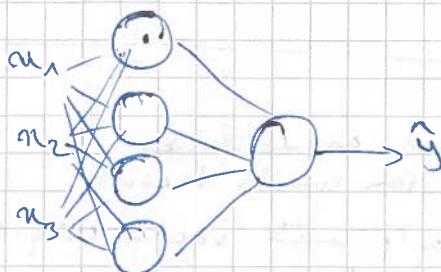
$$\rightarrow z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$(1,1)$        $(1,4)$        $(4,1)$        $(1,1)$

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

$(1,1)$        $(1,1)$

## Vectorizing across Multiple Examples



$$\left\{ \begin{array}{l} z^{[1]} = W^{[1]} a + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{array} \right.$$

$$\begin{aligned} & \xrightarrow{\text{ex}} a^{[2]} = \hat{y} \\ & a^{(1)} \xrightarrow{} a^{(2)(1)} = \hat{y}^{(1)} \\ & a^{(2)} \xrightarrow{} a^{(2)(2)} = \hat{y}^{(2)} \\ & \vdots \\ & a^{(n)} \xrightarrow{} a^{(2)(n)} = \hat{y}^{(n)} \end{aligned}$$

layer 2 training example n

for  $i = 1$  to  $m$ :

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

let's get rid of this for loop by vectorization.

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

$(n_m, m)$

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

Vectorized  
Forward  
Pass

(18)

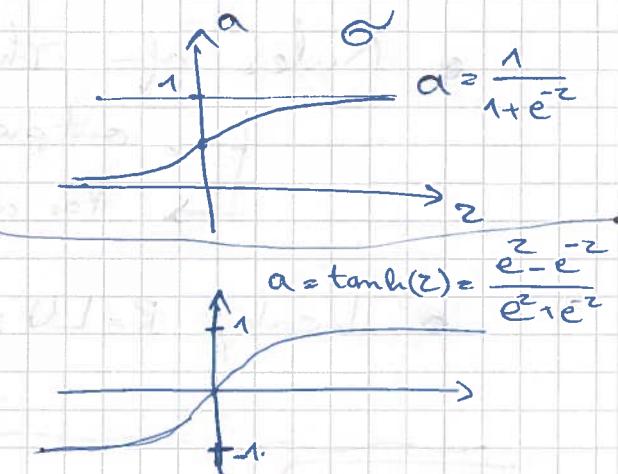
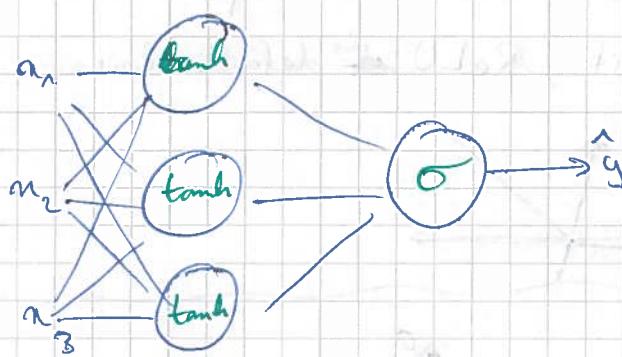
$$\begin{aligned} z^{[1]} &= \begin{bmatrix} z^{[1](1)} & z^{[1](2)} & \dots & z^{[1](m)} \end{bmatrix} \\ A^{[1]} &= \begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(m)} \end{bmatrix} \end{aligned}$$



## Explanation for the Vectorized implementation

↳ Obvious! See the video!

### Activation Function :



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

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

$$z^{[2]} = W^{(2)}x + b^{[2]}$$

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

activation function

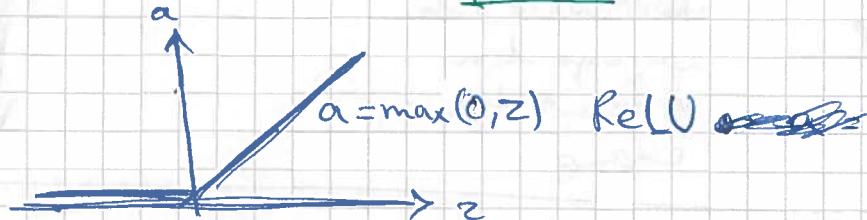
- ④ tanh is like shifted scaled version of  $\sigma$
- ④ **tanh** is almost always superior to  $\sigma$  because it has a mean of 0 and centers the data in hidden layers
- ④ use  $\sigma$  at the output layer when  $y \in \{0,1\} \Rightarrow 0 \leq \hat{y} \leq 1$
- ④ a downside of both tanh and  $\sigma$  is if  $z$  is very large or  $z$  is very small, the slope of the function nears the zero so this slow down the GD

you want



PG

## ← • Rectified Linear Unit (ReLU)



- the derivative is 1 when  $z > 0$ , and 0 when  $z < 0$
- in theory the function is not differentiable, but in practice this is ~~not important~~ at  $z=0$ .

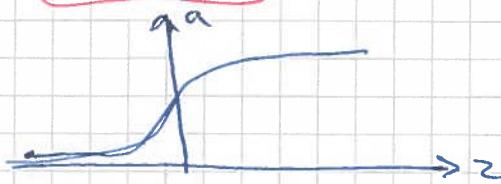
### • Rules of thumb:

- ↳ output is  $\{0, 1\} \rightarrow \tilde{o}$  for output
- ↳ for all other units, ReLU is default choice

### • Leaky ReLU:

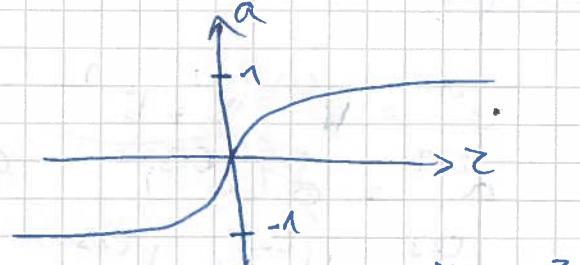


#### Recap:



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

- almost never use except for output layers of a  $\{0, 1\}$  classification



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

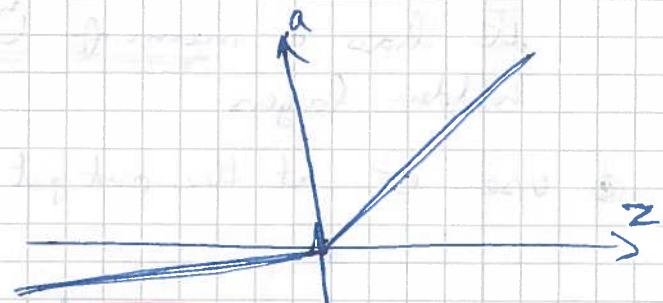
- always superior to  $\tilde{o}$



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

- default

- the most commonly used
- if not sure, use ReLU



$$\text{Leaky ReLU: } a = \max(0.01, z)$$

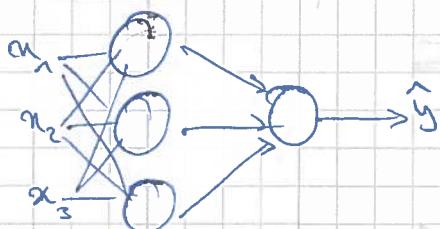
- sometimes better
- not very common
- try it, if better use it

## ① Advise:

- when starting a new NN project, there are a lot of choices on how to build the network:
  - number of hidden units
  - choice of the activation function
  - how to initialize the choices.
- there are common practices in the industry.  
But it always depend on the specific problem.
- so good is to first try different choices on a smaller data set, and see how it works

Development Set

## Why Do we need a sort of Non-linear Activation Function



Given an

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

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

why not get rid of  $g$  function and just say  $a^{[1]} = z^{[1]}$   
or let's say alternatively

$$g(z) = z$$

\* called "linear activation function"

better name would be an "identity activation function"

If we do this, then our model just computes  $\hat{y}$  as a linear function of our input features:

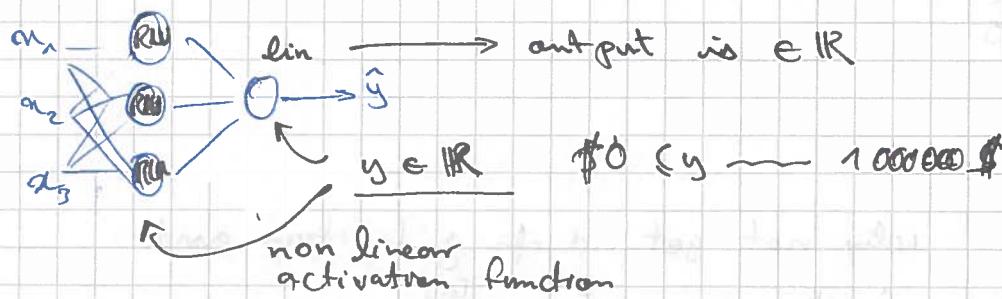
$$\begin{aligned} a^{[1]} &= z^{[1]} = w^{[1]} a + b^{[1]} \\ a^{[2]} &= z^{[2]} = w^{[2]} a + b^{[2]} \end{aligned}$$

$$\Rightarrow a = w^{[2]} (w^{[1]} a + b^{[1]}) + b^{[2]} = \underbrace{(w^{[2]} w^{[1]})}_{W'} a + \underbrace{(w^{[2]} b^{[1]} + b^{[2]})}_{b'} \Rightarrow a^{[2]} = \hat{y} = W' a + b'$$

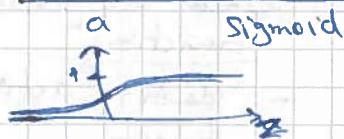


- ① the network only outputs a linear function of the input
- ② no matter how many layers you have, the output is always a linear function of the input, so you could simply achieve it without ~~any~~ any hidden layers
- ③ if you have linear activation in hidden layers and a sigmoid at last node, the the model is not more expressive than a normal logistic regression model

- ④ One Case you want to have a linear activation function is doing ML on a "Regression" problem like estimating the price of a house



### Derivatives of Activation Functions



$$g(z) = \frac{1}{1+e^{-z}} \quad \Rightarrow \quad \frac{dg}{dz} = g(z)(1-g(z))$$

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

tanh:



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad \Rightarrow \quad \frac{dg}{dz} = 1 - (\tanh(z))^2$$

$$g' = 1 - g^2 = 1 - a^2$$

ReLU

$$g' = \begin{cases} 0 \\ 1 \end{cases}$$

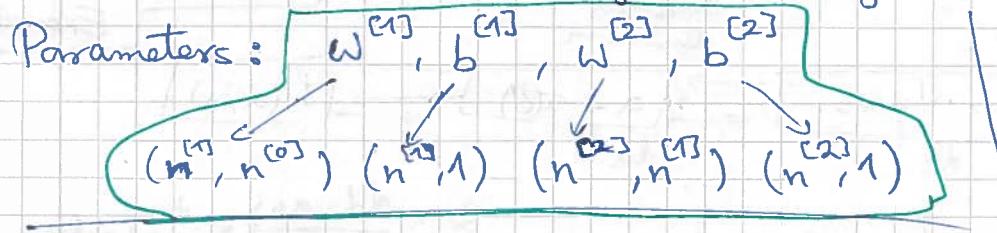


$$g(z) = \max(0, z) \quad \Rightarrow \quad g'(z) = \frac{dg}{dz} = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \\ \text{undef} & 0 \end{cases}$$

for Leaky ReLU:  $g(z) = \max(0.01z, z) \Rightarrow g' = \begin{cases} 0.01 & z < 0 \\ 1 & z \geq 0 \end{cases}$

# Gradient Descent for NN

Neural Network with a single hidden layer:



$$n_x = n^{[0]}, n^{[1]}, n^{[2]}$$

Assuming we are doing a binary classification:

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

Gradient Descent:

- init Params randomly

- Repeat {

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

- compute the derivatives ( $d w^{[1]} = \frac{dJ}{dw^{[1]}}$ ,  $d b^{[1]} = \frac{dJ}{db^{[1]}}$ , ...)

- update gradient descent

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

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

$$w^{[2]} = w^{[2]} - \alpha d w^{[2]}$$

$$b^{[2]} = b^{[2]} - \alpha d b^{[2]}$$

Forward Prop. equations: Back Prop.:

$$z^{[0]} = w^{[0]} X + b^{[0]}$$

$$A^{[0]} = g^{[0]}(z^{[0]})$$

$$z^{[1]} = w^{[1]} A^{[0]} + b^{[1]}$$

$$A^{[1]} = g^{[1]}(z^{[1]}) \rightarrow \text{bin. class.} \Rightarrow g^{[1]} = \sigma$$

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

$$Y = [y^{(0)} \ y^{(1)} \ \dots \ y^{(m)}]$$

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

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

$$d z^{[1]} = \underbrace{w^{[1]} d z^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]}'(z^{[1]})}_{\text{derivative of } g^{[1]}} \quad \begin{matrix} \text{elementwise prod} \\ \text{dotprod} \end{matrix}$$

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

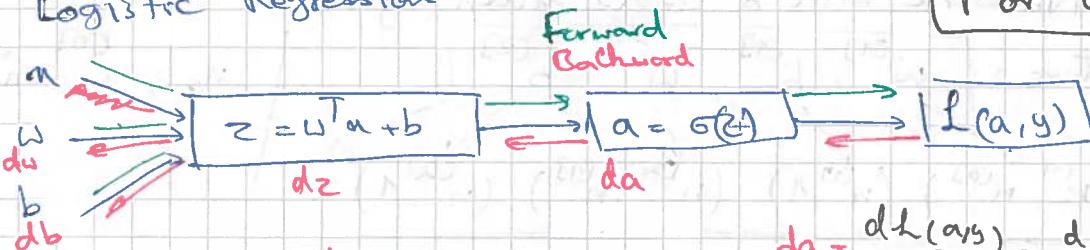
$$d b^{[1]} = \frac{1}{m} \text{np.sum}(d z^{[1]}, \text{axis}=1, \text{keepdims=True})$$

no  $(n^{[1]}, 1)$   
rank 1 array  
23

output no  
rank 1  
 $(n^{[1]}, 1)$   
thi.  
output  
 $(n^{[2]}, 1)$

# Backpropagation Intuition

## Logistic Regression



$$dw = dz \cdot x$$

$$db = dz$$

$$dz = a - y$$

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

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

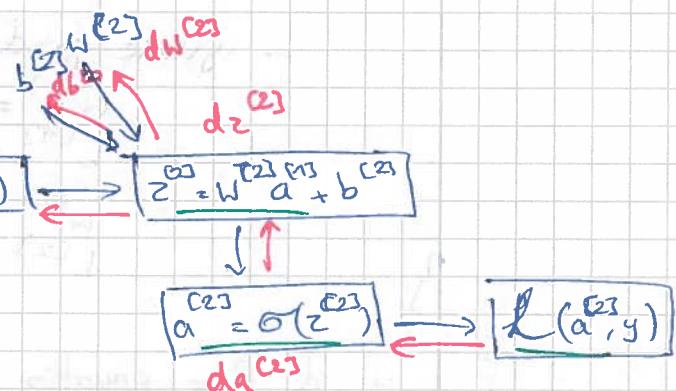
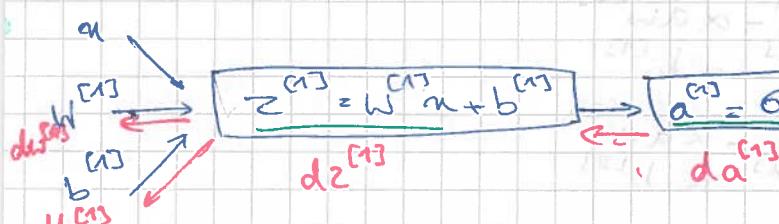
$$g'(z) = a(1-a)$$

$$\frac{da}{dz} = \frac{\partial L(a, y)}{\partial a} = \frac{d}{da} (-y \log a - (1-y) \log(1-a)) \\ = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z}$$

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

- calculation for a NN will be very similar to Log. Regression, unless we do it twice



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

$$dw^{[2]} = dz^{[2]} a^{[1] T}$$

$$db^{[2]} = dz^{[2]}$$

quite similar to log. reg.: "dw = dz  $\cdot$  x"  
now  $a^{[1]}$  plays the role of x here

dw<sup>[2]</sup> is here a row vector dw<sup>[2]</sup> = [ ]  
in the case of log. reg. with a single output.

Whereas dw was a column vector.

so  $a^{[1]}$  has a Transpose

$$(n^{[1]}, 1) \quad (n^{[1]}, n^{[2]}) (n^{[2]}, 1) \quad (n^{[2]}, 1)$$

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

$$w^{[2]} : (n^{[2]}, n^{[1]}) \quad \text{same dimension as } z^{[1]}$$

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

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

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



\* Tip: When implementing back prop. just making sure that the dimensions match up, eliminates sometimes a lot of bugs already.

$$\begin{aligned} dW^{[1]} &= dz^{[1]} \cdot X^T \\ db^{[1]} &= dz^{[1]} \end{aligned}$$

Very similar to  $dW^{[2]}$  and  $db^{[2]}$   
because  $X$  plays the same role here as a  $a^{[1]}$

Recap: for a single training example:

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

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]}$$

$$dW^{[1]} = dz^{[1]} a^{[0]T}$$

$$db^{[1]} = dz^{[1]}$$

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} a^{[1]T}$$

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

Vectorized  
Back Prop.

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

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

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

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

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

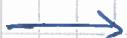
$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

Same for Back prop.

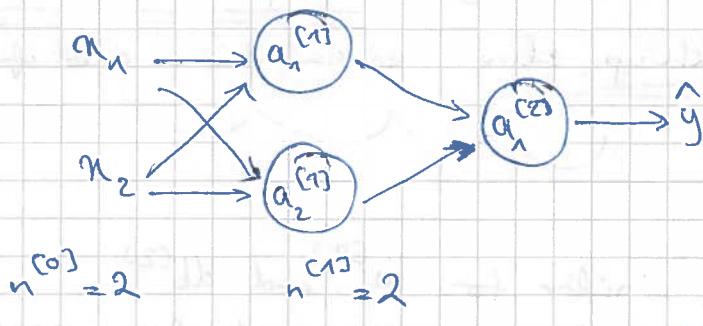
## Initializing the Weights

- for NN the weight must be initialized "randomly"

In logistic Regression we could do it by init to zero, but it does not work in NN. Why?



←



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

$$a_1^{(1)} = a_2^{(1)}$$

$$dz_1^{(1)} = dz_2^{(1)}, \Rightarrow w^{(2)} = [0 \ 0]$$

- if the initial weights for both hidden units are the same, then both units compute the same function.

Symmetric

- by induction after several iteration they still computing the same function  $\Rightarrow$  it is like having only 1 single hidden unit.

- Solution:

→ initialize the weight randomly. You want your hidden units to compute different functions.

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

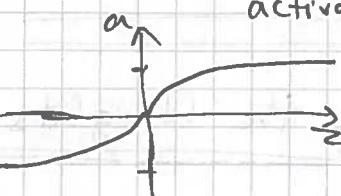
$$b^{(1)} = np.zeros((2, 1))$$

b is not important because we initialized w randomly

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

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

so that weights are small, because if activation function tanh or sigmoid



- $\Rightarrow w$  large.
- $\Rightarrow z$  large
- $\Rightarrow$  slope  $\approx 0$
- $\Rightarrow$  GD very slow
- $\Rightarrow$  learning slow

# Programming Assignment 8 Planar Data Classification with one hidden layer

## ④ Cross entropy loss

## ⑤ Scikit Learn :

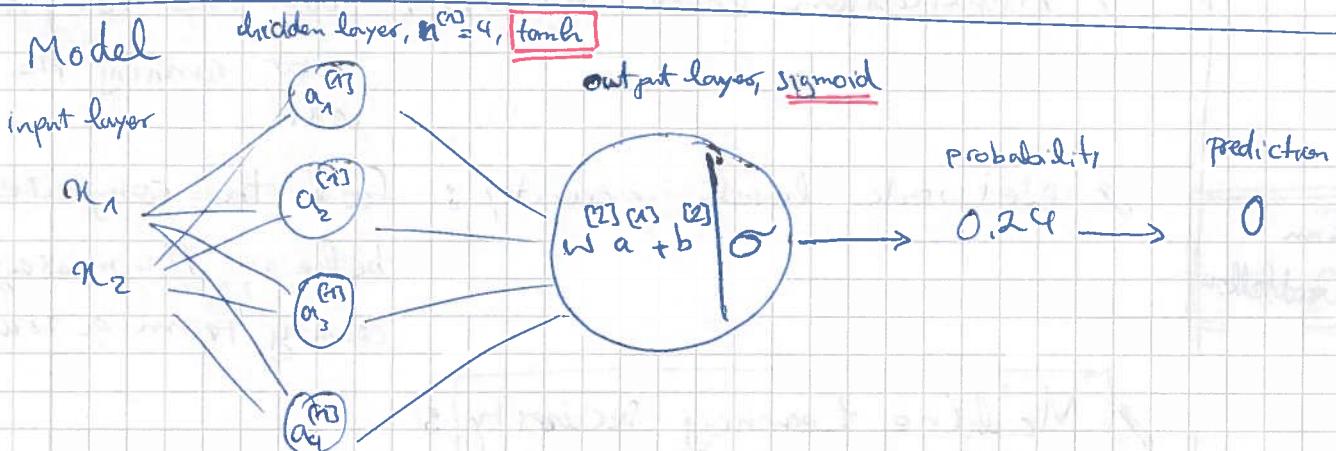
```
clf = sklearn.linear_model.LogisticRegressionCV()
clf.fit(X.T, Y.T)
```

- ⑥ plot decision boundary



- ⑦ linear Regression does not perform well because the data is not linearly separable

NN Model



for one example  $x^{(i)}$

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

$$a^{(1)(i)} = \tanh(z^{(1)(i)})$$

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

$$\hat{y}^{(i)} = a^{(2)(i)} = \sigma(z^{(2)(i)}) ; y^{(i)}_{\text{prediction}} = \begin{cases} 1 & \text{if } a^{(2)(i)} > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

Cost function:

$$J = -\frac{1}{m} \sum_0^m \left( y^{(i)} \log(a^{(2)(i)}) + (1-y^{(i)}) \log(1-a^{(2)(i)}) \right)$$

## General Methodology to Build a NN:

1. Define the NN structure (# input units, # hidden units, etc)
2. Initialize the model parameters
3. Loop
  - Impl forward prop.
  - Compute loss
  - Impl Backward prop. to get the gradients
  - update params (gradient descent)

go!

- Accuracy is really high compared to Logistic Regression (47%)
- Neural Networks are able to learn even slightly non-linear decision boundaries unlike Logistic Regression

• post on Github / ArXiv ↗

// Application level security : fool the computer into running the wrong code

// Network level security : fool the computer into believing the messages are coming from a trusted party.

Machine Learning Security :

↳ Adversarial Networks

## Week 4

### Deep Neural Networks

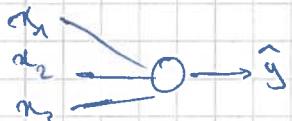
Learning Objectives:

- See deep neural networks as successive blocks put one after each other
- Build and train deep L-layer NNs
- Analyze Matrix & vector dimensions to check NN impls
- Understand how to use a cache to pass information from forward prop. to backprop.
- Understand the role of hyperparameters - d n

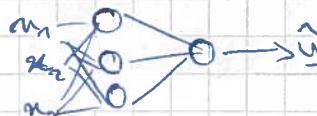
## Deep L-layer neural Networks

- forward prop., back prop., log. regression, NN with single hidden layer, Vectorization, init weights randomly

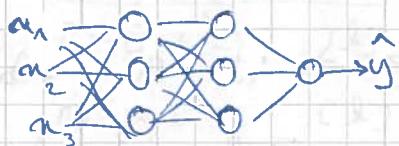
3 - What is a deep NN?



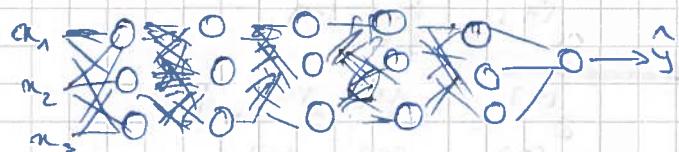
logistic regression  
"shallow model"



1 hidden layer



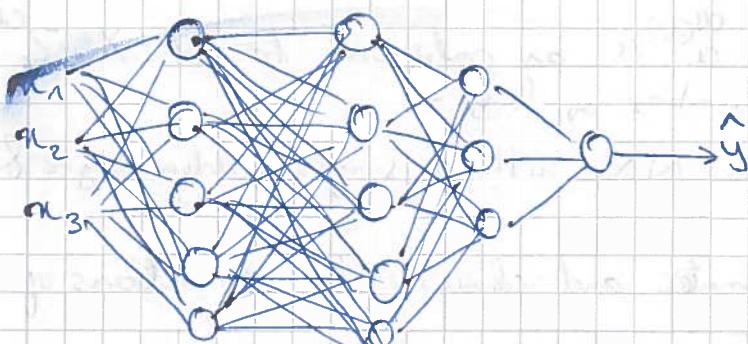
2 hidden layers



5 hidden layers "deep"

- # hidden layers: another example of hyperparameters

## Deep NN Notation



$L = 4$  (# layers)

$n^{[l]}$  (# units in layer  $l$ ):  $n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3, n^{[4]} = n^{[L]} = 1$   
 $n^{[0]} = n_{\alpha} = 3$

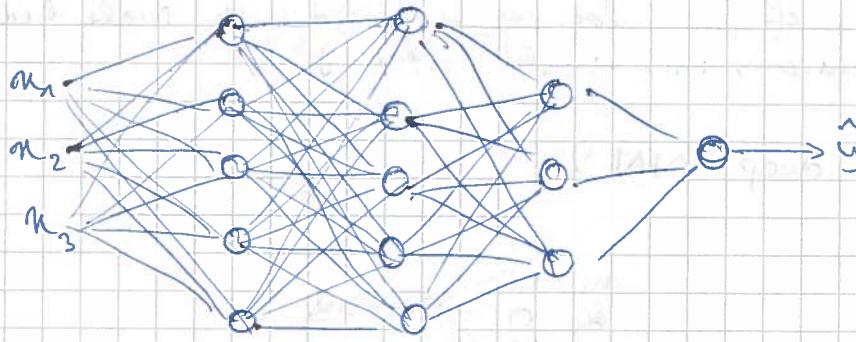
$a^{[l]}$  (activations in layer  $l$ )

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

$a^{[0]} = a$

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

# Forward Propagation in a Deep Network



$n_i$ : single training example:

$$z^{[1]} = W^{[1]} \alpha^{[0]} + b^{[1]}$$

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

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

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

:

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

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

$$! \rightsquigarrow \hat{y} = g^{[4]}(\alpha^{[4]}) = A^{[4]}$$

General forward

Prop. equation:

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

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

Vectorized Version:

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

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

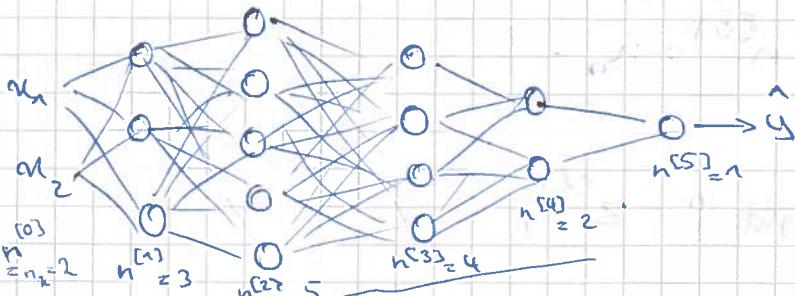
$$X = A^{[0]}$$

for  $l=1-L$

- stacking  $\alpha^{[i]}, z^{[i]}, a^{[i]}$ 's as columns to get  $X, Z^{[l]}, A^{[l]}$
- we need a for of number of layers
- very similar to a NN with a single hidden layer & repeating that.
- we need to systematic and checking the dimensions of matrices

Getting your Matrix Dimensions right:

$$0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5$$



$$L = 5$$

$$z^{[1]} = W^{[1]} \alpha^{[0]} + b^{[1]}$$

> let's ignore  $b$  for now.

$$(3,1) = (3,2) (2,1)$$

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

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



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

$$W^{[2]} : (n^{[2]}, n^{[1]}) \rightarrow (5, 3)$$

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

$$(5, 1) = (5, 3) (3, 1)$$

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

$$W^{[4]} : (2, 4), W^{[5]} : (n, 2)$$

Vectorized:

- dimensions of  $W, b, dW, db$  stay the same

- $Z, A, X =$

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

$$(n^{[l]}, m) = (n^{[l]}, n^{[0]}) (n^{[0]}, m) \quad (n^{[l]}, 1) \xrightarrow{\text{broad casting}} (n^{[l]}, m)$$

General Rule

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

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

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

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

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

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

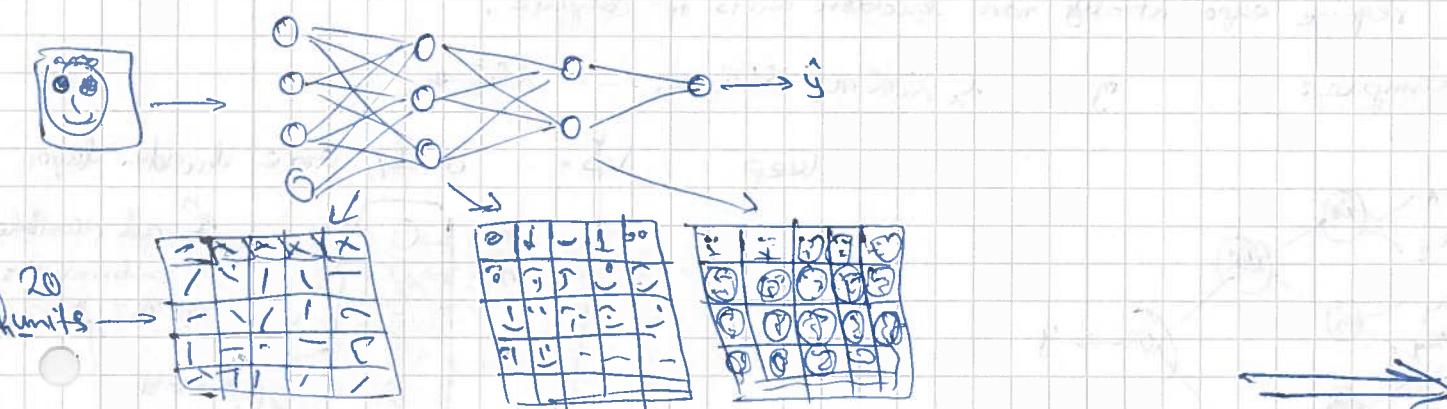
$$Z^{[l]} \rightarrow A^{[l]} : (n^{[l]}, m)$$

$$l=0, A^{[0]} = X : (n^{[0]}, m)$$

$$dZ^{[l]}, dA^{[l]} : (n^{[l]}, m)$$

Why deep representations work?

- The NN does not have to be big but needs to be deep
- what does a NN compute?



first layer  
detects ~~edges~~  
grouping together pixels

detects ~~parts of~~  
graying together  
~~edges~~

different  
types of  
face by  
grouping ~~parts of face~~

(31)

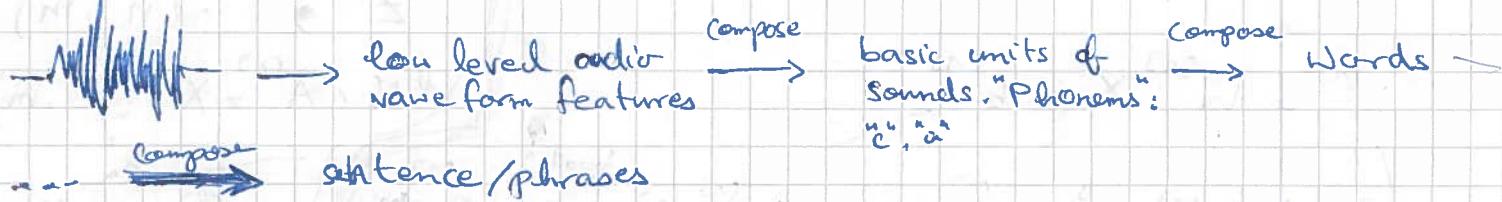
• think:

- Earlier layers of NN compute simpler functions
- later layers compute more complex function by composing those simpler functions.
- the simpler function look at smaller areas of image, complex functions look at much larger areas

edge  
detectors

! Simple to complex, hierarchical, compositional Representation

• This applies to other kinds of data like Audio as well



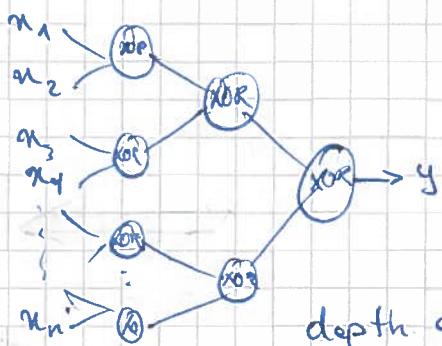
• The deeper you go into the network, the more complex things can be detected.

• Some compare it to human brain. Although that is not sure

Circuit Theory: 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.

Example:

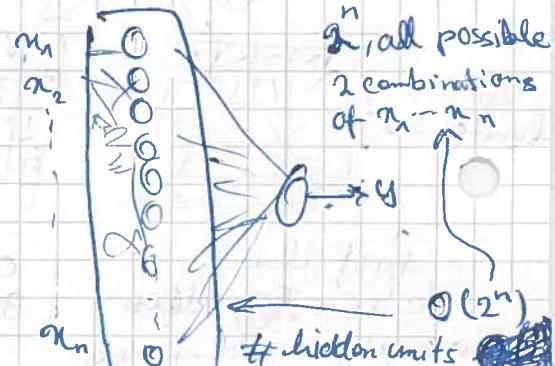
$$y = a_1 \text{XOR} a_2 \text{XOR} a_3 \dots \text{XOR} a_n$$



depth of XOR tree:  $\Theta(\log n)$

# units:  $\Theta(\log n)$

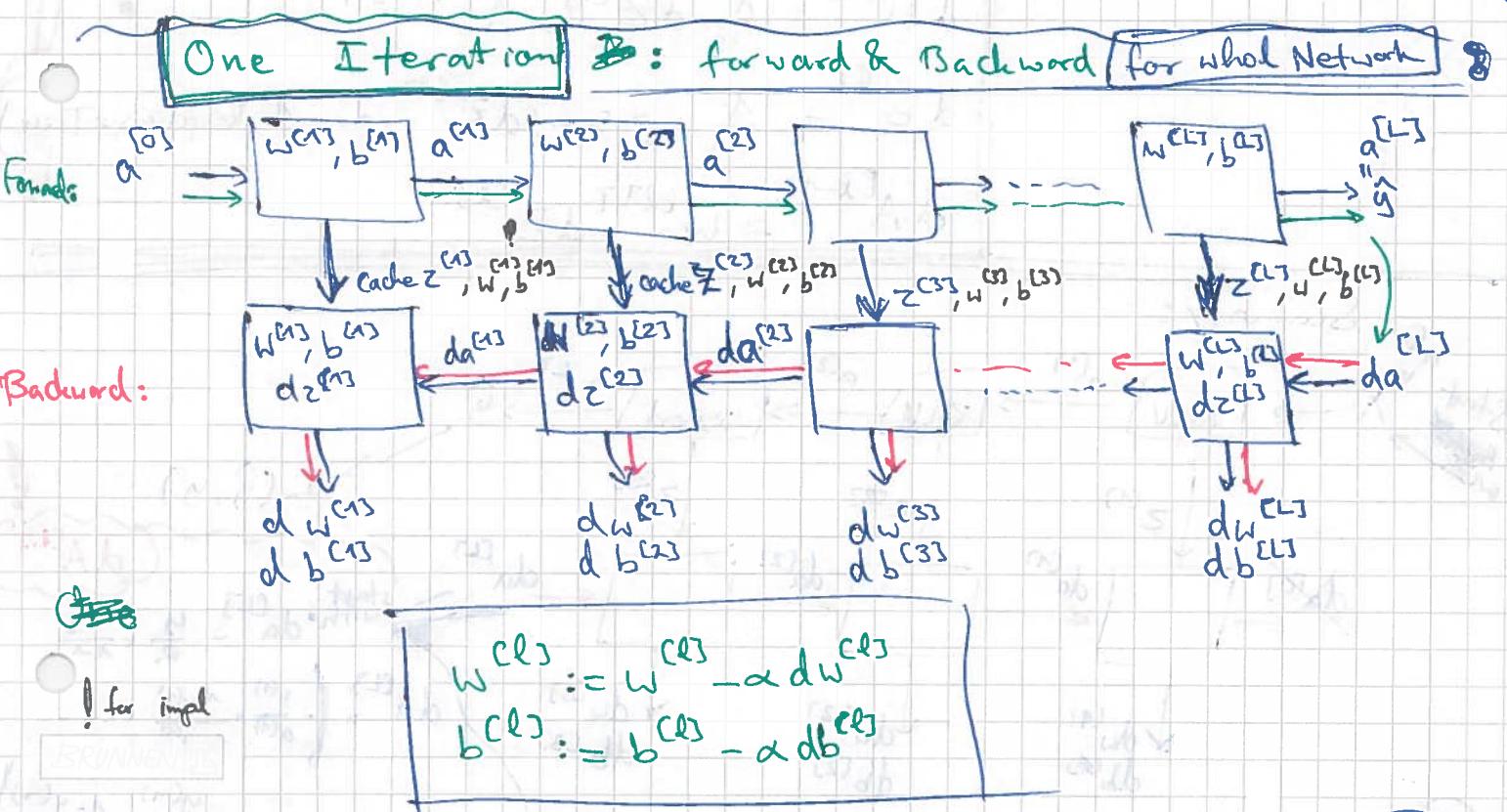
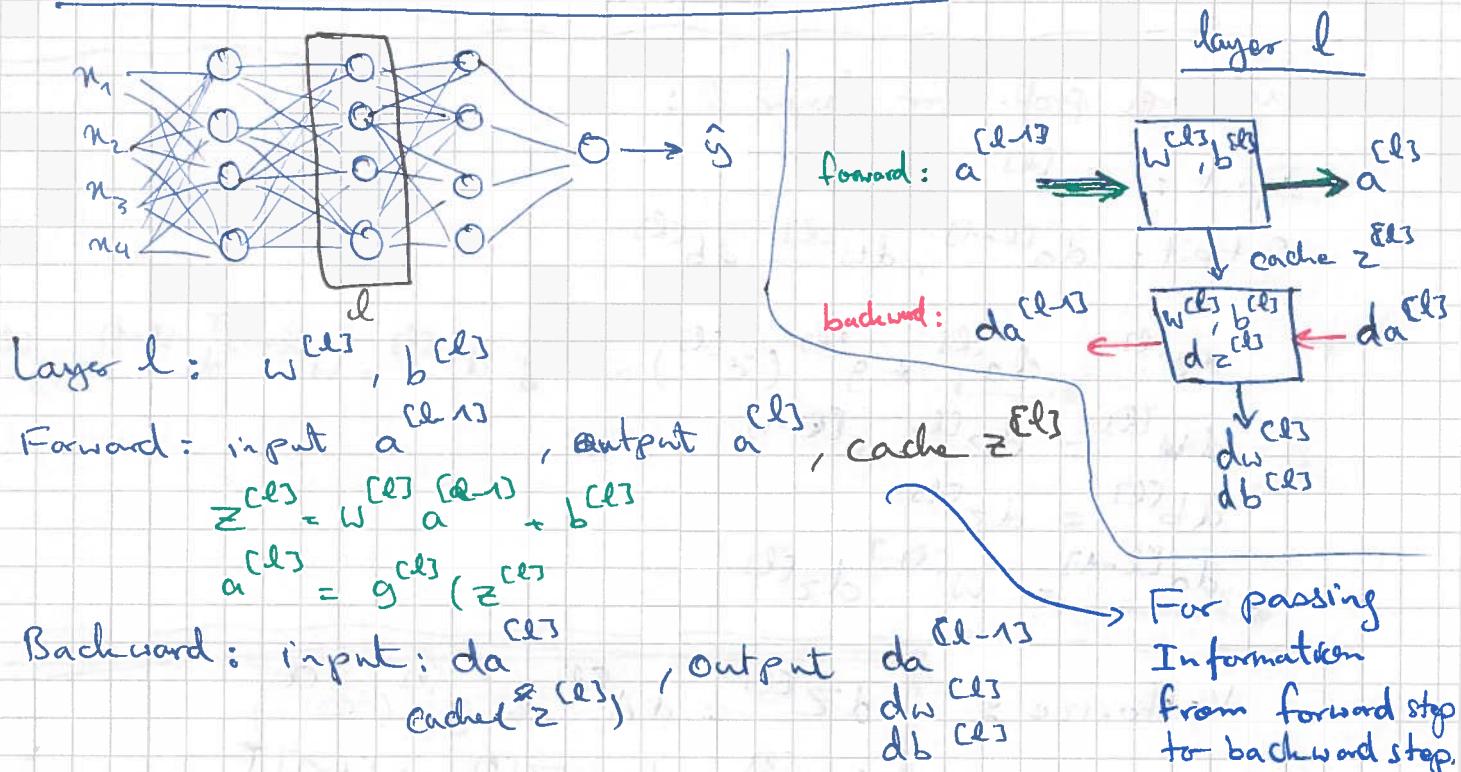
deep VS. only one hidden layer



⑥ other reason for deep learning taking off is "Branding")  
It is basically NN or NN with multiple layers rebranded.

⑦ Tip: don't get excited. Don't start with many many layers.  
Start even with log. regression. Then add other layers.  
and see how your model behaves.

## Building Blocks of Deep Neural Networks



## Forward and Backward Propagation:

Input  $a^{[l-1]}$   
Output  $a^{[l]}$ , cache ( $z^{[l]}$ )  
 $z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$   
 $a^{[l]} = g^{[l]}(z^{[l]})$

### Forward Prop. for layer $l$

Vectorized:

$$Z^{[l]} = w^{[l]} A^{[l-1]} + b^{[l]}$$

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

$$x = A^{[0]} \rightarrow \square \rightarrow \square \rightarrow \dots \rightarrow \square \rightarrow \dots$$

### Backward prop. for layer $l$ :

Input:  $da^{[l]}$

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

$$dz^{[l]} = da^{[l]} * g'(z^{[l]}) \rightsquigarrow dz^{[l]} = w^{[l+1]T} dz^{[l+1]} * g'(z^{[l]})$$

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

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

$$da^{[l-1]} = w^{[l]T} dz^{[l]}$$

Vectorize:

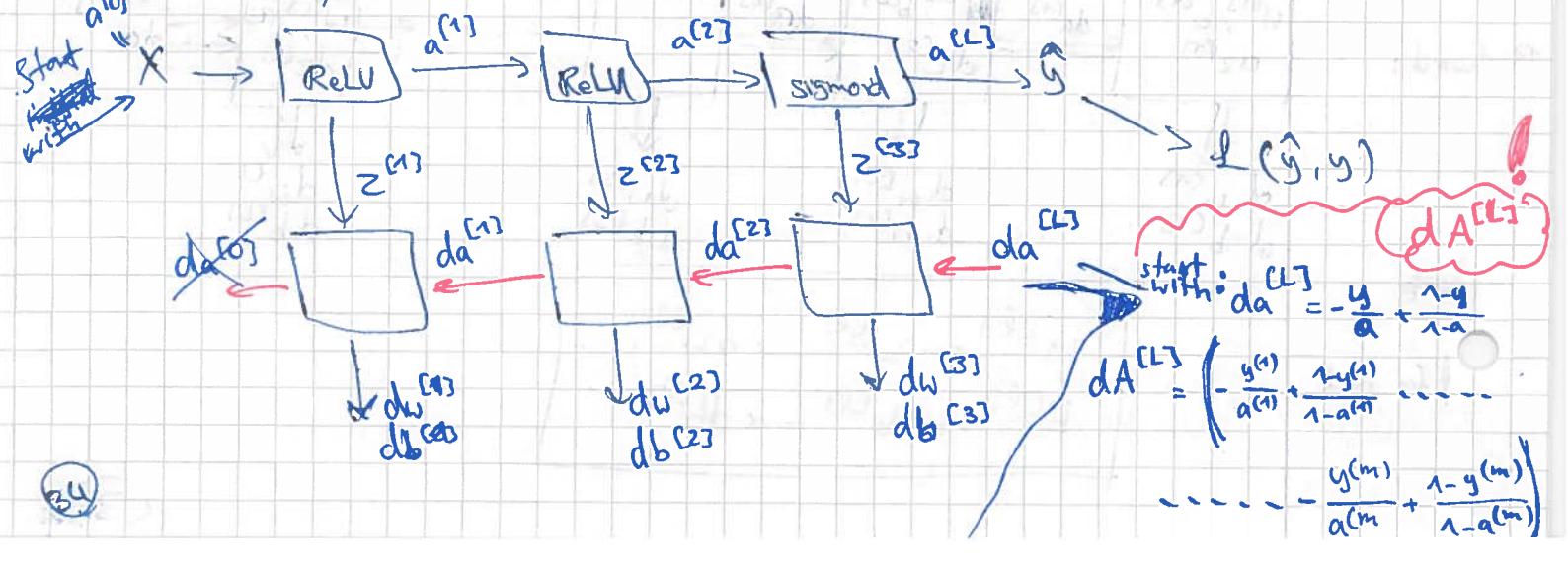
$$dZ^{[l]} = dA^{[l]} * g'(z^{[l]})$$

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

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

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

Summary:



Even Andrew Ng is sometimes surprised when the code works :)

Because a lot of complexity in ML comes from data  
not from thousands of lines of code

### Parameters vs. Hyperparameters :

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

Hyperparameter: learning rate  $\alpha$

# iterations

# hidden layers  $L$

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

choice of activation function

Hyperparameters control the ultimate parameters  $b$ 's and  $W$ 's.

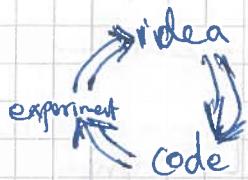
Deep learning has a lot of Hyperparameters:

Other examples:

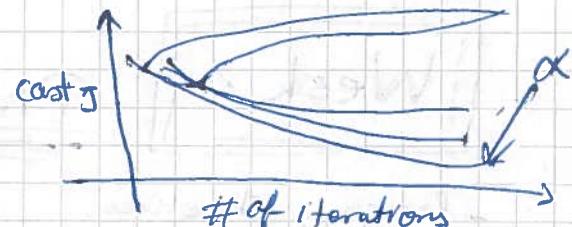
Momentum, mini-batch size, regularization params, ...

Values for

Hyperparameters can have many different settings that we need to try-out!



Applied deep learning is  
very imperical



it is very difficult to  
know in advance what to choose  
for hyperparameters,

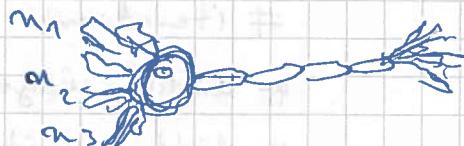
so try a lot of things and see how it works

Some times the intuition for Hyperparam from one application domain does not carry over. So you have to try a range of values

Also as the HW progresses, the Hyperparam that was very good, might no longer be the best

# Deep Learning and Human Brain

- Deep learning is a very good way to learn complex functions  $y = f(x)$ , but it does not have much to do with how brain works.
- Actually we still don't even know how a single neuron works/learn
- there are over simplified analogies of NNs and brains or a single logistic regression unit and a single neuron but they are for general public



neuron does a threshold computing on input and then fires

~~•~~ ~~threshold~~

- We still don't know whether learning in brain is like gradient descent or there is a fundamentally different process?

## Course 2 (C2 - Improving Deep Neural Networks, Hyperparameters tuning, Regularization, and Optimization )

### Week 1

#### Practical Aspects of Deep Learning

- learning objectives:

- Recall that different kinds of initialization lead to different results.
- Recognize the importance of initialization in complex neural networks
- Recognize the difference between train/dev/test sets
- Diagnose the Bias and variance issues in your model
- Learn how to use regularization methods such as dropout & L2 regularization
- Understand experimental issues in deep learning such as vanishing and exploding gradients and learn how to deal with them.
- use gradient checking to verify the correctness of your backprop implementation

# Setting up your Machine Learning Application

## Train / Dev / Test sets :

- Hyperparameters tuning
- setting up data
- make sure the optimization Alg. runs quickly

! Applied ML is a highly iterative process:

- # layers
- # hidden units
- # learning rate
- activation function



## Idea      Experiment      Code

NLP, vision, Speech, Structured data  
→  
Ads    search    security    logistic --

Start with an idea,  
and try it, and improve  
it based on outcome

! Intuitions from

one Application area often don't transfer to other application areas.

! A correct guess of correct/good hyperparameter values ~~at~~ at the first time is almost impossible

⇒ do iterations ⇒ make iterations faster and more efficient.

! - Setting up the data set correctly can make the iteration cycle much more efficient

## Train / dev / Test sets:

Data



- train on your training set
- use dev set to see which of models performs best
- use test set to evaluate the final best model.

BRUNNEN

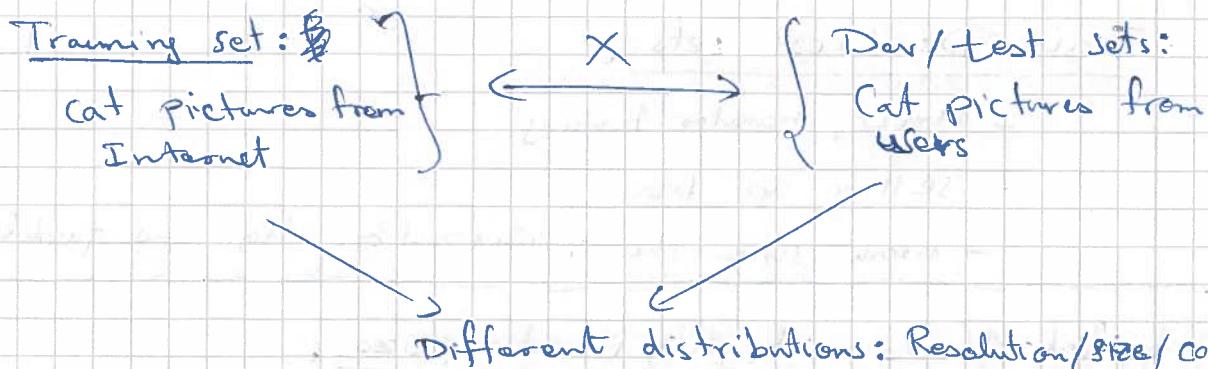
Previous era: 70 / 30 %

60 / 20 / 20 %

Big data era: 1,000,000 examples →  $\frac{10,000}{98\%}$  MIT; or even 99 / 0.5 / 0.5 % (3)

## Mis matched Train / test distribution

Cats

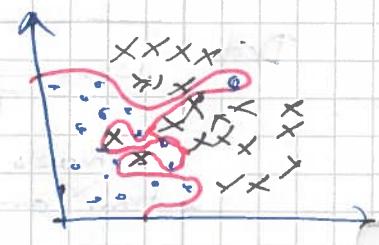
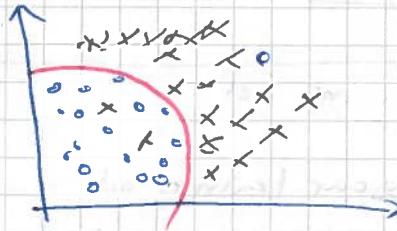
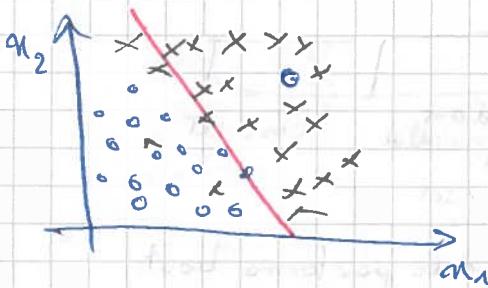


- Make sure that Dev / Test sets come from the same distribution!
- Not having a test set might be ok. (only dev set)
- make sure to distinguish between ~~Dev~~ / test set.  
Dev set is the "hold out" cross-validation set.  
Because sometimes you overfit to dev set.
- Train / Dev / Test sets allow you to iterate more quickly and more efficiently measure the bias and variance of your algorithm

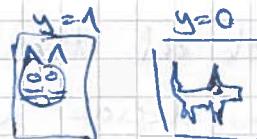
## Bias and Variance :

- good ML/DL practitioners have sophisticated understanding of Bias & Variance
- Bias & Variance are easily learned but difficult to master concepts
- In DL era there is less discussion of
  - less discussion of trade-off

Bias-Variance  
trade-off



# Cat Classification :



Two key values to look at

- Train set error:
- Dev set error:

Assuming humans  $\approx 0\%$  error

Optimal error  $\approx 0\%$   
(Bayes Error)

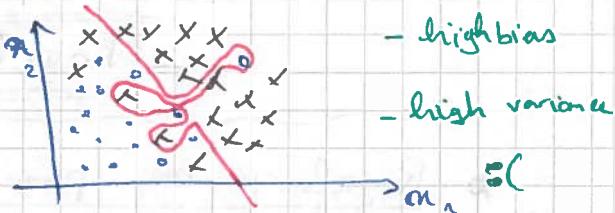
If Bayes Error was 15%,  
then the classifier with  
15% error was doing  
very well

	1 %.	15 %.	15 %.	0.5 %.
Train set error:	1 %.	15 %.	15 %.	0.5 %.
Dev set error:	11 %.	16 %.	30 %.	1 %.
	- very well on training set - poorly on Dev set - overfitting - no generalizing well on the dev set	- not even doing well on train set - under fitting - high bias	- worst of both worlds	
"High Variance"		"High Bias"	"High Bias"	"Low Bias Low Variance"

The relationship between train set & dev set error  
gives you insight into bias & variance of your model

High Bias and High Variance:

This can happen  
easily when you  
have high  
dimensional data!



## Basic Recipe for Machine Learning:

- After training an initial model:

↓  
- high bias?

(Performance on training dat.)

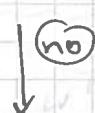
continue iterating  
till you at least  
fit the Training set

→ more layers  
→ more hidden units

→ Bigger Network

- train longer

- try another NN Architecture



no

↓  
- High variance?

(dev set performance)

yes

no

Done

→ more data

- Regularization

- try another NN archit.

← first look at the performance on train & dev sets  
and then decide what to do

! For example: If you have high Bias  
adding more data won't help

In Earlier times of ML there was the discussion about



trade-off

no need for much :)  
trade-off in DL era

But in DL Era we have now tools that let us reducing both Bias and Variance without hurting the other

## Regularization:

- If you suspect overfitting (high variance):
  - add more training data → not always possible  
- might be expensive
  - Regularization often helps

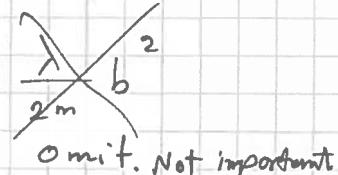
- Regularization penalizes large coefficients of W and shrinks them to zero. This way the coefficients relating to irrelevant features disappear, and model does not react to noise.

## Example Logistic Regression:

$$\min_{w,b} J(w,b) \quad w \in \mathbb{R}^n, b \in \mathbb{R}$$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

Regularization param



$$L_2 \text{ Regularization: } \|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w \quad \leftarrow \begin{array}{l} L_2 \text{ norm} \\ \text{Euclidean norm} \end{array}$$

$$L_1 \text{ Regularization: } \frac{\lambda}{2m} \sum_{j=1}^n |w_j| = \frac{\lambda}{2m} \|w\|_1 \quad \leftarrow \begin{array}{l} \text{- will make } w \text{ sparse} \\ \text{- not often used} \end{array}$$

For a Neural Network:

$$J(w^{(0)}, b^{(0)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$
$$\|w^{(l)}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{(l+1)}} (w_{ij}^{(l)})^2$$

w:  $(n^{(0)}, n^{(1)}, \dots)$

Forbenius Norm

$$\|\cdot\|_F^2$$

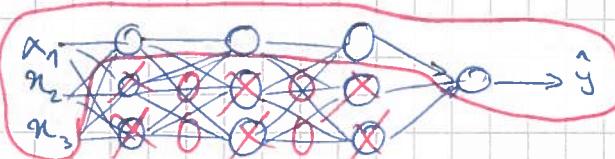
$$\frac{\partial J}{\partial w^{(l)}} = dw^{(l)} = (\text{term from backprop}) + \frac{\lambda}{m} w^{(l)}$$

GD:

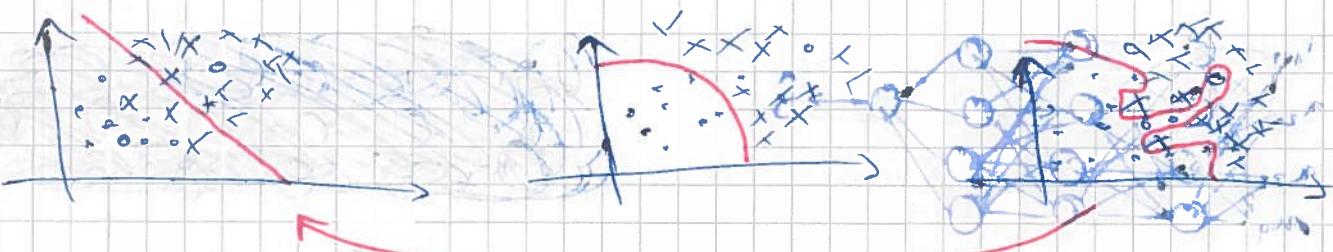
$$w^{(l)} := w^{(l)} - \alpha dw^{(l)} = w^{(l)} - \alpha \left[ (\text{term from backprop}) + \frac{\lambda}{m} w^{(l)} \right]$$
$$= w^{(l)} \left( 1 - \frac{\alpha \lambda}{m} \right) - \alpha (\text{term from backprop})$$

"Weight decay"  $\rightsquigarrow$  Alternative name for L<sub>2</sub> Regularization

Why Regularization Reduces Overfitting?



overfitting happens? Two ways



$$J(w^{(0)}, b^{(0)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

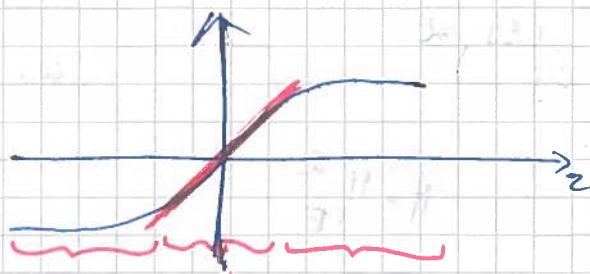
$\lambda \rightarrow \infty \Rightarrow w^{(l)} \rightarrow 0$

$\lambda$  large  $\rightsquigarrow w^{(l)} \approx 0 \rightsquigarrow$  ~~from zeroing out~~ many of the hidden units  $\rightsquigarrow$  model becomes roughly linear  $\rightsquigarrow$  going from high variance to high bias  $\rightsquigarrow$  not capable of estimating complex <sup>non-linear</sup> decision boundaries.



## Another Intuition

- assume the activation function is tanh



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

$$\lambda \uparrow \quad w^{(l)} \downarrow \quad z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$

$\lambda$  high  $\rightarrow w^{(l)}$  low  $\rightarrow z^{(l)}$  small  $\rightarrow$  "linear Regime" of tanh

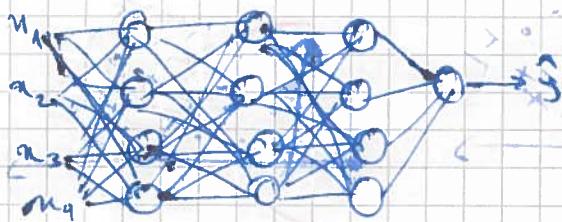
Every layer  $\approx$  linear  $\rightarrow$  model roughly linear

$\downarrow$  rot

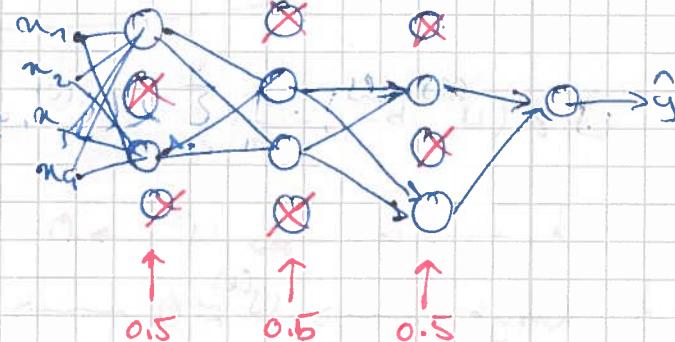
$\rightarrow$  not capable of computing complex non-linear functions

- L2 regularization  $\rightarrow$  used most often

## Dropout Regularization



- Dropout:
  - Randomly just eliminate some of the nodes and their corresponding weights.
  - By a given probability



$\rightarrow$

$\Rightarrow$  Results in a much smaller much diminished Network. Then you do backprop. Then diminish

# Implementing : ("Inverted Dropout")

Take layer 3 for example :  $d = 3$ ,  $\text{keep-prop} = 0.8$   
a boolean array

$$d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) \Rightarrow \text{eliminating chance} = 0.2$$

$$a_3 = \text{np.multiply}(a_3, d_3)$$

!  $a_3 \neq \text{keep-prop}$

why?

Imagine you have 50 hidden units at layer 3

$\Rightarrow$  on average at each iteration, you shut-off 10 units

$\Rightarrow z^{[4]} = w^{[4]} a^{[3]} + b^{[4]}$  will be reduced by  $1/20$  via value

$\Rightarrow$  in order not to change the expected value of  $a^{[3]}$  we need to bump it up by 20%.

$\Rightarrow$  hence dividing by keep-prop.

$\Rightarrow$  inverted dropout

! At each iteration of GD, you zero out a different set of units.

! At Test time, prediction:

No Dropout.

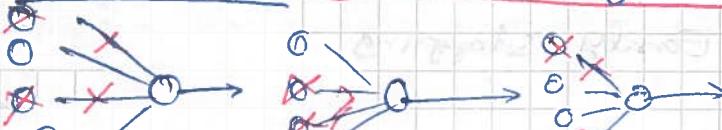
Because you don't want your output to be random.  
This will add noise to predictions.

## Why does Dropout work?

• it seems to be crazy randomly knocking out ~~the~~ units in the network.

• why does it work as a regularizer?

Intuition: can't rely on any one feature. So has to spread out the weights.

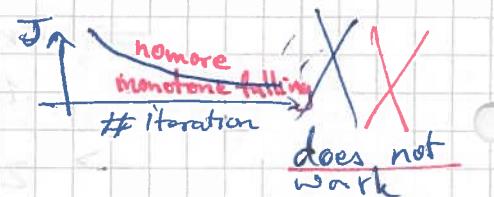


$\Rightarrow$  shrinking  $L_2$  norm of weights.

(~~so small~~)  $\rightarrow$   $L_2$  norm

← ◉ In fact Dropout is very similar to L<sub>2</sub> Regularization.

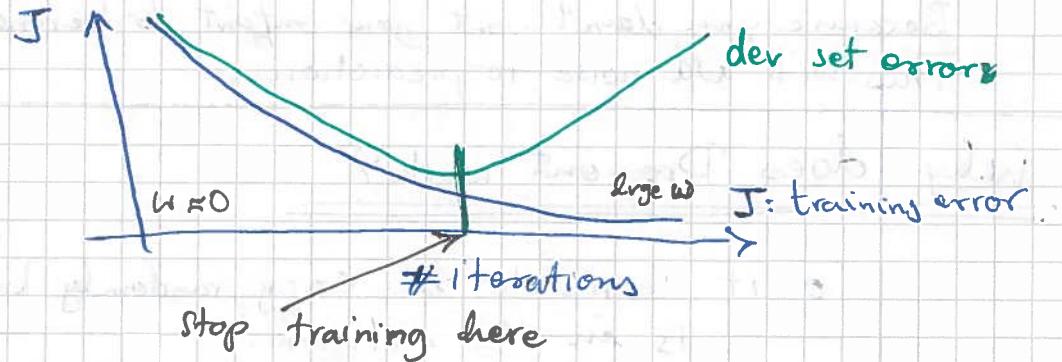
- ◉ You can also have different values for keep-prop for different layers. If a layer is too big ( $\Rightarrow$  may overfit) make keep-prop lower.
- ◉ Normally you don't dropout input layer. (keep-prop = 1)
- ◉ Dropout is often used in computer vision, where input size is too big ( $\Rightarrow$  you almost never have enough data)
- ◉ But does not always generalize to other domains
- ◉ Down-side of Dropout:  
Cost func( $J$ ) is no longer well-defined  $\Rightarrow$ 
  - $\Rightarrow$  first turn-off drop-out
    - $\Rightarrow$  make sure your  $J$  is monotonically decreasing \*
    - $\Rightarrow$  then turn on drop-out



## Other Regularization Methods :

- Data Augmentation !
  - ↳ e.g. flipping images horizontally
  - ↳ random distortion & transformations

### - Early Stopping



⇒ this works like L<sub>2</sub> regularization

$$\text{mid-size } w : \|w\|_F^2$$

Down-side of Early stopping



- ML is a process comprising of several steps:

{ - Optimizing the cost function  $J$   
 { - Gradient descent. —  
 { - RMS. ---  
 { - Not overfitting  
 { - regularization

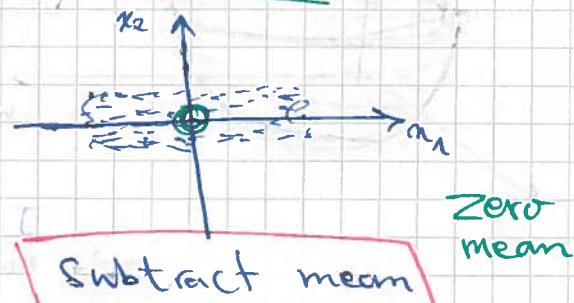
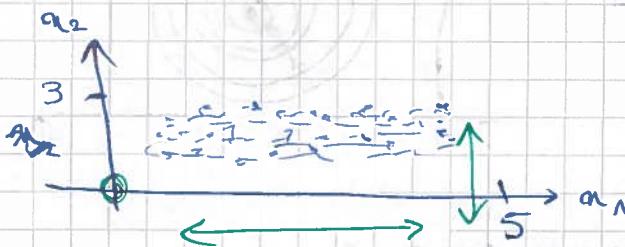
- ML is easier if we focus on each problem/step separately (orthogonalization)

- Early stopping mixes the two tasks. It makes the set of tools to use more complicated.
- An alternative can be using Regularization using different values of  $\lambda$

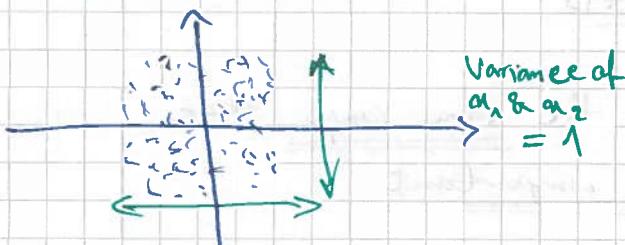
## Setting up the Optimization Problem :

### Normalizing Inputs :

- a technique to speed-up the training



### Normalize Variance



$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} \otimes 2$$

$$X \equiv \sigma^2 \text{ for all samples}$$

Subtract mean

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x^{(i)} := x^{(i)} - \mu \text{ for all the samples.}$$

! Use the same  $\mu, \sigma^2$

to normalize your training set and test set

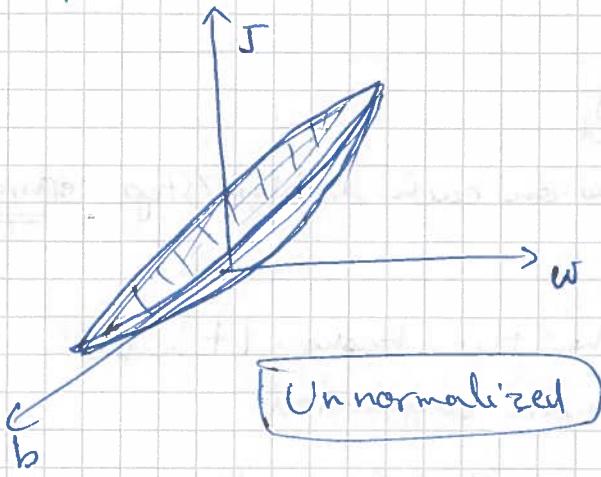
features have very different range

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

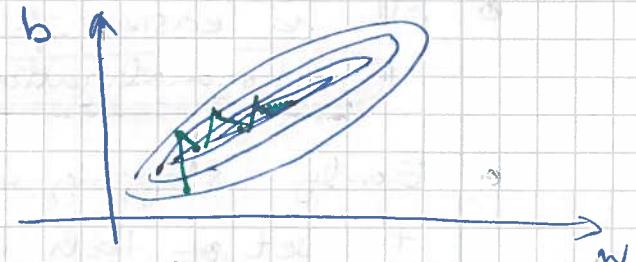
$m_1: 1 \dots 1000$

$m_2: 0 \dots 1$

$\Rightarrow w_1 \text{ & } w_2 \text{ will have very different values}$



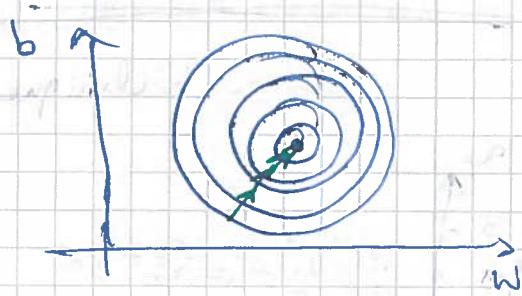
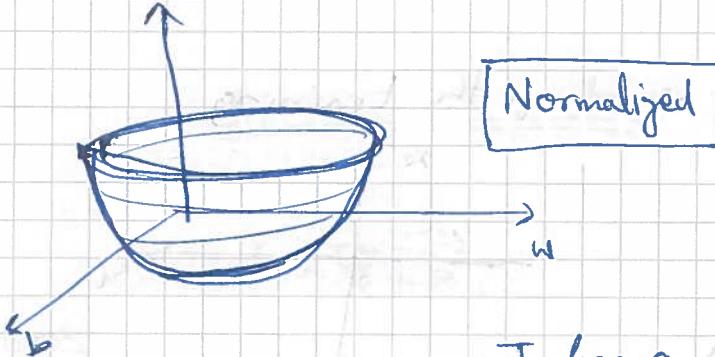
$\Rightarrow$  the cost function will be very elongated bowl



$\Rightarrow$  GD takes much longer to find the minimum  
 ↳ have to use a very small learning rate

Whereas Normalization:

- features become very symmetric

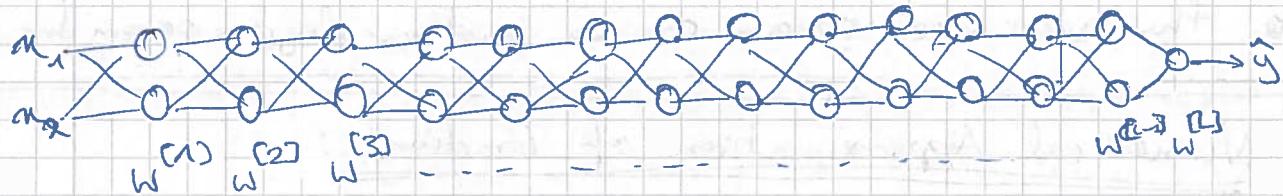


$J$  has a  
 $\Rightarrow$  more spherical contours  
 $\Rightarrow$  you can choose larger learning rate  
 $\Rightarrow$  faster GD

- If features have roughly the same range then normalization becomes less important.

## Vanishing / Exploding Gradients:

- Here an example with weights  $w$ , but same argumentation is also valid for derivatives.



for simplicity, assume  $g(z) = z$  (linear activation) &  $b = 0$

$$\Rightarrow \hat{y} = w^{[L]} w^{[L-1]} w^{[L-2]} \dots w^{[3]} w^{[2]} w^{[1]} x$$

- now imagine each  $w^{[l]}$  slightly larger than  $I$  (identity matrix)

$$w^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

$$\Rightarrow \hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x$$

Weights grow exponentially!! !

- on the other hand imagine

$$w^{[l]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$



Weights decrease exponentially!! !

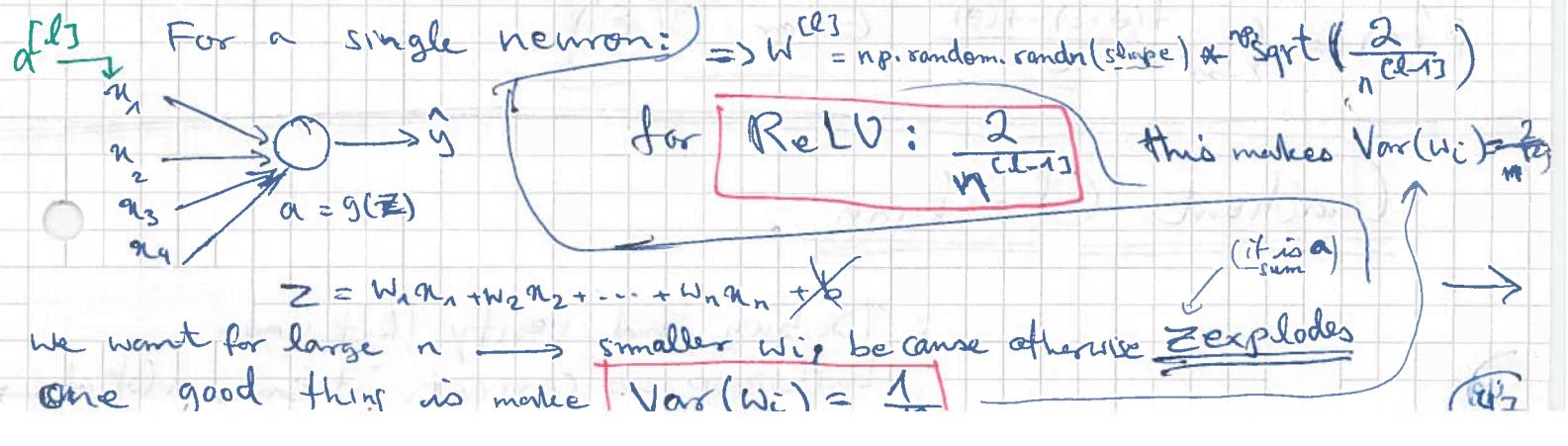
- Imagine now you had  $L=150$  layers!

- this problem was for a long time a huge barrier to train very deep networks.

- there is a partial solution:

careful choice of initial weights

## Weights Initialization Deep Networks

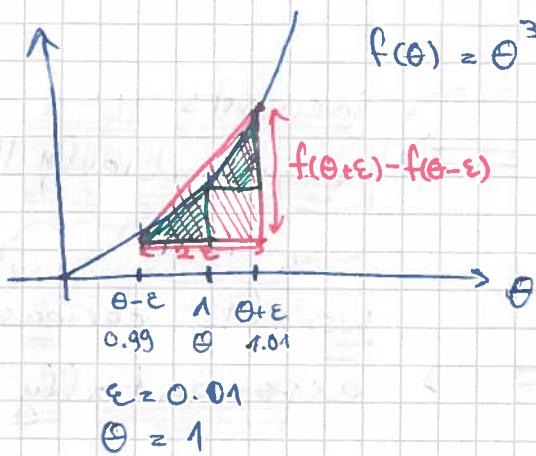


- if using a tanh activations  $\sqrt{\frac{1}{n^{(d-1)}}$
- another one:  $\sqrt{\frac{2}{n^{(d-1)} + n^{(d)}}$
- this variance param can be another hyper-param for model

Xavier  
Initialization

## Numerical Approximation of Gradient:

- **Gradient Checking** can help make sure that the impl of Backprop is correct.
- First let's see how to numerically approximate the gradient.



- the larger triangle gives a much better approximation of derivative of  $f$ .

- two sided difference: derivative

$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta) = 3\theta^2 = 3$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

approx. error: 0.0001

approx. error: 0.03 !

Whereas :

- one sided difference

$$\hookrightarrow \cancel{\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}} = 3.0301$$

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$$

Error:  $O(\epsilon^2)$   
 $\epsilon = 0.01 \rightsquigarrow$  Error: 0.0001

whereas :

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\epsilon)}{\epsilon}$$

Error:  $O(\epsilon)$

Two-sided difference derivation is much more accurate.

## Gradient Checking:

- Debug and verify that your backprop is correct. It saves a lot of time.

- Take  $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$  and reshape into a big vector  $\theta$   
 $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$   
concatenate  
 $J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$
- Same for  $dw^{[1]}, db^{[1]}, \dots, dw^{[L]}, db^{[L]}$ : reshape into a big vector  $d\theta$   
 $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$   
concatenate

Question is now:

Is  $d\theta$  the gradient of  $J(\theta)$ ?  
(slope)

Gradient Checking

Implementation of Grad Check:

$$J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots)$$

for each  $i$ :

$$d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \theta_2, \dots, \underbrace{\theta_i + \epsilon}, \dots) - J(\theta_1, \theta_2, \dots, \underbrace{\theta_i - \epsilon}, \dots)}{2\epsilon} \quad \text{rest remains the same}$$

this should be  $\approx \frac{\partial J}{\partial \theta_i}$

Compute  $d\theta_{\text{approx}}$  for all  $i$ 's  $\Rightarrow$  at the end you have two vectors and need to check:

Is  $d\theta_{\text{approx}} \approx d\theta$ ?

Use the Euclidean distance

Check:  $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} \approx \begin{cases} 10^{-7} & \rightarrow \text{great! :)} \\ 10^{-5} & \rightarrow \text{a little concerned} \\ 10^{-3} & \rightarrow \text{Worry} \\ 10^0 & \rightarrow \text{double check everything} \\ 10^3 & \rightarrow \text{most probably a bug somewhere (as)} \end{cases}$

just normalize by length of both vectors (Just for the case the vectors are very different in size)

## Grad Check implementation notes :

- Don't use in training - Only for debug.
  - $d\theta_{approx}[i]$  for all  $i$  is very costly
  - compute the derivatives and  $d\theta$  during backprop. only when debugging turn on grad check and compute  $d\theta_{approx}[i]$  for checking
- If algorithm fails grad check, look at components to identify the bug.
  - ~~value~~  $db_{cls}$  → very far off in  $d\theta_{approx}$
  - but  $dW_{cls}$  → ok.
  - ⇒ problem in  $db_{cls}$ .
- Don't forget regularization in definition of  $J$  if turned on
  - $d\theta$  : grad of  $J$  w.r.t.  $\theta$
- Doesn't work with dropout.
  - implement without dropout, then turn on dropout grad check when you are sure every thing is ok.
- Maybe (sometimes): Run at random init and then again after some training.

## Course 2 - Improving Deep NN - Week 1 - Programming 1

### Goal:

- Understand different regularization methods
- Implement dropout
- Realize the role of regularization: Without regularization the model performs better on training set but worse on test set.
- Try / Understand both dropout and regularization on your model.

### Initialization

#### Initialization:

A well chosen initialization ~~can speed up~~

- can speed up the convergence of GD
- increase the odds that GD converges to a lower ~~loss~~ training and generalization error.

- ④ zero init.
- ⑤ random init.
- ⑥ **He init**: random values scaled according to a paper by He et al., 2015

- ⑦ Zero init:

- ⑧ networks **fails** to break **symmetry**
- ⑨ each neuron in each layer computes the same thing as if an NN with  $n^{(l)} = 1$  for all layers
- ⑩ no more powerful than a **linear classifier**.
- ⑪ The weights  $W^{(l)}$  should be initialized randomly to break symmetry
- ⑫ It is ok to init  $b^{(l)}$  to zeros.

- ⑬ Random init with large values:

- ⑭ cost starts very high. This is because large random weights make the last activation (sigmoid) output very **near to 1 or 0**; when that example is wrong, the **Loss is very high**. Indeed when  $\log(a^{(s)}) = \log(0)$  the loss goes to **infinity**.

"inf" because  
"numerical round off"

- ⑮ Poor init leads to vanishing/exploding grads, which slows down the optimization.
- ⑯ Initializing weights to very large random values does not work
- ⑰ But how small should the values be?

⑱ He init: similar to Xavier init  $(\sqrt{\frac{1}{\text{layer-dims}(l-1)}})$

$$\sqrt{\frac{2}{\text{layers-dims}(l-1)}}$$

dimension of previous layer.

- ⑲ Recommended for ReLU layers
- ⑳ very good accuracy.

model	train accuracy	problem / comment
3-layer NN with zero init	50%	fails to break symmetry
3-layer NN with large random init	83 %	too large weights
3-layer NN with He init	99 %	recommended method

- Different init leads to different results
- Random init is used to break symmetry and make sure different hidden units can learn different things
- Don't initialize to values that are too large.
- He initialization works well for networks with ReLU activation.

## Course 2 - Improving Deep NN - Week 1 - Programming 2

### Regularization

- Deep learning models have so much flexibility and capacity that overfitting can become a serious problem if training data is not big enough.
- It will perform well on the training data set, but the learned Network does not generalize well on new examples.

- Baseline model, Trainset accuracy: 94.8%  
Test set : 91.5% !

- The base line non-regularized mode overfits the training set.  
It is fitting the noisy points 
- L2 Regularization to avoid overfitting & Modify the normal cost function

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

to

$$J_{\text{regul.}} = \underbrace{-\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(a^{[L]^{(i)}}) + (1-y^{(i)}) \log(1-a^{[L]^{(i)}}))}_{\text{cross-entropy cost}} + \underbrace{\frac{1}{m} \frac{\lambda}{2} \sum_{l=1}^L \sum_{k=1}^{n_l} \sum_{j=1}^{n_{l+1}} w_{kj}^{[l+1]2}}_{\text{L2 Regularization cost}}$$

- We must now change the backprop as well:  
All the gradients must be computed w.r.t. new cost function

- To each  $dW_1, dW_2, dW_3$  we have to add the gradient of the regularization term.

$$\left( \frac{d}{dW} \left( \frac{1}{2} \frac{\lambda}{m} W^2 \right) \right) = \frac{\lambda}{m} W$$

- ← ⚫ After <sup>L2</sup> regularization: Train set accuracy: 93.83%. Test set accuracy: 93 %. } !
- ( $\lambda = 0.7$ )
- The model is not fitting to noise anymore



## Notes

- the value of  $\lambda$  is a hyperparameter that we can tune using a dev set
- L2 Regularization makes the decision boundary smoother. If  $\lambda$  is too large, it is also possible to "oversmooth" resulting in a model with high bias
- L2 Regularization relies on the assumption that models with smaller weights is simpler than a model with larger weights. Thus, by penalizing the square values of the weights in the cost function you drive all the weights to smaller values. It becomes too costly for the model to have large weights. This leads to a smoother model, where the output changes more slowly as the input changes.
- The implications of the L2 Regularization are:
  - The cost function: A Regularization term is added to the cost.
  - The backprop function: There are extra terms in gradients w.r.t. weight matrices.
  - Weights end up smaller ("Weight decay") Weights are pushed to smaller values.

## Dropout

- randomly shut down some neurons in each iteration, with the probability (1 - keep-prob)
- The dropped out neurons don't contribute to training in both the forward and backward propagation of iteration.
  - We actually modify the model at each iteration and train a different model that uses only a subset of neurons. This way, the neurons become less sensitive to activation of one other specific neuron, because that other neuron might be shut-down at any time.

$$\text{Impl: } D^{[l]} = [d^{[l](1)} \ d^{[l](2)} \ \dots \ d^{[l](m)}]$$

$d^{[l]}$  same shape as  $A^{[l]}$   
 $d^{[l]}$  mask vector  
 $d^{[l]} = np.random.rand(n^{[l]}, 1)$   
 $d^{[l]} := (d^{[l]} < (1 - \text{keep-prop}))$

- Mask matrix with the same shape as  $A^{[l]}$
- For each iteration, and for each training sample, we shut down a different set of neurons of layer  $l$ .
- $A^{[l]} := A^{[l]} * D^{[l]}$  (Shutting down some neurons)
- $A^{[l]} := A^{[l]} / \text{keep-prop}$  (Assure that the result of the cost ~~will~~ will have still the same expected value as without Dropout. ("inverted dropout")

### ④ Back prop:

- In forward pass ~~we~~ we shut down some neurons using  $D^{[l]}$ .  
 In backward pass, we should shut down exactly the same neurons by reapplying  $D^{[l]}$  to  $dA^{[l]}$
- Since  $A^{[l]}$  was scaled by keep-prop in forward pass, we have to scale  $dA^{[l]}$  by the same value in the backward pass

- With dropout: Trainset accuracy: 92.8% (keep-prop = 0.86) Test set accuracy: 95% ! Dropout works great :)
- Dropout only in Training ! Not in testing !

model	train accuracy	test accuracy
3-layer NN without regularization	7.95	91.5 %
3-layer NN with L2 regularization	1.94	1.93
3-layer NN with Dropout	1.93	1.95

## Course 2 - Improving deep NNs - Week 1, programming 3

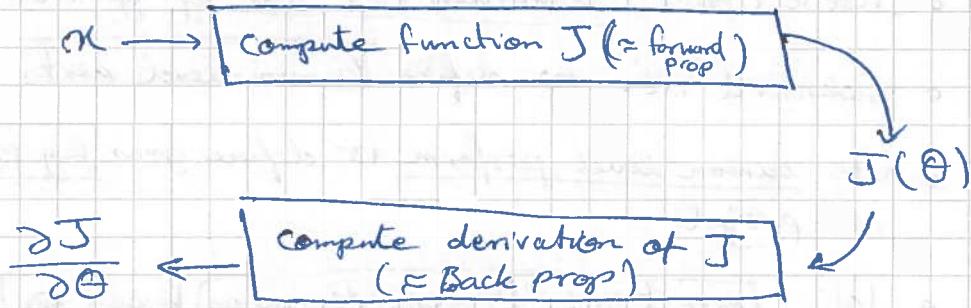
### Gradient Checking

- implement Gradient checking from scratch.
- understand how to use differences formula to check the backprop impl
- Back prop and the differences formula should give the same results.
- Find out which common abs is Gradient error committed

- Backprop computes gradients  $\frac{\partial J}{\partial \theta}$ , where  $\theta$  denotes the parameters of the model.  $J$  is computed using forward prop and the loss function.
- We use the Forward prop to verify the implementation of backprop. Because forward prop can be implemented easier and we assume it does not have ~~a~~ a bug.

$$\frac{\partial J}{\partial \theta} = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \approx \frac{J(\theta^+) - J(\theta^-)}{2\epsilon} \approx \frac{J^+ - J^-}{2\epsilon}$$

### 1-Dimensional Grad check :

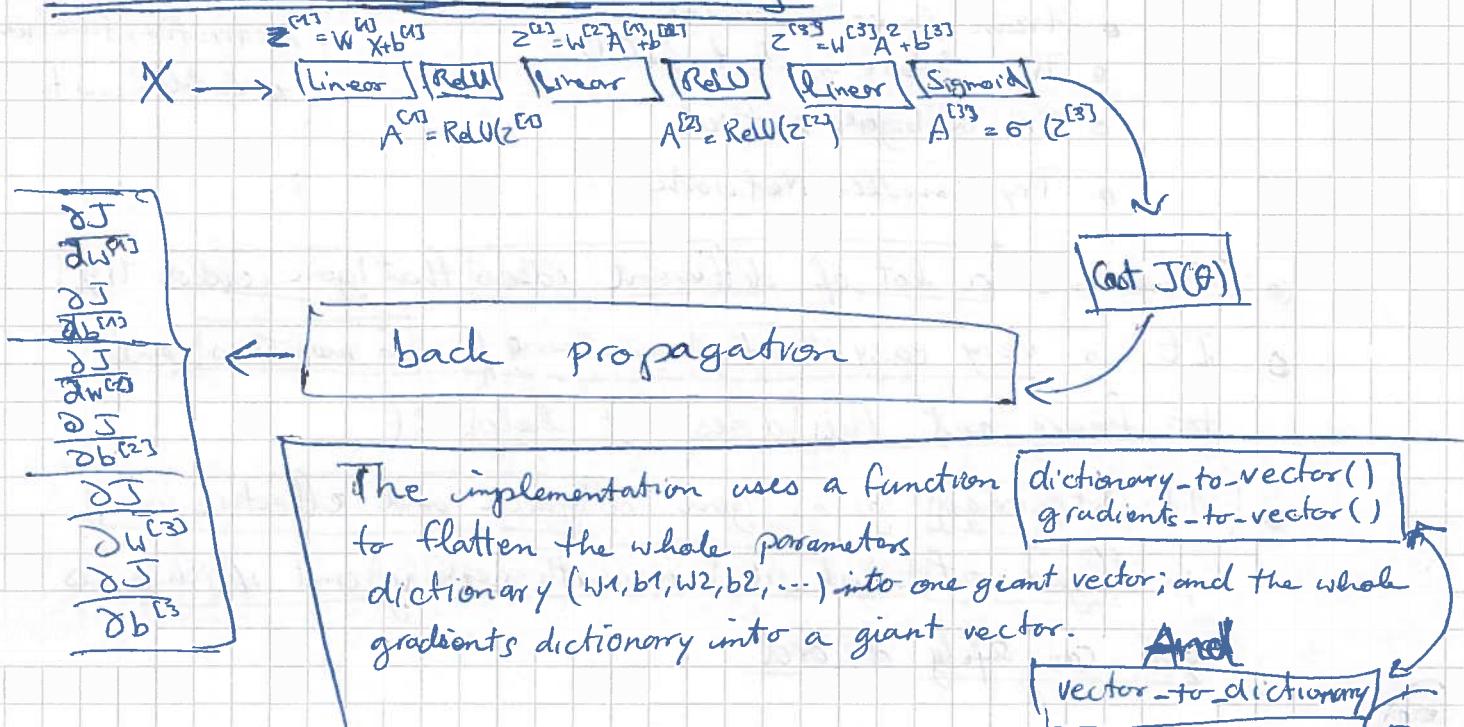


1. Compute gradapprox =  $\frac{J^+ - J^-}{2\epsilon}$
2. Compute grad using back prop.
3. Compute the relative difference between gradapprox and grad:

$$\text{difference} = \frac{\| \text{grad} - \text{gradapprox} \|_2}{\| \text{grad} \|_2 + \| \text{gradapprox} \|_2}$$

4. if difference is small, say  $< 10^{-7}$  then you are ok.

### N-Dimensional grad checking :



# Course 3 - C3 Structuring ML Projects

## Week 1

### ML Strategy (1)

#### Learning Objectives:

- Understanding why a ML strategy is important.
- Apply satisfying and optimizing metrics to set up your goal for ML project
- Choose correct train/dev/test split of your dataset
- Understand how to define human-level performance
- Use human level perform to define your key priorities in ML project
- Take correct ML strategic decisions based on observations of performances and dataset

#### Introduction : Motivation

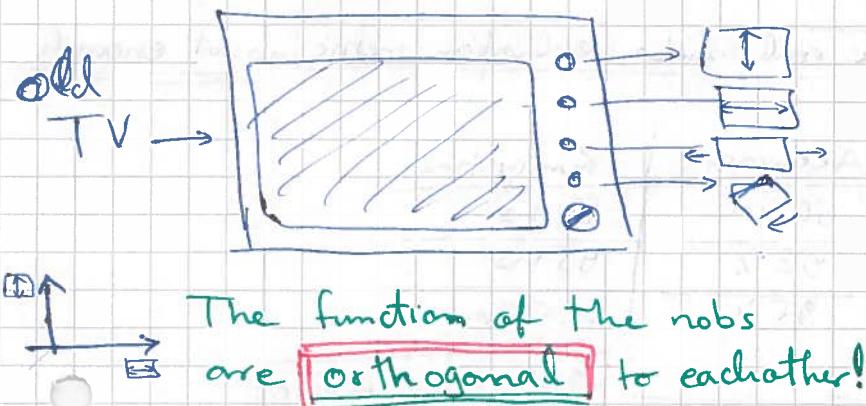
- Imagine you cat classifier achieves 90% accuracy, but this is not enough.
- What to do? you could:
  - collect more data
  - collect more diverse training set
  - train longer with GD
  - Try Adam instead of GD
  - Try a bigger Network
  - Try smaller Network
- There are a lot of different ideas that you could try.
- It is very easy to try something for 6 months only to figure out this does not help. :(
- ML Strategy gives you a quick and effective way to figure out which idea is worth pursuing and which ones you can safely discard.

- Try dropout
- Try L2 Regularization
- Change Network Architecture
  - Activation function
  - # hidden units

## Orthogonalization:

### Challenge:

- In ML System you could try many different things, such as hyperparameters.
- Having a very clear view on "what to change" in order to achieve an effect.



• Each nob has a distinct function

• If you had ~~one~~ one single nob that did:

$$(0.1 \times \downarrow + 0.5 \rightarrow + 1.7 \leftarrow + 0.2 \square)$$

It would be impossible or very hard to tune the tv.

## Chain of Assumptions in ML:

- Fit training set well on cost function.  
(≈ human-level performance)
- ↓
- Fit dev set well on cost function.
- ↓
- Fit test set well on cost function.
- ↓
- Perform well in real world.  
(Happy users)

- nobs to tune
  - bigger network
  - Adam
  - Regularization
  - Bigger training set
  - Bigger dev set
  - Change either dev set or the cost function

- Each nob/function addresses a distinct/separate problem
- Early Stopping could be a bad nob. Because it simultaneously effects two things:
  - How good you fit training set (you stop before you reach the best training fitting)
  - Improve dev set performance.

## Setting up your goals: Single number evaluation metric

- Having 1 single real number evaluation metric speeds up iteration and helps choosing algorithms easier.

$$\underline{\text{Precision}} = \frac{\text{TP}}{\text{TP} + \text{FP}} \%$$

$$\underline{\text{Recall}} = \frac{\text{TP}}{\text{TP} + \text{FN}} \%$$

$$\boxed{\text{F1-score}} = \frac{2}{\frac{1}{P} + \frac{1}{R}} : \text{Average of precision and recall} \\ (\text{harmonic mean})$$

## Satisficing and Optimizing metrics:

Single Real number evaluation metric

- Some times just one **single real number evaluation metric** is not enough

Example:

Classifiers	Accuracy	Running time
A	90 %	80 ms
B	92 %	95 ms
C	95 %	1,500 ms

{ maximize Accuracy ← Optimizing metric  
 subject to running time ≤ 100 ms ← Satisficing metric  
 ↳ (Just has to be good enough)

Generally: If you have N metrics: 1 Optimizing  
 N-1 Satisficing

Another example: Wake Words / Trigger words

↳ Interested in accuracy & # false positive.  
 ⇒ { maximize accuracy ← Optimizing  
 s.t. ≤ 1 FP in every 24 hours. ← Satisficing  
 ↳ (at most)

## Train / dev / Test distributions:

- The way we set up train / dev / test set can have a huge impact on the progress. It can also slowdown progress dramatically.
- make sure your dev & test sets come from "the same" distribution  
 ↳ you can for example randomly shuffle the data into dev / test sets

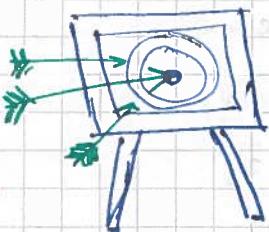
True Story: Wasted 3 month of work:



- Optimized on dev set on loan approval for medium income zip codes
- tested on low income zip codes.

- Choose the dev set and test set
  - from the same distribution
  - and reflect the data you expect to get in the future and consider important to do well on.

training set  
helps you  
hit well.  
**the target**



{ dev set } / { test set }

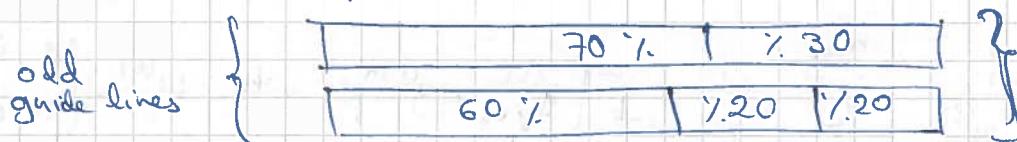
+ evaluation metric

Defines the target!



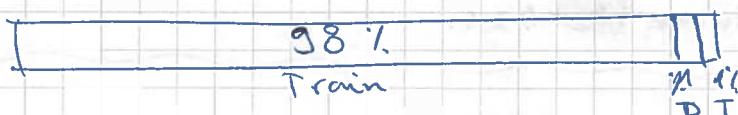
### Size of the dev and test sets :

- guidelines of how to setup your dev and test sets are changing in the DL era.
- old way of splitting data



Small DS sizes  
1000 samples  
10 000 samples

- In DL era we have much larger DSs



large DS  
1,000,000 samples

- Size of your test set :
  - test set gives you ~~confidence on how~~
  - good your final system is after you have developed it.
- Big enough to give you high confidence in the overall performance of your system.

- Some times you don't even need a high confidence in the overall performance of the final system  
 $\Rightarrow$  all you need is train/dev set



- Some time test / dev is used exchangeably.

- Not recommended not having a test set

## When to change dev/test sets and metrics?

Cat classifier for example:

→ Metric classification error

Algorithm A: 30% error

Algorithm B: 5% error



But lets pornographic images wrongly classified as cat!

⇒ { metric + Dev set : prefers A  
you / users : prefers B

⇒ We must change the error metric, because it does not reflect our preference.

for example: If error metric was indicator func.

$$\text{Error} = \frac{1}{m_{\text{dev}}} \sum_{i=1}^{m_{\text{dev}}} I \{ y_{\text{pred}}^{(i)} \neq y^{(i)} \}$$

Just simply count the number of misclassifications

we change it to

$$\text{Error} = \frac{1}{\sum_i w^{(i)}} \sum_{i=1}^{m_{\text{dev}}} w^{(i)}$$

$I \{ y_{\text{pred}}^{(i)} \neq y^{(i)} \}$  with  $w^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ non-porn} \\ 10 & \text{if } x^{(i)} \text{ porn} \end{cases}$

normalize  
so that, the error remains between 0 & 1.  
penalize wrong porn misclassification

## Orthogonalization:

1. First define a metric to evaluate the two classifiers  
 2. Then worry separately about how to do well on this metric

- First define the metric
- then figure out how to achieve it

1. Define the target

2. Worry about achieving it well SEPARATELY

Guideline:

If doing well on your metric + dev/test set does not correspond to doing well on your application, change your metric and/or dev/test set.



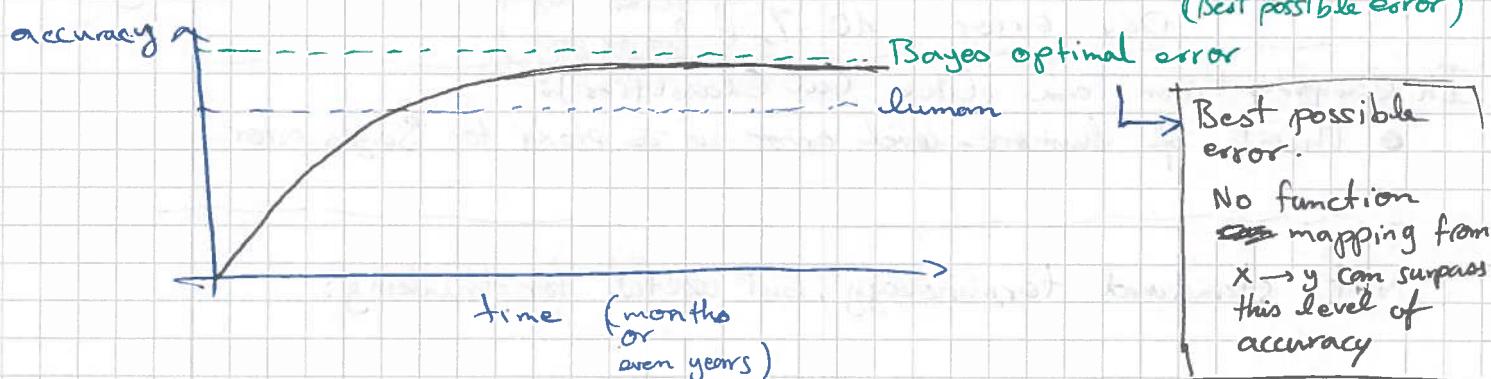
to  
speed  
up  
progress

- set up a dev/test set and evaluation metric quickly and then iterate.
- do not work for too long without an evaluation metric and dev set.

## Comparing to human-level Performance:

### Why human level Performance?

- Because of advances in Deep learning, ML is working much much better and has even become competitive to human performance in some areas.
- The workflow for setting up ML solution has become very efficient, so it is now natural to think of human-level performance.



- The progress is normally fast, as long as we haven't surpassed human-level performance. But after that the progress becomes very slow. It has two reasons:
  - for many tasks the human-level performance is already very near to Bayes optimal error.
  - After passing human-level performance, certain tools can hardly, or no longer be used to improve performance:
    - Getting labeled data from humans.
    - Gain insights from manual error analysis. (Why did a person get this right?)
    - Better analysis of bias/Variance.

## Avoidable Bias

- We always want our training algorithm to do well on the training set. But sometimes we don't want it to do too well.
- Knowing the human-level performance can give us a good limit on ~~how well~~ we want our algorithm to perform on the training data. → But not too well.

### Cat Classification Example:

Human error	1%	big gap	⇒ Focus on reducing bias.
Training error	8%		
Dev error	10%		⇒ Probably high <u>bias</u> .

Human error	7.5%	relatively small gap	⇒ Probably high <u>Variance</u>
Training error	8%		
Dev error	10%		⇒ Focus on reducing variance.

In Computer Vision Tasks (Like Cat Classification):

- Think of human-level error as a proxy for Bayes error.

Not standard terminology, but useful for thinking:



- By using the term "Avoidable Bias" we actually indicate that there is a level of error that we can't or basically don't want to surpass.

So, in above example we can say "Avoidable Bias is 0.5%!"

### Understanding Human-level performance:

- Human-level error as a proxy for Bayes error.



Example of a medical image classification Task

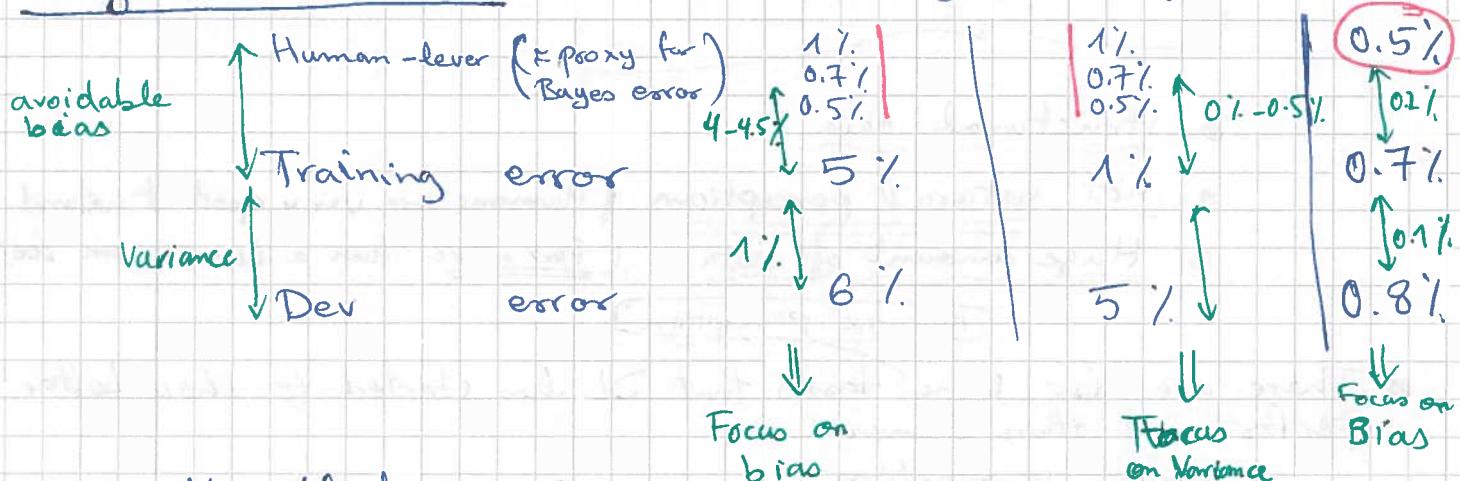
- |                       |            |  |
|-----------------------|------------|--|
| a) untrained human    | 3% error   | b) typical doctor 1% error                       |
| c) experienced doctor | 0.7% error | d) consensus minimum of a team of expert doctors |

← Now with these values, what is the value for "human-level error"?

From d)  $\Rightarrow$  Bayes error  $\leq 0.5\%$ , error

$\Rightarrow$  We define human-level error (as a proxy for Bayes error) to be 0.5%

Why this matters? → Error Analysis Example.



- ① in the third example it really matters to set the human-level error to 0.5%. Because otherwise, (for example if you have set Human-level error to 0.7%) you would have missed the fact that you need to work on bias and would have worked (wasted time) on variance instead.

- ② Interestingly, as we approach to human-level performance (third example) it becomes harder to figure out what to do to improve performance.

### Surpassing human-level performance:

Team of humans

One human

Training error

Dev. error

0.5%

~~0.1%~~

0.6%

0.2%

0.8%

0.5%

1%

0.3%

0.4%

Relatively  
easy to figure  
out what to  
do.

now what is the  
avoidable bias?

What is Bayes error?

{0.1%?  
0.2%?  
0.3%?}

- In second example we cannot know what to do. What to focus on?

~~reducing bias? or reducing the variance?~~

- Moreover, as we have surpassed the human-level performance

↳ now it's time to find out how to improve.

- Problems where ML significantly surpasses human-level performance:

- online advertising
- product recommendation
- Logistics (predicting transit times)
- loan approval

### • Structural data

- not natural perception (Humans are very good at natural perception task)
- Huge amount of data (Far more than a human can see/process)

(natural perception)

• There are also some tasks that DL has started to show better performance than humans:

- speech recognition
- some image recognition tasks
- some medical tasks, such as ECG, ...

## Improving your model performance:

The two fundamental assumptions of supervised learning:

1. You can fit the training set pretty well.

≈ Achieving low avoidable bias.

2. The training set performance generalizes pretty well to the dev/test sets.

≈ Achieving good low variance.

## Reducing (avoidable) bias and Variance:

Human-level

↑  
available  
bias →

Training error

↑  
Variance →

Dev error

- Train a bigger model
- Train longer/ better optimization algorithm
  - ↳ momentum, RMSProp, Adam, ...
- NN architecture / hyperparameter search
  - ↳ RNN
  - ↳ CNN
  - ↳ # layers
  - ↳ # hidden units
- More data
- Regularization
  - ↳ L<sub>2</sub>, Dropout, data augmentation
- NN architecture / hyperparameter search.

## From Quiz:



It is not a problem to have different training and dev distribution.

On the other hand, it would be very problematic to have different dev and test set distribution.

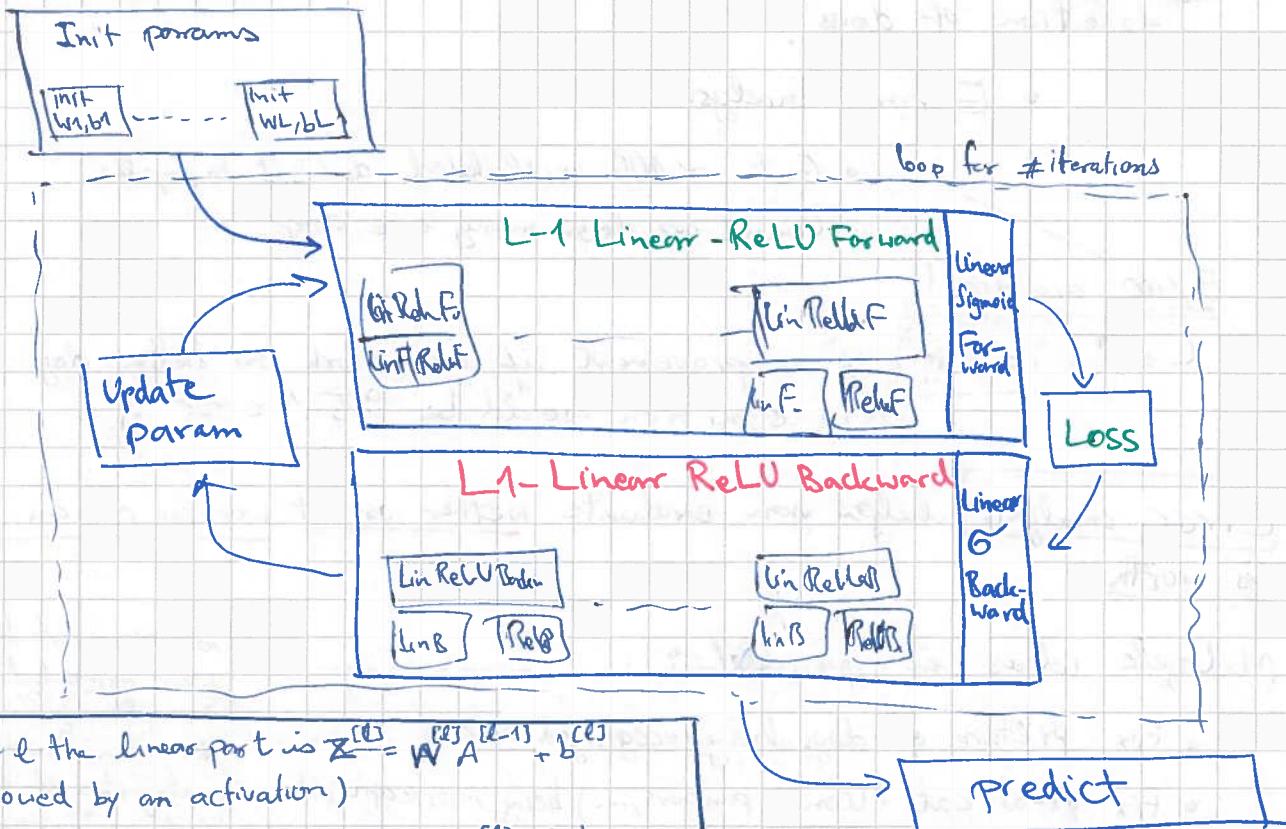
- training distribution  $\neq$  dev distribution O.K.
- dev distribution  $\neq$  test distribution Not O.K.

## Course 1 - NNs and Deep Learning - Week 4 - Programming

Building your deep Neural Network Step-by-step:

- develop an intuition for the overall structure of a neural Network
- write functions (e.g. forward prop, backward prop, logistic loss, etc...) that help you decompose your code and ease the process of building a neural Network.
- Initialize / update parameters according to your desired structure

First a  
2-layer



For layer  $l$  the linear part is  $Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$   
(followed by an activation)

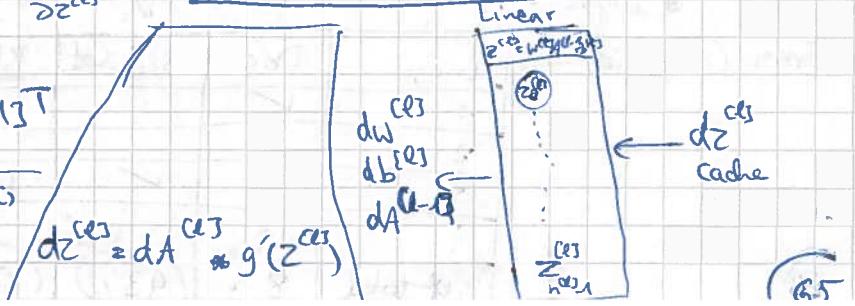
Suppose you have already calculated  $dZ^{[l]} = \frac{\partial L}{\partial Z^{[l]}}$

You want to get  $(dW^{[l]}, db^{[l]}, dA^{[l-1]})$

$$dW^{[l]} = \frac{\partial L}{\partial W^{[l]}} = \frac{1}{m} dZ^{[l]} A^{[l-1]T}$$

$$db^{[l]} = \frac{\partial L}{\partial b^{[l]}} = \frac{1}{m} \sum_{i=1}^m dZ^{[l](i)}$$

$$dA^{[l-1]} = \frac{\partial L}{\partial A^{[l-1]}} = W^{[l]T} dZ^{[l]}$$



# Course 3 - Structuring ML Projects - Week 2

## Week 2

### ML Strategy 2:

#### Error Analysis: Carrying out Error Analysis

- If you want your ML to do a task that humans can do, but your performance is still not as good as you wished,

(Manually examining the errors in your classification.

can help giving you intuition on what to do next.

Example:

Cat Classifier: Look at "Dev Examples" to evaluate ideas.

→ Not working good: { accuracy: 90%  
error: 10% }

→ mislabeled.

→ some of dog pictures are misclassified as cats. Should you work on your classifier to improve detection of dogs?

#### Error Analysis

- Get ~100 mislabeled dev set examples
- Count up how many are dogs.

5/100 are dogs!

→ "Cieling" of improvement if you work on better dog recognition would be 9.5% error =

- Error analysis helps you evaluate whether or not working on an idea is worth.

- Multiple ideas in parallel:

- Fix picture of dogs being recognized as cats
- Fix great cats (lions, panthers,...) being misrecognized.
- Improve performance on blurry images

As you work through your mislabeled dev examples, you may even encounter new categories of error, such as "Instagram filters".

Error Categories →

Image	Dog	Big Cat	Blurry	Instagram	Comment
1	✓				Pitbull
2			✓		
3		✓	✓		Rainy dog at zoo
:	?	?	?	?	
% Mislabeled	8%	42%	61%	12%	