

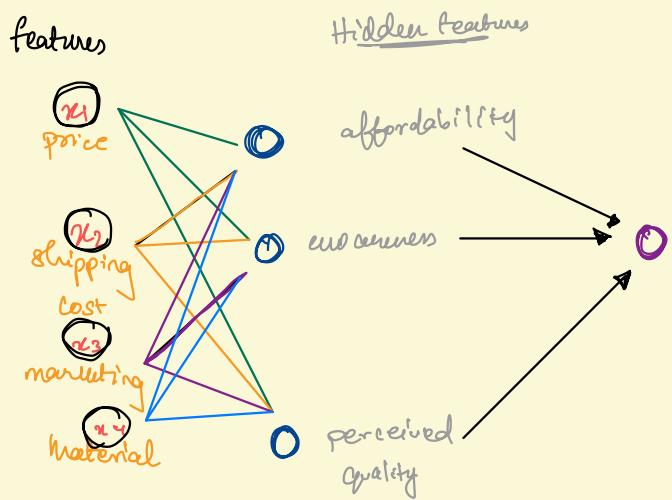
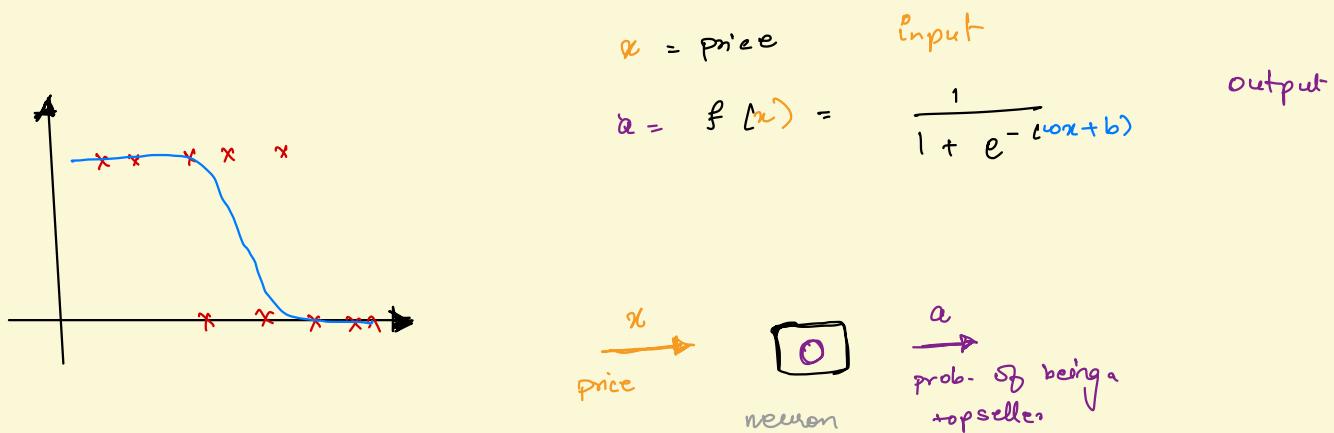
## Topics - covered

Source:  
NPTEL: Nitin Khapra  
Deep learning.  
coursera Andrew NG  
CodingCrane - YouTube  
Towards Data Science - blog

- ① Neural Network - Introduction & Analogies
- ② Representations & Notations Shrouded
- ③ Feed Forward & Back propagation derivation (optional) and Algorithm.
- ④ Activation functions
- ⑤ Example 1: Implementing XOR gate with a Neural Network.
- ⑥ Example 2: MNIST dataset

- Neural Networks: Algorithms that involves a network of neurons & trains on data to predict the correct set of parameters to predict the correct o/p.

- Demand Prediction:



L1:  
Inputs 4 numbers  
Oips 3 numbers

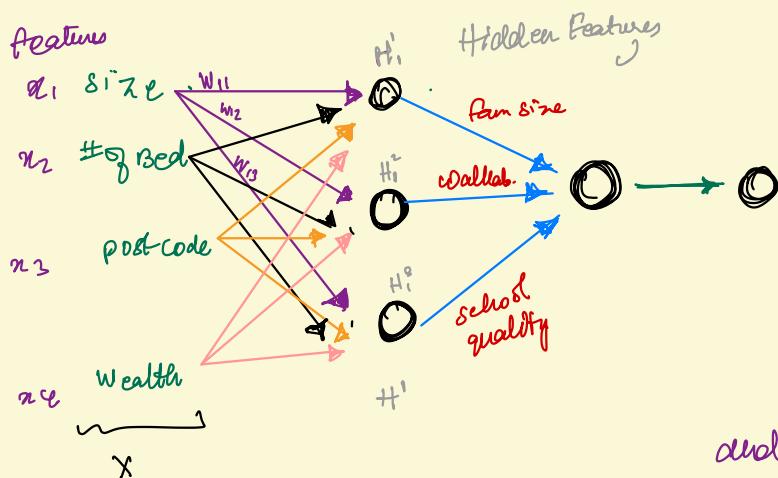
L2:  
Iops 3 n.o.s  
Olp 1

Speciality of NN:  
One of nice prop. of NN is when we train NN from data, we don't have to explicitly decide what other features such as affordability, awareness & so on.  
The ANN figures out all if based on training data

A neural network as it trains automatically extracts the hidden features.

## \* Housing Price Prediction

Price based on size, # of bedrooms, postal code.



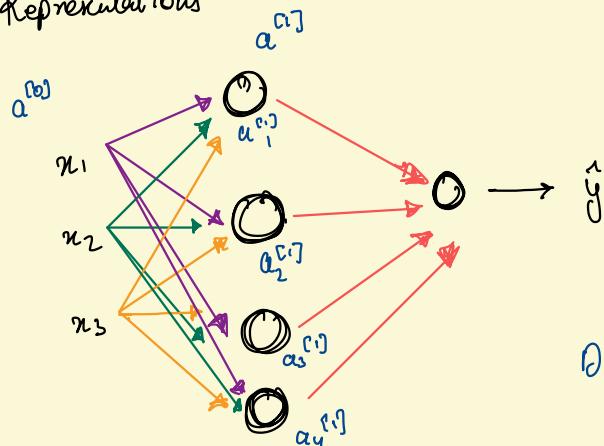
Dlp layer connected with every hidden layer.

Instead of user specifying that the feature term size will depend on feature  $x_1$  &  $x_2$ , we provide all 4 features to family size for the neural net to decide on which feature it will depend upon.

and how does it decide which features will it depend on? (By understanding from training set it decides what weights  $w_{11}, w_{12}, w_{13}$  are!)

We can have another hidden layer (which means what else parameters it might depend on like direction of the house, if it has balcony etc etc)

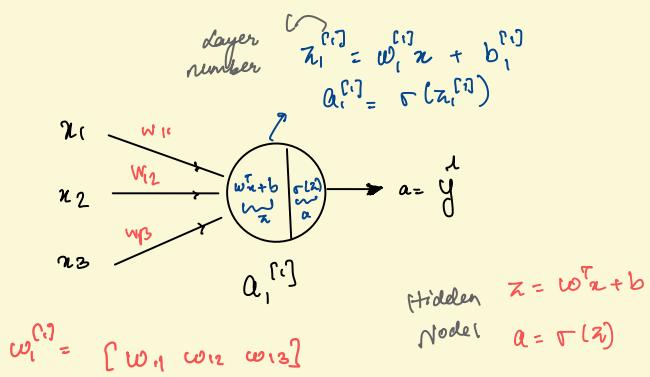
## \* Neural N/w Representations



$$a^{[l]} = \begin{bmatrix} a_1^{[l]} \\ a_2^{[l]} \\ \vdots \\ a_n^{[l]} \end{bmatrix}$$

2 Layer NN

Don't count the dlp layer



Parallelly,

$$\begin{aligned} z_2^{[l]} &= w_2^{[l]} x + b_2^{[l]} \\ a_2^{[l]} &= \sigma(z_2^{[l]}) \\ z_3^{[l]} &= w_3^{[l]} x + b_3^{[l]} \\ a_3^{[l]} &= \sigma(z_3^{[l]}) \end{aligned}$$

↳ Vectorization:

$$\begin{aligned} z_1^{[l]} &= w_1^{[l]} x + b_1^{[l]} \\ z_2^{[l]} &= w_2^{[l]} x + b_2^{[l]} \\ z_3^{[l]} &= w_3^{[l]} x + b_3^{[l]} \\ z_4^{[l]} &= w_4^{[l]} x + b_4^{[l]} \end{aligned}$$

For four nodes in the hidden layer.

$$z^{[l]}$$

$$\begin{bmatrix} w_{11}^{[l]} & w_{12}^{[l]} & w_{13}^{[l]} \\ w_{21}^{[l]} & w_{22}^{[l]} & w_{23}^{[l]} \\ w_{31}^{[l]} & w_{32}^{[l]} & w_{33}^{[l]} \\ w_{41}^{[l]} & w_{42}^{[l]} & w_{43}^{[l]} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[l]} \\ b_2^{[l]} \\ b_3^{[l]} \\ b_4^{[l]} \end{bmatrix} = \begin{bmatrix} w_1^{[l]} x + b_1^{[l]} \\ w_2^{[l]} x + b_2^{[l]} \\ w_3^{[l]} x + b_3^{[l]} \\ w_4^{[l]} x + b_4^{[l]} \end{bmatrix} = \begin{bmatrix} z_1^{[l]} \\ z_2^{[l]} \\ z_3^{[l]} \\ z_4^{[l]} \end{bmatrix}$$

$$a_i^{[l]} = \sigma(z_i^{[l]}) \quad a^{[l]} = \begin{bmatrix} a_1^{[l]} \\ a_2^{[l]} \\ a_3^{[l]} \\ a_4^{[l]} \end{bmatrix}$$

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

$\downarrow$

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

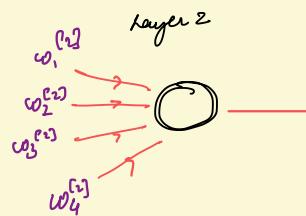
$4$  nodes in hidden layer each with  $3$  inputs

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

$\downarrow$

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

$1$  node in the present layer with  $4$  inputs

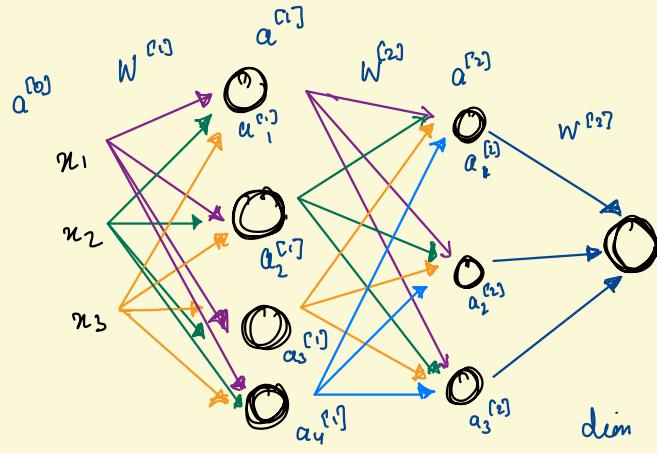


Each  $w_1^{[2]}, w_2^{[2]}, w_3^{[2]}, w_4^{[2]}$  is a scalar.

$$w^{[2]} = [w_1^{[2]} \ w_2^{[2]} \ w_3^{[2]} \ w_4^{[2]}]$$

$a_j^{[l]} = j^{\text{th}} \text{ op of } l^{\text{th}} \text{ layer}$

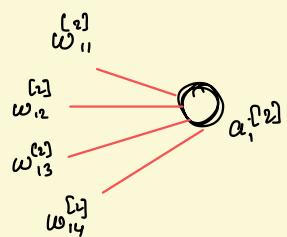
$$a_j^{[l]} = \sigma(z_j^{[l]}) = \sigma(\vec{w}_j^{[l]} \vec{a}^{[l-1]} + b_j^{[l]})$$



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

Aim: To find  $a^{[2]}$

:



$$a_1^{[2]} = \sigma(z_1^{[2]})$$

$$z_1^{[2]} = [w_{11}^{[2]} \quad w_{12}^{[2]} \quad w_{13}^{[2]} \quad w_{14}^{[2]}] \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} + b_1^{[2]}$$

$$w_1^{[2]} = \begin{bmatrix} w_{11}^{[2]} \\ w_{12}^{[2]} \\ w_{13}^{[2]} \\ w_{14}^{[2]} \end{bmatrix}$$

$$\therefore z_1^{[2]} = w_1^{[2] T} a^{[1]} + b_1^{[2]}$$

$$z_2^{[2]} = w_2^{[2] T} a^{[1]} + b_2^{[2]}$$

$$z_2^{[2]} = w_2^{[2] T} a^{[1]} + b_2^{[2]}$$

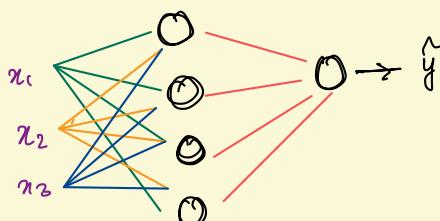
$$z^{[2]} = \begin{bmatrix} z_1^{[2]} \\ z_2^{[2]} \\ z_3^{[2]} \end{bmatrix} = \begin{bmatrix} -w_1^{[2]} \\ -w_2^{[2]} \\ -w_3^{[2]} \end{bmatrix} \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} + \begin{bmatrix} b_1^{[2]} \\ b_2^{[2]} \\ b_3^{[2]} \end{bmatrix} = \underbrace{\begin{bmatrix} w_{11}^{[2]} & w_{12}^{[2]} & w_{13}^{[2]} & w_{14}^{[2]} \\ w_{21}^{[2]} & w_{22}^{[2]} & w_{23}^{[2]} & w_{24}^{[2]} \\ w_{31}^{[2]} & w_{32}^{[2]} & w_{33}^{[2]} & w_{34}^{[2]} \end{bmatrix}}_{W^{[2]}} \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} + \begin{bmatrix} b_1^{[2]} \\ b_2^{[2]} \\ b_3^{[2]} \end{bmatrix}$$

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

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

4 inputs, 3 layers = Inputs x # of layers

## VECTORIZATION



$$x^{(1)} = \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{bmatrix}$$

Training examples

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

$$x^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)}$$

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

for  $i = 1 \text{ to } m$ :  $\# \text{ of samples}$

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

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

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

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



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

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

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

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

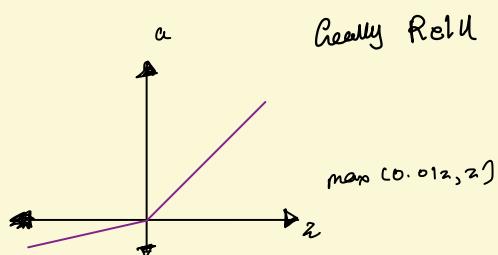
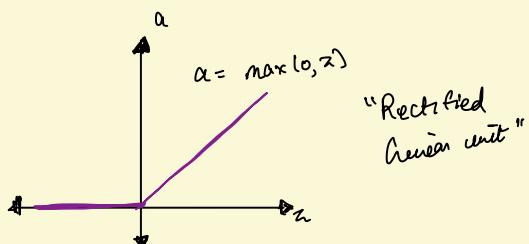
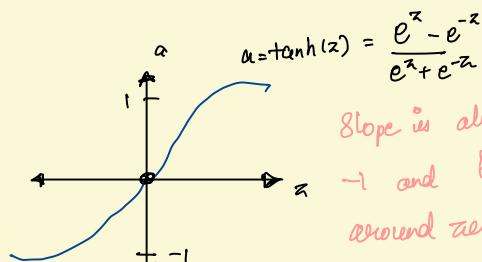
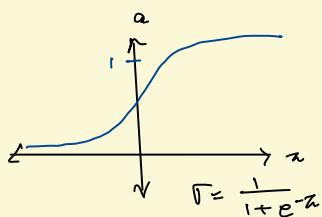
$$X = A^{[0]}$$

$$A^{[0]} = \left[ \begin{array}{c} \text{Sample} \\ a_{1}^{[0](1)} \\ a_{1}^{[0](2)} \\ \vdots \\ a_{1}^{[0](n)} \end{array} \right]$$

$$A^{[1]} = \left[ \begin{array}{c} a_{1}^{[1](1)} \\ a_{1}^{[1](2)} \\ a_{1}^{[1](3)} \\ a_{1}^{[1](4)} \end{array} \right]$$

$$X = \left[ \begin{array}{c} x_{1}^{(1)} \\ x_{1}^{(2)} \\ x_{1}^{(3)} \\ x_{1}^{(4)} \\ \vdots \\ x_{1}^{(m)} \end{array} \right]$$

### \* Activation Functions



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

$g(z) = z$  "Linear activation function"

Aim: To compute Weights and biases

### Neural Network Classification:

#### Binary Classification

i) OLP  $y$  is 0 or 1

e.g.

Check if digit is 5 or 6

OLP:  $OLP = 0$  it's 5

$OLP = 1$  it's 6

#### Multiclass Classification

ii) OLP  $y \in \mathbb{R}^k$

Classify digits from 0 to 9

$$OLP = \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \end{bmatrix}$$

Digit is 1

In summary NN:

- Data:  $\{x_i, y_i\}_{i=1}^N$
- Model:  $\hat{y}_i = f(x_i) = \Theta [W^{[0]} + (W^{[1]} \sigma (W^{[0]} x + b^{[1]}) + b^{[2]}) + b^{[3]}]$   
 $\Theta$  = some output func.  
 acting on  $W^{[3]} A^{[2]} + b^{[2]}$
- Parameters:  $\Theta = W^{[0]}, W^{[1]}, \dots, W^{[n]}, b^{[1]}, b^{[2]}, \dots, b^{[n]}$   
 $n = \# \text{ of hidden layers} + 1$
- Algorithm: Gradient descent with Back propagation.

Aim: To find an algorithm for learning parameters of the above model.

\* Now that we have managed to find an o/p with the given initial weights and biases, How do we improve upon the weights?

With the given weights we have predicted the o/p. Comparing it with the original o/p's, we need to find weights & biases s.t. the error is minimum

\* Backpropagation: Algorithm used to improve the weights & biases using gradient descent to minimize the errors.

Gradient descent:

$t = 0$

max-iter = 1000

repeat

calculate  $\frac{\partial J}{\partial w_i}$ ,  $\frac{\partial J}{\partial b}$

$$w_i = w_i - \eta \frac{\partial J}{\partial w_i}, \quad b_i$$

$$b_i = b_i - \eta \frac{\partial J}{\partial b}$$

3

Strategy: Model

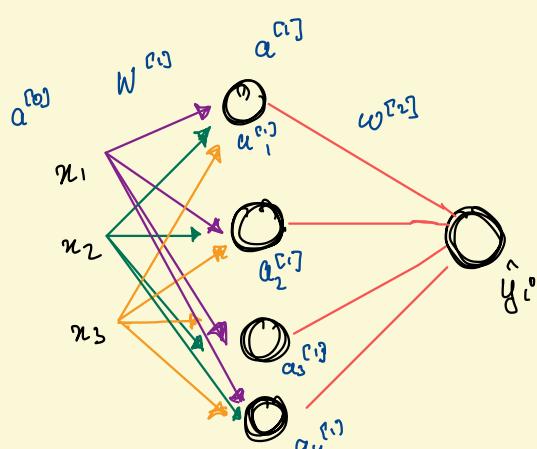


loss function



Find Parameters using Grad

Descent



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

$$\begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{22} \\ w_{13} & w_{23} & w_{23} \\ w_{14} & w_{24} & w_{24} \end{bmatrix}$$

We have seen we choose loss functions based on the problem at hand.

## \* BACKPROPAGATION USING GD:

Cross Cost func :

$$\text{Cost } C = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) \quad \text{where } \hat{y}_i = \sigma(W^{(l)} A^{(l-1)} + b^{(l)})$$

$L$  = Loss function

$$w = w - \alpha \frac{\partial C}{\partial w}$$

Minimize  $C$

$w, b$

$$b = b - \alpha \frac{\partial C}{\partial b}$$

Aim: To find  $\frac{\partial C}{\partial w}$  and  $\frac{\partial C}{\partial b}$  in each layer. to update the weights in each layer.

## Backpropagation:

Layer 2

$$z = W^{(2)} A^{(1)} + b \quad A^{(2)} = \sigma(z) \quad \frac{\partial z}{\partial w^{(2)}} = A^{(1)T}$$

$$C = \frac{1}{m} \sum_{i=1}^m L(y_i, A_i^{(2)})$$

$$\frac{\partial C}{\partial A^{(2)}} \text{ is given.}$$

$$\frac{\partial C}{\partial z^{(2)}} = \frac{\partial C}{\partial A^{(2)}} \frac{\partial A^{(2)}}{\partial z^{(2)}}$$

$$\bullet \frac{\partial C}{\partial w^{(2)}} = \frac{\partial C}{\partial A^{(2)}} \frac{\partial A^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial w^{(2)}} = \frac{\partial C}{\partial A^{(2)}} \frac{\partial A^{(2)}}{\partial z^{(2)}} A^{(1)T} = \frac{\partial C}{\partial z^{(2)}} A^{(1)T}$$

$$\bullet \frac{\partial C}{\partial b^{(2)}} = \frac{\partial C}{\partial A^{(2)}} \frac{\partial A^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial b^{(2)}} = \frac{\partial C}{\partial z^{(2)}}$$

$$\therefore W^{(2)} = W^{(2)} - \alpha \frac{\partial C}{\partial w^{(2)}} = W^{(2)} - \alpha \frac{\partial C}{\partial z^{(2)}} A^{(1)T}$$

$$b^{(2)} = b^{(2)} - \alpha \frac{\partial C}{\partial b^{(2)}} = b^{(2)} - \alpha \frac{\partial C}{\partial z^{(2)}}$$

# Updation of  $W^{(2)}$

Layer 1

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

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

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

$$\frac{\partial C}{\partial A^{[2]}} = \frac{\partial C}{\partial Z^{[2]}} \frac{\partial A^{[2]}}{\partial Z^{[2]}} \frac{\partial Z^{[2]}}{\partial A^{[2]}} = (W^{[2]})^T \frac{\partial C}{\partial Z^{[2]}}$$

$$\frac{\partial C}{\partial Z^{[1]}} = \frac{\partial C}{\partial A^{[1]}} \cdot \sigma'(Z^{[1]})$$

$$\frac{\partial C}{\partial W^{[1]}} = \frac{\partial C}{\partial A^{[1]}} \frac{\partial A^{[1]}}{\partial Z^{[1]}} \frac{\partial Z^{[1]}}{\partial W^{[1]}} = \frac{\partial C}{\partial Z^{[1]}} A^{[1]}$$

$$\frac{\partial C}{\partial b^{[1]}} = \frac{\partial C}{\partial Z^{[1]}} \frac{\partial Z^{[1]}}{\partial b^{[1]}} = \frac{\partial C}{\partial Z^{[1]}}$$

$$W^{[1]} = W^{[1]} - \alpha \frac{\partial C}{\partial W^{[1]}} = W^{[1]} - \alpha \frac{\partial C}{\partial Z^{[1]}} A^{[0]T}$$

$$b^{[1]} = b^{[1]} - \alpha \frac{\partial C}{\partial b^{[1]}} = W^{[1]} - \alpha \frac{\partial C}{\partial Z^{[1]}}$$

\* Summarizing, we have,

for each epoch, the backpropagation involves calc. of weights & biases of each layer. i.e. from o/p to i/p.

From forward prop. we have value of  $A^{[l]}$ . and  $C$  is a func. dependip on  $A^{[l]}$ .

Algorithm:

layer = outermost layer  $l$ ,  $\frac{\partial C}{\partial A^{[l]}}$  is given.

repeat  $l \geq 1$

$$\frac{\partial C}{\partial Z^{[l]}} = \sigma'(Z^{[l]}) \frac{\partial C}{\partial A^{[l]}}$$

$$\frac{\partial C}{\partial W^{[l]}} = \frac{\partial C}{\partial Z^{[l]}} (A^{[l-1]})^T$$

$$\frac{\partial C}{\partial b^{[l]}} = \frac{\partial C}{\partial Z^{[l]}}$$

$$W^l = W^l - \alpha \frac{\partial C}{\partial W^l}$$

$$b^l = b^l - \alpha \frac{\partial C}{\partial b^l}$$

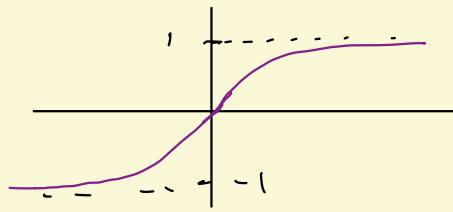
} update step.

$$\frac{\partial C}{\partial A^{[l-1]}} = (W^{[l]})^T \frac{\partial C}{\partial Z^{[l]}}$$

$$l = l - 1 \quad \{$$

Activation functions:

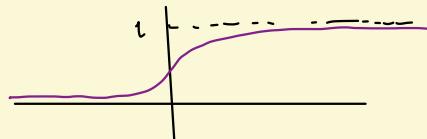
① Tanh x :  $\frac{e^x - e^{-x}}{e^x + e^{-x}}$



derivative:

$$\begin{aligned} \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2} &= 1 - \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2} \quad (\text{separating denominators}) \\ &= 1 - (\tanh x)^2 \end{aligned}$$

② Sigmoid:  $\frac{1}{1 + e^{-x}}$



derivative :

$$\begin{aligned} \frac{-(-e^{-x})}{(1 + e^{-x})^2} &= \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1 + e^{-x} - 1}{(1 + e^{-x})^2} = \frac{1}{(1 + e^{-x})} - \frac{1}{(1 + e^{-x})^2} \\ &= \text{sig} (1 - \text{sig}) \end{aligned}$$

Back propagation:

1)  $\frac{\partial C}{\partial A^l}$  is given

Calculate  $\frac{\partial C}{\partial Z^l} = \sigma'(z^l) \frac{\partial C}{\partial A^l}$

$\sigma'(z^l)$  element wise  
multipl.

2)  $\frac{\partial C}{\partial W^l} = \frac{\partial C}{\partial Z^l} \cdot (A^{l-1})^T$

$\frac{\partial C}{\partial b^l} = \frac{\partial C}{\partial Z^l}$

3)  $W^l = W^l - \alpha \frac{\partial C}{\partial W^l}$   
 $b^l = b^l - \alpha \frac{\partial C}{\partial b^l}$

4) Update

$$\frac{\partial C}{\partial A^{l-1}} = (W^l)^T \cdot \frac{\partial C}{\partial Z^l}$$

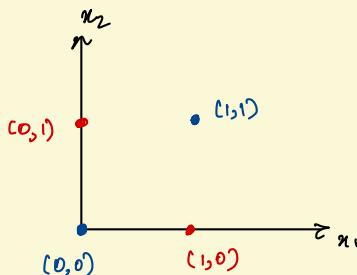
$$f^{(l)} = \sigma(z^{(l)})$$

Hence,  $\frac{\partial A^{(l)}}{\partial z^{(l)}} = \sigma'(z^{(l)})$

# \* Neural Network Implementation:

The XOR Gate:

with a single perceptron, it is not possible find a linear hyperplane.



clearly the points are not linearly separable.  
Why?

This is an example of an XOR gate.

$x_1$	$x_2$	X OR
0	0	0
1	0	1
0	1	1
1	1	0

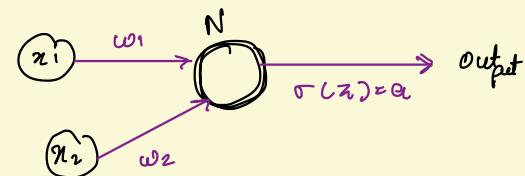
$$\begin{aligned} w_0 + \sum_{i=1}^2 w_i x_i &< 0 \\ w_0 + \sum_{i=1}^2 w_i x_i &\geq 0 \\ w_0 + \sum_{i=1}^2 w_i x_i &\geq 0 \\ w_0 + \sum_{i=1}^2 w_i x_i &< 0 \end{aligned}$$

- ①  $w_0 + w_1(0) + w_2(0) < 0 \Rightarrow w_0 < 0$
- ②  $w_0 + w_1(0) + w_2(1) \geq 0 \Rightarrow w_2 \geq -w_0$
- ③  $w_0 + w_1(1) + w_2(1) \geq 0 \Rightarrow w_1 \geq -w_0$
- ④  $w_0 + w_1(1) + w_2(0) > 0 \Rightarrow w_1 + w_2 < -w_0$

2nd & 3rd cond' contradict the 4th condition.

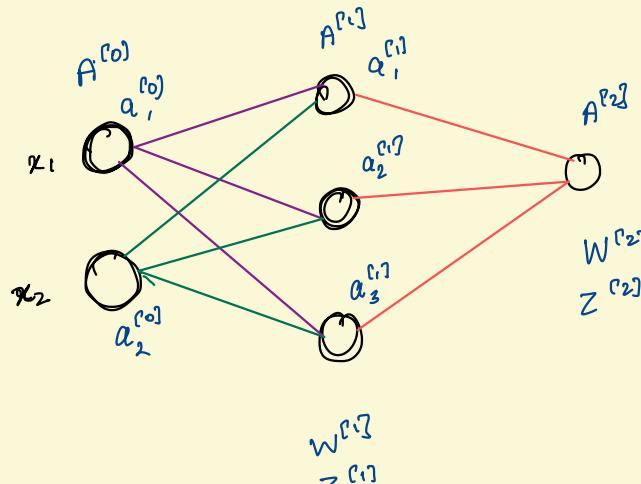
∴ We can't have a sol'n to this set of inequalities.

∴ it is possible to find a solution with a network of perceptrons.



No such  $w_0, w_1, w_2$   
s.t output of TV =  
output of XOR gate

- ① Input : 2 features  $x_1, x_2$  and 4 samples



In this example,  
layer 1 has 3 neurons  
layer 2 is output layer  
with one neuron

$$\dim W^{[l]} = (\# \text{ of neurons in layer } l, \# \text{ of inputs})$$

$$\dim b^{[l]} = (\# \text{ of neurons in layer } l, 1)$$

$$\begin{bmatrix} w_{11}^{[0]} & w_{12}^{[0]} \\ w_{21}^{[0]} & w_{22}^{[0]} \\ w_{31}^{[0]} & w_{32}^{[0]} \end{bmatrix} \begin{bmatrix} a_1^{[0]} \\ a_2^{[0]} \end{bmatrix} + \begin{bmatrix} b_1^{[0]} \\ b_2^{[0]} \\ b_3^{[0]} \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$W^{[0]}$

- \* Given training set with correct output, In each iteration or epoch, we try to find a set of weight & bias so that the error (diff b/w actual value & the one predicted by the set of weights) is minimum.
- \* In each epoch, we train the network with all training samples. So in one epoch, we will have to adjust the weights the # of sample times. This is inefficient.
- Hence, one epoch, we will try to input all the samples as a Batch.

\* Input:

Feature  $x_1$ : [ 0 0 + 1 ]

Feature  $x_2$ : [ 0 1 0 1 ]

training data set is same as testing data set.

$$x_{\text{train}} = \begin{bmatrix} [0 0 1 1] \\ [0 1 0 1] \end{bmatrix}$$

$$y_{\text{train}} = [ [0 1 1 0] ]$$

$x_1$	$x_2$	X OR
0	0	0
1	0	1
0	1	1
1	1	0

\* Defining Each layer through a class:

Each layer has 2 variables ① # of neurons in that layer

② # of inputs.

③ weights : dim: (# of neurons, # of inputs)

④ bias : dim: (# of neurons, 1)

## • Iterating through Epochs.

In each epoch following takes place:

① Feed forward.

Instead of  $z^{[0]} = W^{[0]}X + b^{[0]}$ ,  $X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \text{Single Input.}$ ,

We now have a batch of inputs.

$$X = \begin{bmatrix} [x_1] \\ [x_2] \end{bmatrix} = \begin{bmatrix} [x_1^1, x_1^2, \dots, x_1^m], \\ [x_2^1, x_2^2, \dots, x_2^m] \end{bmatrix} = A^{[0]} = \begin{bmatrix} [0, 0, 1, 1] \\ [0, 1, 0, 1] \end{bmatrix} \xrightarrow{\substack{\leftarrow 4 \text{ samples} \\ \uparrow 2 \text{ features}}}$$

Code flow

$$\textcircled{1} \quad z^{[0]} = W^{[0]} X + b^{[0]} \quad \begin{array}{l} (\text{# of neurons in layer}), (\text{# of neurons, # of inputs}) \\ (\text{# of features, # of samples}) \\ \# \text{ of samples} \end{array} \quad \begin{array}{l} b^{[0]} \\ (\text{# of neurons, 1}) \end{array}$$

$$z^{[0]} = \begin{bmatrix} z_1^{[0]} \\ z_2^{[0]} \\ z_3^{[0]} \end{bmatrix} = \begin{bmatrix} [z_{11}^{[0]}, z_{12}^{[0]}, \dots, z_{1m}^{[0]}] \\ [z_{21}^{[0]}, z_{22}^{[0]}, \dots, z_{2m}^{[0]}] \\ [z_{31}^{[0]}, z_{32}^{[0]}, \dots, z_{3m}^{[0]}] \end{bmatrix}$$

$$A^{[0]} = \sigma(z^{[0]}) = \begin{bmatrix} \{\sigma(z_{11}^{[0]}) \dots\} \\ \vdots \\ \{\sigma(z_{31}^{[0]}) \dots\} \end{bmatrix} \quad \begin{array}{l} (\text{# of neurons in layer 1}) \\ \# \text{ of samples} \end{array}$$

Note:

→ This addition is possible with np.array.

Each row of  $b^{[0]}$  gets added to all columns of same row of  $W^{[0]}$ .

Layer I has 3 neurons

$$W^{[0]} = \text{dim } (3, 2)$$

$$A^{[0]} = \text{dim } (3, 4)$$

$$\# \text{ of samples} = m = 4$$

$$\textcircled{2} \quad z^{[1]} = W^{[1]} A^{[0]} + b^{[1]} \quad \begin{array}{l} (1, m) \quad (1, 3) \quad (3, m) \quad (1, 1) \end{array}$$

$$W^{[1]} = \text{dim } (1, 3)$$

$$z^{[1]} = [z_1^{[1]}] = [\{z_{11}^{[1]}, z_{12}^{[1]}, \dots, z_{1m}^{[1]}\}]$$

$$A^{[1]} = \sigma(z^{[1]}) = [\{\sigma(z_{11}^{[1]}), \sigma(z_{12}^{[1]}), \dots, \sigma(z_{1m}^{[1]})\}]$$

$$\text{dim } A^{[1]} = (1, m)$$

— FORWARD PROP —

## ② Calculating the loss

The cost function is

$$C(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m (y_i \log(a_i) + (1-y_i) \log(1-a_i)) \quad \xrightarrow{\text{Summ over } m \text{ samples}} \text{loss for all samples } (x_1^{(i)}, x_2^{(i)}) \\ 1 \leq i \leq m$$

$$a_i = A^{[i]} \quad i = \text{last layer} = \text{predicted output}$$

In our case,

$$A^{[i]} = A^{[2]} = [\sigma(z_{11}^{(2)}), \sigma(z_{12}^{(2)}), \dots, \sigma(z_{1m}^{(2)})]$$

↓  
 predicted op of 1st sample      ↓  
 pred. op of 2nd sample      ↓  
 pred. op of mth sample

Actual op:  $y^{[i]} = [y_1, y_2, \dots, y_m]$

$$1 - y^{[i]} = [(1-y_1), (1-y_2), \dots, (1-y_m)]$$

$$\log A^{[2]} = [\log(\sigma(z_{11}^{(2)})), \dots, \log(\sigma(z_{1m}^{(2)}))]$$

$$\log(1 - A^{[2]}) = [\log(1 - \sigma(z_{11}^{(2)})), \dots, \log(1 - \sigma(z_{1m}^{(2)}))]$$

$$\begin{aligned}
 & - [y_i * \log A^{[2]}] + [(1-y_i) * \log(1 - A^{[2]})] = \\
 & - [y_1 * \log(\sigma(z_{11}^{(2)})), \dots, y_m * \log(\sigma(z_{1m}^{(2)}))] + [(1-y_1) * \log(1 - \sigma(z_{11}^{(2)})), \dots, (1-y_m) * \log(1 - \sigma(z_{1m}^{(2)}))] \\
 & = - [y_1 * \log(\sigma(z_{11}^{(2)})) + (1-y_1) * \log(1 - \sigma(z_{11}^{(2)})), \dots, y_m * \log(\sigma(z_{1m}^{(2)})) + (1-y_m) * \log(1 - \sigma(z_{1m}^{(2)}))]
 \end{aligned}$$

$$\text{Cost} = \text{np.sum}(\text{above array})$$

Now, have to update the weights through back prop-

$$\dim A^{[l]} = (\# \text{ of neurons in layer } l, m \text{ samples})$$

$$\dim Z^{[l]} = (\# \text{ of neurons in layer } l, m \text{ samples})$$

(i) calc.

$$\frac{\partial C}{\partial A^{[2]}}$$

$$C(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m (y^i \log(a_i^{[2]}) + (1-y^i) \log(1-a_i^{[2]}))$$

$$A^{[2]} = [a_1^{[2]}, a_2^{[2]}, \dots, a_m^{[2]}]$$

$$\frac{\partial C}{\partial A^{[2]}} = \left[ \frac{\partial C}{\partial a_1^{[2]}}, \frac{\partial C}{\partial a_2^{[2]}}, \dots, \frac{\partial C}{\partial a_m^{[2]}} \right]$$

$$\frac{\partial C}{\partial a_i^{[2]}} = -\frac{1}{m} \left( \frac{y^i}{a_i^{[2]}} - \frac{(1-y^i)}{1-a_i^{[2]}} \right)$$

$$\frac{\partial C}{\partial a_i^{[2]}} = -\frac{1}{m} \left[ \frac{y^i - a_i^{[2]} y^i - a_i^{[2]} + a_i^{[2]} y^i}{a_i^{[2]} (1-a_i^{[2]})} \right]$$

$$\frac{\partial C}{\partial a_i^{[2]}} = \frac{1}{m} \left( \frac{a_i^{[2]} - y^i}{a_i^{[2]} (1-a_i^{[2]})} \right)$$

$$\frac{\partial C}{\partial A^{[2]}} = \left[ \left[ \frac{1}{m} \left( \frac{a_1^{[2]} - y^i}{a_1^{[2]} (1-a_1^{[2]})} \right), \frac{(a_2^{[2]} - y^i)}{m a_2^{[2]} (1-a_2^{[2]})}, \dots, \frac{1}{m} \left( \frac{a_m^{[2]} - y^i}{a_m^{[2]} (1-a_m^{[2]})} \right) \right] \right]$$

Q How can this be computed?

For Vectorization:  $(A^{[2]} - y_{\text{train}}) \rightarrow \text{array dim: } (1, \# \text{ of samples})$  gives  $\uparrow$   
 $A^{[2]} \times (1 - A^{[2]}) \rightarrow \text{array dim: } (1, \# \text{ of samples})$

(ii) Backprop for layer 2

$$\frac{\partial C}{\partial Z^{[2]}} = \sigma'(Z^{[2]}). \frac{\partial C}{\partial A^{[2]}}$$

↓      ↓      ↓  
 $(1, \# \text{ of samples})$      $(1, \# \text{ of samples})$      $(1, \# \text{ of samples})$

$$\frac{\partial C}{\partial W^{[2]}} = \frac{1}{m} \cdot \frac{\partial C}{\partial Z^{[2]}} \cdot (A^{[1]})^T$$

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

$$\frac{\partial C}{\partial b^{[2]}} = \frac{1}{m} \cdot \left( \frac{\partial C}{\partial Z^{[2]}} \right) \text{ summing across columns} \rightarrow \text{Why?}$$

$(1, 1)$        $(1, 1)$

Algorithm:

layer = outermost layer  $l$ ,  $\frac{\partial C}{\partial A^{[l]}}$  is given

repeat  $l \geq 1$

$$\left\{ \begin{aligned} \frac{\partial C}{\partial Z^{[l]}} &= \sigma'(Z^{[l]}). \frac{\partial C}{\partial A^{[l]}} \\ \frac{\partial C}{\partial W^{[l]}} &= \frac{\partial C}{\partial Z^{[l]}} \cdot (A^{[l-1]})^T \left[ \frac{1}{m} \right] \\ \frac{\partial C}{\partial b^{[l]}} &= \frac{\partial C}{\partial Z^{[l]}} \end{aligned} \right.$$

$$\left. \begin{aligned} w^l &= w^l - \alpha \frac{\partial C}{\partial w^l} \\ b^l &= b^l - \alpha \frac{\partial C}{\partial b^l} \end{aligned} \right\} \text{Update step}$$

$$\frac{\partial C}{\partial A^{[l-1]}} = (W^{[l]})^T \frac{\partial C}{\partial Z^{[l]}}$$

$$l = l-1 \quad \exists$$

$$Z^{[2]} = (1, 4)$$

$$\sigma'(Z^{[2]}) = \underbrace{\text{sig}(Z^{[2]})}_{\downarrow} \underbrace{(1 - \text{sig}(Z^{[2]}))}_{\text{Elementwise}}$$

$$\left[ \sigma'(Z_1^{[2]}), \sigma'(Z_2^{[2]}), \dots, \sigma'(Z_n^{[2]}) \right]$$

$\downarrow$   
 $\text{sig}(Z_1^{[2]}) (1 - \text{sig}(Z_1^{[2]})) \rightarrow (1, \# \text{ of samples})$

$$\bullet A^{[1]} = (\# \text{ of neurons}, \# \text{ of samples})$$

Layer 1  
 $(3, 4)$

Update the weights

$$W^{[2]} = W^{[2]} - \alpha \frac{\partial C}{\partial W^{[2]}}$$

$$b^{[2]} = b^{[2]} - \alpha \frac{\partial C}{\partial b^{[2]}}$$

$$\frac{\partial C}{\partial A^{[1]}} = (W^{[2]})^T \cdot \left( \frac{\partial C}{\partial z^{[2]}} \right)$$

$\downarrow$

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

return  $\frac{\partial C}{\partial A^{[1]}}$

why  $(3,4)^T$

$$A^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \end{bmatrix} = \begin{bmatrix} [a_{11}^{[1]}, a_{12}^{[1]}, \dots, a_{1m}^{[1]}] \\ [a_{21}^{[1]}, a_{22}^{[1]}, \dots, a_{2m}^{[1]}] \\ [a_{31}^{[1]}, a_{32}^{[1]}, \dots, a_{3m}^{[1]}] \end{bmatrix}$$

We found  
need  $W^{[2]}$ ,  $b^{[2]}$

$$\frac{\partial C}{\partial A^{[1]}} = \begin{bmatrix} \left[ \frac{\partial C}{\partial a_{11}^{[1]}}, \frac{\partial C}{\partial a_{12}^{[1]}}, \dots, \frac{\partial C}{\partial a_{1m}^{[1]}} \right]^T \\ \left[ \frac{\partial C}{\partial a_{21}^{[1]}}, \dots, \frac{\partial C}{\partial a_{2m}^{[1]}} \right]^T \\ \left[ \frac{\partial C}{\partial a_{31}^{[1]}}, \dots, \frac{\partial C}{\partial a_{3m}^{[1]}} \right]^T \end{bmatrix}$$

(iii) Backpropagation for layer 1.

$\frac{\partial C}{\partial A^{[1]}}$  is calculated

$$\frac{\partial C}{\partial A^{[1]}} = (3,4)$$

$$\frac{\partial C}{\partial z^{[1]}} = \sigma'(z^{[1]}). \frac{\partial C}{\partial A^{[1]}}$$

$\downarrow$

$(3, \# \text{samples}) \quad (3, \# \text{samples}) \quad (3, \# \text{samples})$

$$A^{[0]} = (2, \# \text{samples})$$

$(2, 4)$

$$\frac{\partial C}{\partial W^{[1]}} = \frac{1}{m} \cdot \frac{\partial C}{\partial z^{[1]}} \cdot (A^{[0]})^T$$

$(3,4) \quad (4,2)$

$$W^{[1]} = (3, 2)$$

$$\frac{\partial C}{\partial b^{[1]}} = \frac{1}{m} \cdot \left( \frac{\partial C}{\partial z^{[1]}} \right) \text{ summing across columns}$$

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

Intuition:

for every sample we do,

$$b^{[2]} = b^{[2]} - \alpha \left[ \frac{\partial C}{\partial b^{[2]}} \right]$$

i.e. for each sample sample above happens if we don't vectorize

Vectorizing, it happens when we add all columns together

Update the weights

$$W^{[l]} = W^{[l]} - \alpha \frac{\partial C}{\partial W^{[l]}}$$

$$b^{[l]} = b^{[l]} - \alpha \frac{\partial C}{\partial b}$$

$$\frac{\partial C}{\partial A^{[0]}} = (W^{[1]})^T \cdot \left( \frac{\partial C}{\partial z^{[1]}} \right)$$

$\downarrow$                      $\downarrow$   
 $(2,4)$                  $(2,3)$                  $(3,4)$

return  $\frac{\partial C}{\partial A^{[0]}}$

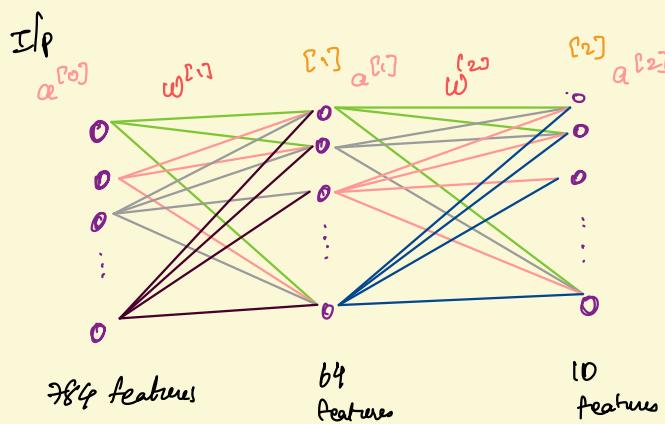
Hence after the number of epochs we have the minimized weights & biases.

## \* MNIST dataset Example

Problem: Build a neural network to classify digits from 0 to 9

Data Set: MNIST dataset

Each Image is of  $28 \times 28 = 784$  pixels  $\therefore \# \text{ of features} = 784$   
 $m = \# \text{ of samples}$ .  $x_{\text{train}} = [m, \# \text{ of features}]$



$$\text{If } a^{[2]} = \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ \vdots \\ 9 \end{bmatrix} \quad \text{then there is 80% chance that the given sample is 1.}$$

2 hidden layers

$$\text{Layer 1: } (\# \text{ of neurons}, \# \text{ of inputs}) = (64, 784)$$

$$\text{Layer 2: } (\# \text{ of neurons}, \# \text{ of inputs}) = (10, 64)$$

Activation functions used:

$$\text{Layer 1: } a^{[1]} = \text{ReLU}(z^{[1]})$$

$$\text{Layer 2: } a^{[2]} = \text{softmax}(z^{[2]})$$

We need to run gradient descent in every epoch to improve the weights & bias of each neuron in every layer.

$$W^{[1]} = \begin{matrix} \text{feat}_1 & \text{feat}_2 & \dots & \text{feat}_{784} \\ \text{neu}_1 & w_{11} & w_{12} & \dots & w_{1,784} \\ \text{neu}_2 & w_{21} & w_{22} & \dots & w_{2,784} \\ \vdots & & & \ddots & \vdots \\ \text{neu}_{64} & w_{64,1} & w_{64,2} & \dots & w_{64,784} \end{matrix}$$

$$A^{[0]} = \begin{matrix} s_1 & s_2 & \dots & s_m \\ \text{fe}_1 & x_{1,1} & x_{2,1} & \dots & x_{m,1} \\ \text{fe}_2 & x_{1,2} & x_{2,2} & \dots & x_{m,2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{fe}_{784} & x_{1,784} & x_{2,784} & \dots & x_{m,784} \\ \downarrow & \downarrow & \downarrow & \ddots & \downarrow \\ \text{Samp}_1 & \text{Samp}_2 & \dots & \text{Samp}_m \end{matrix}$$

$$B^{[1]} = \begin{matrix} \text{neu}_1 \\ \text{neu}_2 \\ \vdots \\ \text{neu}_{64} \end{matrix}$$

$$A^{[0]} = x_{\text{train}}. T$$

$$y_{\text{train}} = [7 \ 8 \ 0 \ 1 \ \dots \ 9]_{1 \times m}$$

$\downarrow$        $\downarrow$        $\downarrow$        $\downarrow$

Olp samp<sub>1</sub>      Olp samp<sub>2</sub>      ...      Olp samp<sub>m</sub>

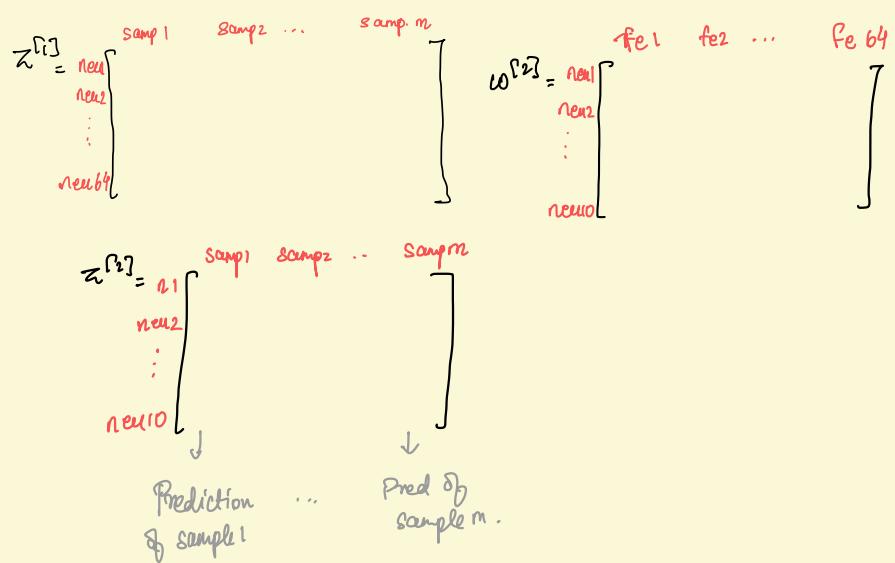
### ① Forward Propagation.

$$z^{[1]} = W^{[1]} A^{[0]} + B^{[1]}$$

$$A^{[1]} = \text{ReLU}(z^{[1]})$$

$$z^{[2]} = W^{[2]} A^{[1]} + B^{[2]}$$

$$A^{[2]} = \text{softmax}(z^{[2]})$$

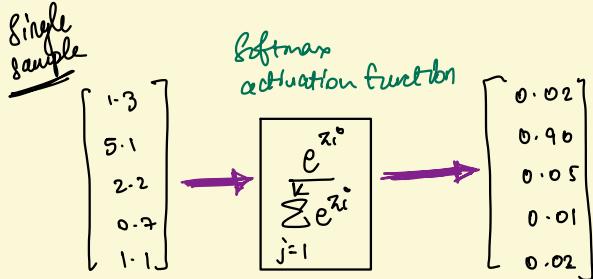


### ② Calculating the loss

$$\text{loss} = - \sum_{i=1}^n y_i \log a_i \quad \text{for one sample}$$

1 sample

$$\text{Cost} = - \sum_{i=1}^m \sum_{j=1}^n y_j^{(i)} \log(a_j^{(i)})$$



$$\text{If } y_1 = 5$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$a_1 = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.8 \\ 0.1 \\ \vdots \\ 0.1 \end{bmatrix}$$

$$L = \text{loss}(a_1, \dots, a_N, y) = \begin{cases} -\log a_1, & y = 1 \\ -\log a_2, & y = 2 \\ \vdots \\ -\log a_N, & y = N \end{cases}$$

### ③ Backpropagation.

#### a) Layer II

$$\text{Find } \frac{\partial C}{\partial z^{[2]}} \quad \text{Cost} = - \sum_{i=1}^m \sum_{j=1}^n y_j^{(i)} \log(a_j^{(i)})$$

$$\frac{\partial C}{\partial z^{[2]}} = \left[ \frac{\partial C}{\partial z_1^{[2]}}, \frac{\partial C}{\partial z_2^{[2]}}, \dots, \frac{\partial C}{\partial z_{10}^{[2]}} \right]$$

Expanding for all  $m$  samples,

Algorithm:

layer = outermost layer  $l$ ,  $\frac{\partial C}{\partial A^{[l]}}$  is given

repeat  $l \leftarrow l-1$

$$\frac{\partial C}{\partial z^{[l]}} = \sigma'(z^{[l]}) \frac{\partial C}{\partial A^{[l]}}$$

$$\frac{\partial C}{\partial w^{[l]}} = \frac{\partial C}{\partial z^{[l]}} (A^{[l-1]})^T \left[ \frac{1}{m} \right]$$

$$\frac{\partial C}{\partial b^{[l]}} = \frac{\partial C}{\partial z^{[l]}}$$

$$\begin{aligned} w^l &= w^l - \alpha \frac{\partial C}{\partial w^{[l]}} \\ b^l &= b^l - \alpha \frac{\partial C}{\partial b^{[l]}} \end{aligned} \quad \left. \begin{array}{l} \text{update step} \\ \text{up to } l=1 \end{array} \right\}$$

$$\frac{\partial C}{\partial A^{[l-1]}} = (W^{[l]})^T \frac{\partial C}{\partial z^{[l]}}$$

$$l = l-1 \quad \exists$$

$$\frac{\partial C}{\partial z^{[2]}} = \begin{bmatrix} \frac{\partial C}{\partial z_{11}^{[2]}} & \frac{\partial C}{\partial z_{12}^{[2]}} & \dots & \frac{\partial C}{\partial z_{101}^{[2]}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial C}{\partial z_{m1}^{[2]}} & \frac{\partial C}{\partial z_{m2}^{[2]}} & \dots & \frac{\partial C}{\partial z_{10m}^{[2]}} \end{bmatrix}^T$$

\* For a single sample,

$$A = \begin{bmatrix} a_1^{[2]} \\ a_2^{[2]} \\ a_3^{[2]} \\ \vdots \\ a_{10}^{[2]} \end{bmatrix} \quad a_i = \frac{e^{z_i^{[2]}}}{\sum_{j=1}^{10} e^{z_j^{[2]}}}$$

$$\bar{x}^{[2]} = \begin{bmatrix} z_1^{[2]} \\ z_2^{[2]} \\ \vdots \\ z_{10}^{[2]} \end{bmatrix}$$

$$\sigma'(\bar{x}^{[2]}) \quad \frac{\partial A^{[2]}}{\partial \bar{x}^{[2]}} = \begin{bmatrix} \frac{\partial a_1^{[2]}}{\partial z_1^{[2]}} & \frac{\partial a_1^{[2]}}{\partial z_2^{[2]}} & \dots & \frac{\partial a_1^{[2]}}{\partial z_{10}^{[2]}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial a_{10}^{[2]}}{\partial z_1^{[2]}} & \frac{\partial a_{10}^{[2]}}{\partial z_2^{[2]}} & \dots & \frac{\partial a_{10}^{[2]}}{\partial z_{10}^{[2]}} \end{bmatrix}$$

Recall  $a_i^o = \frac{e^{z_i^o}}{\sum_{j=1}^n e^{z_j^o}}$ , softmax.

Find  $\frac{\partial a_i^o}{\partial z_j}$ :

$$\text{Case I } i \neq j \quad \frac{\partial}{\partial z_j} \left[ \frac{e^{z_i}}{e^{z_i} + \sum_{k \neq i} e^{z_k}} \right] = \left( \frac{(e^{z_j} + \sum_{k \neq j} e^{z_k}) \frac{\partial}{\partial z_j} (e^{z_i}) - e^{z_i} e^{z_j}}{(\sum_{k \neq i} e^{z_k})^2} \right) = \frac{-\frac{e^{z_i}}{\sum_k e^{z_k}} \frac{e^{z_j}}{\sum_{k \neq i} e^{z_k}}}{(\sum_k e^{z_k}) (\sum_{k \neq i} e^{z_k})}$$

$$= -a_i^o a_j^o$$

$$\text{Case II } i = j \quad \frac{\partial a_i^o}{\partial z_i} = \frac{\partial}{\partial z_i} \left( \frac{e^{z_i}}{e^{z_i} + \sum_{k \neq i} e^{z_k}} \right) = \frac{(\sum_k e^{z_k}) e^{z_i} - (e^{z_i}) (e^{z_i})}{(\sum_k e^{z_k})^2}$$

$$= \frac{e^{z_i} (\sum_{k \neq i} e^{z_k})}{(\sum_k e^{z_k}) (\sum_{k \neq i} e^{z_k})} = a_i^o (1 - a_i^o)$$

$$\star \quad \frac{\partial L}{\partial z} = \left[ \frac{\partial L}{\partial z_1^{[2]}}, \frac{\partial L}{\partial z_2^{[2]}}, \dots, \frac{\partial L}{\partial z_{10}^{[2]}} \right]$$

Calculate  $\frac{\partial L}{\partial z_i^o}$  w.r.t each  $z_i^o$

$$\frac{\partial L}{\partial z_i} = \frac{\partial}{\partial z_i} \left[ -y_i \log(a_i) - \sum_{k \neq i} y_k \log(a_k) \right] = -\frac{y_i}{a_i} \left( \frac{\partial a_i}{\partial z_i} \right) - \sum_{k \neq i} \frac{y_k}{a_k} \frac{\partial a_k}{\partial z_i}$$

i ranging from  
1 to n.

$$= -\frac{y_i}{a_i} a_i (1-a_i) + \sum_{k \neq i} \frac{y_k}{a_k} a_k \cdot a_i = -y_i + y_i a_i + a_i \sum_{k \neq i} y_k$$

$$= -y_i + a_i (y_i + \sum_{k \neq i} y_k) = -y_i + a_i (\sum_{k=1}^n y_k) = a_i - y_i.$$

$$\therefore \frac{\partial L}{\partial z_i} = a_i - y_i \quad \forall i=1, 2, \dots, n.$$

$$\frac{\partial C}{\partial z^{[2]}} = \begin{bmatrix} \frac{\partial C}{\partial z_{11}^{[2]}} & \frac{\partial C}{\partial z_{12}^{[2]}} & \dots & \frac{\partial C}{\partial z_{10}^{[2]}} \\ \frac{\partial C}{\partial z_{21}^{[2]}} & \frac{\partial C}{\partial z_{22}^{[2]}} & \dots & \frac{\partial C}{\partial z_{20}^{[2]}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial C}{\partial z_{m1}^{[2]}} & \frac{\partial C}{\partial z_{m2}^{[2]}} & \dots & \frac{\partial C}{\partial z_{m0}^{[2]}} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial z_1} & \frac{\partial L}{\partial z_2} & \dots & \frac{\partial L}{\partial z_{10}} \end{bmatrix}$$

$$A^{[2]} = \begin{bmatrix} \text{feat 1} & \text{feat 2} & \dots & \text{feat m} \\ \text{feat 1} & \alpha_{11}^{[2]} & \alpha_{12}^{[2]} & \dots & \alpha_{10}^{[2]} \\ \text{feat 2} & \alpha_{21}^{[2]} & \alpha_{22}^{[2]} & \dots & \alpha_{20}^{[2]} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{feat m} & \alpha_{m1}^{[2]} & \alpha_{m2}^{[2]} & \dots & \alpha_{m0}^{[2]} \end{bmatrix}$$

From the above, we get,

$$\frac{\partial L_m}{\partial z} = \begin{bmatrix} \frac{\partial L_m}{\partial z_1} & \frac{\partial L_m}{\partial z_2} & \dots & \frac{\partial L_m}{\partial z_{10}} \end{bmatrix}$$

$$\frac{\partial L_m}{\partial z} = \begin{bmatrix} a_{m1} - y_{m1} & a_{m2} - y_{m2} & \dots & a_{m10} - y_{m10} \end{bmatrix}$$

$$\frac{\partial C}{\partial z^{[2]}} = \begin{bmatrix} \frac{\partial L_1}{\partial z_1} & \frac{\partial L_1}{\partial z_2} & \dots & \frac{\partial L_1}{\partial z_{10}} \\ \frac{\partial L_2}{\partial z_1} & \frac{\partial L_2}{\partial z_2} & \dots & \frac{\partial L_2}{\partial z_{10}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial L_m}{\partial z_1} & \frac{\partial L_m}{\partial z_2} & \dots & \frac{\partial L_m}{\partial z_{10}} \end{bmatrix} = \begin{bmatrix} a_{11} - y_{11} & a_{12} - y_{12} & \dots & a_{110} - y_{110} \\ a_{21} - y_{21} & a_{22} - y_{22} & \dots & a_{210} - y_{210} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} - y_{m1} & a_{m2} - y_{m2} & \dots & a_{m10} - y_{m10} \end{bmatrix}^T$$

$$= A^{[2]} - y$$

size of

$$\frac{\partial C}{\partial z^{[2]}} = (\# \text{ neurons in layer 2}, \# \text{ samples})$$

Hence we have, for layer II

$$\frac{\partial C}{\partial z^{[2]}} = A^{[2]} - y$$

(10, m)

$$\frac{\partial C}{\partial w^{[2]}} = \left(\frac{1}{m}\right) \frac{\partial C}{\partial z^{[2]}} \cdot (A^{[1]})^T$$

(10, m) (m, 64)

$$\frac{\partial C}{\partial b^{[2]}} = \left(\frac{1}{m}\right) \frac{\partial C}{\partial z^{[2]}}$$

(10, 1)

$$w^{[2]} = w^{[2]} - \alpha \frac{\partial C}{\partial w^{[2]}}$$

of size (10, 64)

$$b^{[2]} = b^{[2]} - \alpha \frac{\partial C}{\partial b^{[2]}}$$

of size (10, 1)

$$\frac{\partial C}{\partial A^{[1]}} = (w^{[2]})^T \frac{\partial C}{\partial z^{[2]}}$$

(64, m) (64, 10) (10, m)

$$\frac{\partial C}{\partial z^{[1]}} = \sigma'(z^{[1]}) \frac{\partial C}{\partial A^{[1]}}$$

$$\frac{\partial C}{\partial w^{[1]}} = \frac{\partial C}{\partial z^{[1]}} (A^{[1]})^T \left(\frac{1}{m}\right)$$

$$\frac{\partial C}{\partial b^{[1]}} = \frac{\partial C}{\partial z^{[1]}} \left(\frac{1}{m}\right) \quad \text{vector got by sum along columns}$$

$$\begin{aligned} w^l &= w^l - \alpha \frac{\partial C}{\partial w^l} \\ b^l &= b^l - \alpha \frac{\partial C}{\partial b^l} \end{aligned} \quad \text{updation step}$$

$$\frac{\partial C}{\partial A^{[2]}} = (w^{[2]})^T \frac{\partial C}{\partial z^{[2]}}$$

$$l = l-1 \quad \exists$$

} Updating weights of the second layer.

We have successfully updated the layer 2 weights.

## B) Back propagation for layer I

Given:  $\frac{\partial C}{\partial A^{[2]}}$  from the previous layer.

To calc:  $\frac{\partial C}{\partial z^{[1]}}$

$$\frac{\partial C}{\partial z^{[1]}} = \sigma'(z^{[1]}) \frac{\partial C}{\partial A^{[1]}}$$

Recall  $\sigma = \text{ReLU}$  for 1st layer.

$$\text{ReLU}(z) = \begin{cases} z, & z \geq 0 \\ 0, & \text{else.} \end{cases}$$

$$\therefore \sigma'(z) = \begin{cases} 1, & z \geq 0 \\ 0, & \text{else.} \end{cases}$$

$$\sigma'(z^{[1]}) = \text{size } (64, m)$$

Caution: This is an element wise multiplication.

$$\frac{\partial C}{\partial z^{[1]}} = \sigma'(z^{[1]}) \cdot \frac{\partial C}{\partial A^{[1]}}$$

(64, m) (64, m) (64, m)

$$\therefore \frac{\partial C}{\partial z^{[1]}} = \sigma'(z^{[1]}). \frac{\partial C}{\partial A^{[0]}}$$

$\downarrow$   
 $(64, m)$

$\downarrow$   
 $(64, m)$

$$\frac{\partial C}{\partial w^{[1]}} = \left(\frac{1}{m}\right) \frac{\partial C}{\partial z^{[1]}} (A^{[0]})^T$$

$(64, m)$        $(m, 784)$   
 $(64, 784)$

$$\frac{\partial C}{\partial b^{[1]}} = \left(\frac{1}{m}\right) \frac{\partial C}{\partial z^{[1]}} \rightarrow \text{Vector summed along the columns.}$$

$[64, 1]$

$$W^{[1]} = W^{[1]} - \alpha \frac{\partial C}{\partial w^{[1]}}$$

$\left. \begin{array}{l} \\ \end{array} \right\} \text{Updating the weights of the layer 1}$

$$b^{[1]} = b^{[1]} - \alpha \frac{\partial C}{\partial b^{[1]}}$$

End of Back propagation.

In Summary, we have,

Our NN will have a simple two-layer architecture. Input layer  $a^{[0]}$  will have 784 units corresponding to the 784 pixels in each  $28 \times 28$  input image. A hidden layer  $a^{[1]}$  will have 10 units with ReLU activation, and finally our output layer  $a^{[2]}$  will have 10 units corresponding to the ten digit classes with softmax activation.

#### Forward propagation

$$\begin{aligned} Z^{[1]} &= W^{[1]} X + b^{[1]} \\ A^{[1]} &= g_{\text{ReLU}}(Z^{[1]}) \\ Z^{[2]} &= W^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} &= g_{\text{softmax}}(Z^{[2]}) \end{aligned}$$

#### Backward propagation

$$\begin{aligned} dZ^{[2]} &= A^{[2]} - Y \\ dW^{[2]} &= \frac{1}{m} dZ^{[2]} A^{[1]T} \\ dB^{[2]} &= \frac{1}{m} \sum dZ^{[2]} \\ dZ^{[1]} &= W^{[2]T} dZ^{[2]} \cdot g^{[1]'}(z^{[1]}) \\ dW^{[1]} &= \frac{1}{m} dZ^{[1]} A^{[0]T} \\ dB^{[1]} &= \frac{1}{m} \sum dZ^{[1]} \end{aligned}$$

#### Parameter updates

$$\begin{aligned} W^{[2]} &:= W^{[2]} - \alpha dW^{[2]} \\ b^{[2]} &:= b^{[2]} - \alpha db^{[2]} \\ W^{[1]} &:= W^{[1]} - \alpha dW^{[1]} \\ b^{[1]} &:= b^{[1]} - \alpha db^{[1]} \end{aligned}$$

Forward prop

- $A^{[0]} = X: 784 \times m$
- $Z^{[1]} \sim A^{[1]}: 10 \times m$
- $W^{[1]}: 10 \times 784$  (as  $W^{[1]}A^{[0]} \sim Z^{[1]}$ )
- $B^{[1]}: 10 \times 1$
- $Z^{[2]} \sim A^{[2]}: 10 \times m$
- $W^{[2]}: 10 \times 10$  (as  $W^{[2]}A^{[1]} \sim Z^{[2]}$ )
- $B^{[2]}: 10 \times 1$

Backprop

- $dZ^{[2]}: 10 \times m$  ( $A^{[2]}$ )
- $dW^{[2]}: 10 \times 10$
- $dB^{[2]}: 10 \times 1$
- $dZ^{[1]}: 10 \times m$  ( $A^{[1]}$ )
- $dW^{[1]}: 10 \times 10$
- $dB^{[1]}: 10 \times 1$

In our case,  
the first layer has  
64 neurons so  
modify  
accordingly

Source: Many many youtube videos for derivation. (Coding have mainly)  
for the code, primarily Samson Zhang

Implementation from scratch.

Final Note : If  $y_{\text{train}} = [4 \ 2 \ 5 \ \dots \ 0 \ 6]_{1 \times m}$

then Y matrix will look like:

$$Y = \begin{matrix} & 1 & 2 & 3 & \dots & m+m \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \end{matrix} & \left[ \begin{matrix} 0 & 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \\ 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & \dots & \vdots & 1 \\ 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \end{matrix} \right] \end{matrix}$$

How can we make our python code run faster?

Vectorizing  
Broadcasting

## Broadcasting

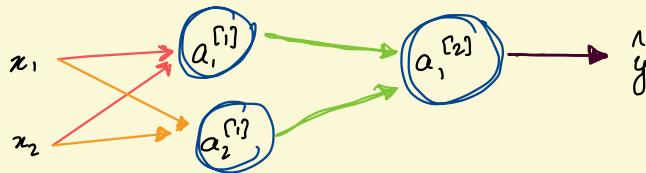
	Apples	Beef	Eggs	Potatoes
Carb	56	0	4.4	68
Protein	1.2	104	52	8
Fat	1.8	185	99	0.9

Calculate %age of calories from Carb, protein and Fat. Can you do this without the explicit for loop?

$$\text{Cal} = \text{np.sum}(A, \text{axis}=0)$$

$$\text{Percentage} = 100 * A / \text{cal} \cdot \text{reshape}(1, 4)$$

## \* Random Initialization



Initializing b to zero is okay  
but initializing w to zero is not okay.

$$\text{If we initialise } W^{[1]} = [0 \ 0], \quad W^{[2]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]}$$

while Back propagation,  $dz^{[1]} = dz^{[2]}$ . When we initialize to zero, the hidden units are completely identical. In other words, completely symmetric so even after n iterations, all the neurons compute the same function as in,

$$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

$$W^{[1]} = \text{np.random.rand}(2, 2) * 0.01 \quad \text{Small value}$$

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

If w is large z is big so it may cause  $\sigma(z) \approx 1$  or to be saturated slowing down the learning rate hence gradient descent

# Logistic regression derivatives

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

→ for  $i = 1$  to  $m$ :

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

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

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

$$\frac{\partial z^{(i)}}{\partial w} = a^{(i)} - y^{(i)}$$

$$\left. \begin{array}{l} \frac{\partial J}{\partial w_1} += x_1^{(i)} dz^{(i)} \\ \frac{\partial J}{\partial w_2} += x_2^{(i)} dz^{(i)} \\ \frac{\partial J}{\partial b} += dz^{(i)} \end{array} \right| n_x=2$$

$$J = J/m, dw_1 = \frac{\partial J}{\partial w_1}/m, dw_2 = \frac{\partial J}{\partial w_2}/m, db = \frac{\partial J}{\partial b}/m$$

# Logistic regression derivatives

$$J = 0, \boxed{dw_1 = 0, dw_2 = 0}, db = 0$$

$$dw = np.zeros((n_x, 1))$$

→ for  $i = 1$  to  $m$ :

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

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

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

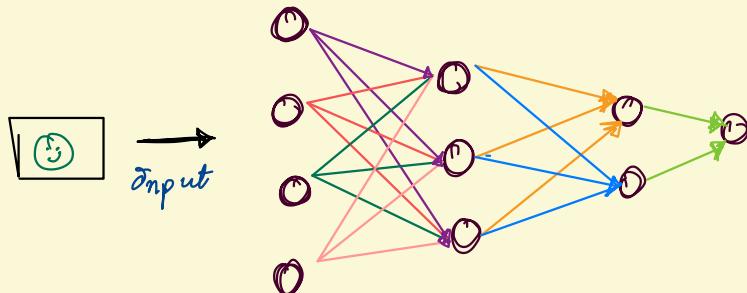
$$\frac{\partial z^{(i)}}{\partial w} = a^{(i)} - y^{(i)}$$

$$\left. \begin{array}{l} \frac{\partial J}{\partial w_1} += x_1^{(i)} dz^{(i)} \\ \frac{\partial J}{\partial w_2} += x_2^{(i)} dz^{(i)} \\ \frac{\partial J}{\partial b} += dz^{(i)} \end{array} \right| n_x=2$$

$$J = J/m, \boxed{dw_1 = \frac{\partial J}{\partial w_1}/m, dw_2 = \frac{\partial J}{\partial w_2}/m}, db = \frac{\partial J}{\partial b}/m$$

$$dw /= m$$

↳ Intuition about deep representations:



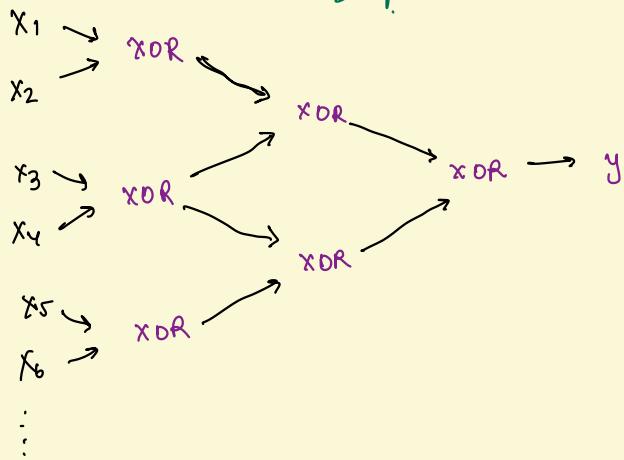
First layer starts off detecting simple things like edges, as the layers increase, the deep layers detect more complex things like eyes, faces etc.

## ★ Circuit theory and deep learning

These are functions you can compute with a small  $L$ -layer deep neural network than shallower networks that require exponentially more hidden units to compute

$$y = (x_1 \text{ XOR } x_2) \text{ XOR } (x_3 \text{ XOR } x_4) \dots \text{ XOR } x_N$$

## Deep Network



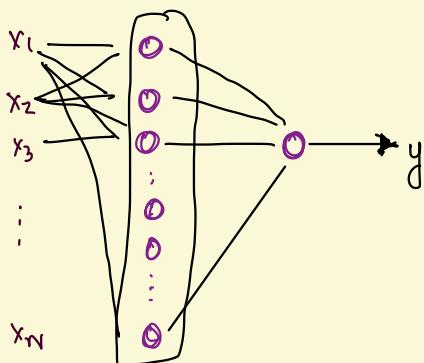
A binary tree with  $n$  leaves has a height of  $a$

$$\text{With 4 leaves, } 4 = 2^{\lceil \frac{n}{2} \rceil} - 2+1$$

$$\text{with 6, } 6 = 2^{\lceil \frac{n}{2} \rceil} - 3+2+1 = 6$$

with 8,

## Shallow Network



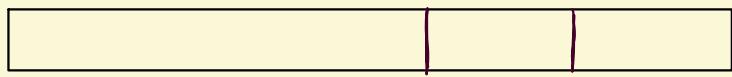
\* Parameters:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]}, \dots$

Hyperparameters: learning rate  $\alpha$   
 $\downarrow$   
# of iterations  $n$

Parameters that  
ultimately control  
values of  $W$  &  $B$   
# of hidden layers  $L$   
# of hidden units  $n^{[1]}, n^{[2]}, \dots$

Choice of activation functions

\* Training set, Dev Set, Test Set



Training Set

Dev Set

test Set

to see which of  
many different models  
perform best on dev set

role of it is we test different  
algorithms on it & see which works better

60% train 20% Dev 20% test → Okay for a small dataset

98% train 1% Dev 1% test → Okay for large datasets like 1 million

### \* Mismatched train/test distribution:

You are building an app that lets users upload a lot of pictures & the goal is to find pictures of cats in order to show the user.

Training set

Cat pics from webpages

high resolution, clear

Dev/test set

Cat pics from users your app.

low resolution, blurry

Both the sets come from different distributions.

Make sure dev/test set come from same distribution

Train set error

1%

15%

Dev set error

11%

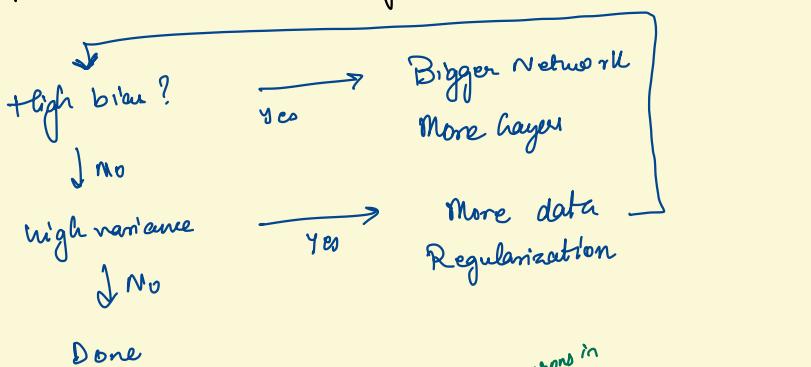
16%

high variance

high bias

Bayes error = 0% (true error)

### \* Basic recipe for machine learning:



$$\|w^{[l]} \|^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2$$

# of neurons in current layer      # of neurons in previous layer

### \* Frobenius Norm :

$$\|w^{[l]} \|^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2$$

## \* Logistic regression

- $L^2$  regularization.

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

- $L^1$  regularization

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m h(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_1$$

## \* Neural Network

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

$\downarrow$   
froeb norm

$\# \text{of neurons in}$   
 $\text{current layer}$   
 $\rightarrow$   
 $n^{(l)}$   
 $\# \text{of neurons in}$   
 $\text{previous layer}$   
 $\rightarrow$   
 $n^{(l-1)}$

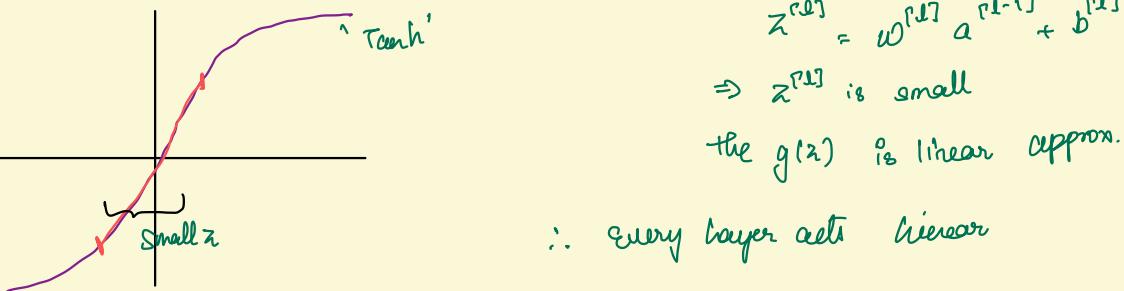
$$\|w^{(l)}\|^2 = \sum_{i=1}^{n^{(l)}} \sum_{j=1}^{n^{(l-1)}} (w_{ij}^{(l)})^2$$

Now,  $d w^{(l)}$  becomes,

$$d w^{(l)} = (\text{from backprop}) + \frac{\lambda}{m} w^{(l)}$$

$$w^{(l)} = w^{(l)} - \alpha d w^{(l)}$$

## \* Why Regularization reduces overfitting?



## \* DROPOUT REGULARIZATION:

Go through each layer go through each node & toss a coin and have 0.5 chance of keeping each node & 0.5 chance of removing each node.

At each iteration remove different set of nodes and observe the performance of the remaining neural network. On each iteration, working with a smaller neural network

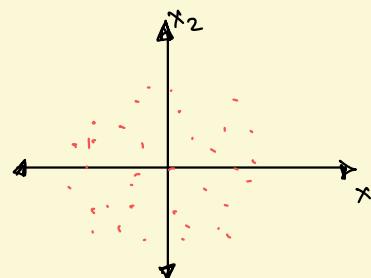
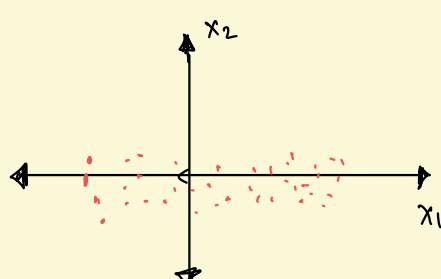
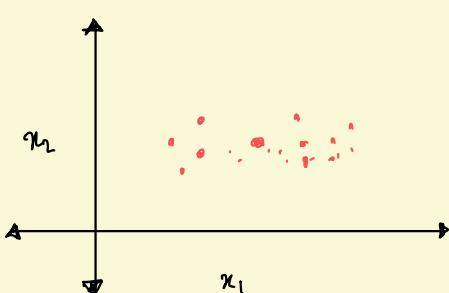
\* Normalizing training sets:

Subtract mean:  $\bar{x} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$

$$x^{(i)}_j = x^{(i), j} - \bar{x}_j$$

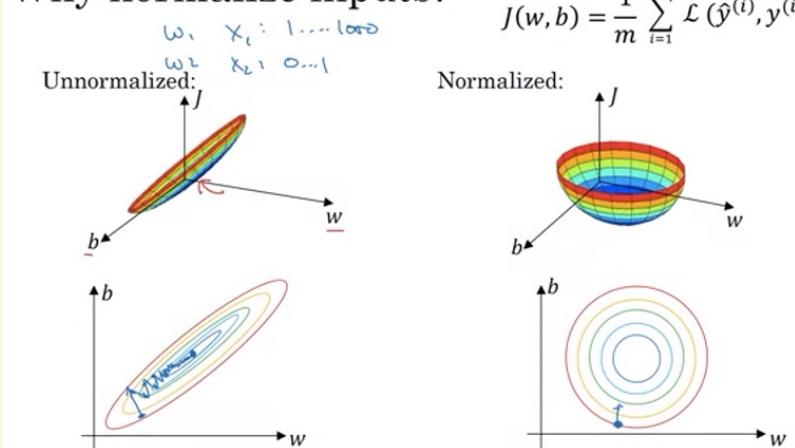
Variance:  $\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i), j})^2$

$$x^{(i), j} = x^{(i), j} / \sigma$$

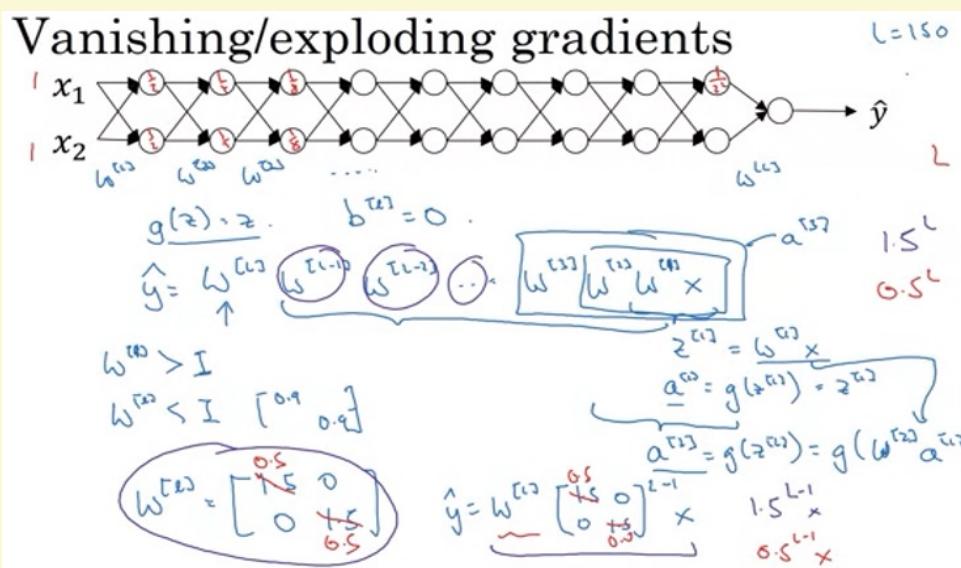


Use same  $\bar{x}$  &  $\sigma$  to normalize the test & train set.

Why normalize inputs?



\* Vanishing / Exploding Gradients: when training neural network, when it's very deep slopes can sometimes get either very big or very very small

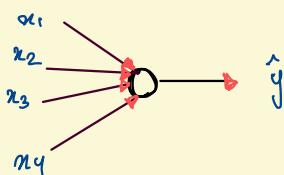


If the  $w^{(l)} \geq 1$  then the activation o/p hence the gradient can explode

& if the  $w^{(l)} < 1$  then activation o/p & hence grad. can vanish.

Careful choice of initializing weights might help  
 ) solution

## \* Choice of Weight Initialization for the Neural Network.



$$z = w_1x_1 + w_2x_2 + \dots + w_n x_n + b$$

larger  $n$ , smaller  $w_i$

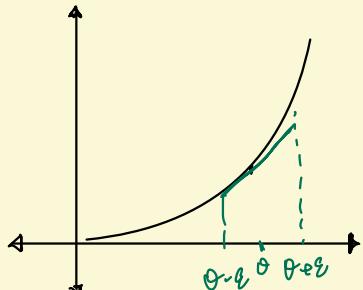
$$\text{Var}(w) := \frac{1}{n} \quad n = \# \text{ of inputs}$$

$$\text{Var}(w) = \frac{2}{n}$$

$$W^{[l]} = np.random.randn(\text{shape}) * np.sqrt(\frac{2}{n^{l-1}})$$

In case of tanh, one can use  $\text{Var}(w) = \frac{1}{n^{l-1}}$   $\Leftarrow$  Xavier initialization.

## \* Numerical Approximation of Gradients.



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

$$\frac{(1.01)^2 - (0.99)^2}{2(0.01)} \approx 8.0001 \approx 3$$

$$g(\theta) = 3\theta^2 = 3$$

$$\text{Approx. error: } 0.0001$$

Gradient Checker      For each  $i$ ,  $d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$

$$\approx d\theta[i] = \frac{\partial J}{\partial \theta^i}$$

Check  $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$   $\epsilon = 10^{-7}$

## \* Optimization Algorithms:

### Mini Batch Gradient Descent

Vectorization allows you to efficiently compute on  $m$  examples.

What if  $m$  was large?  $m = 5 \text{ million} = 5,000,000$

Split the training data into mini batches into 1000 egs each.

$$X = [x^{(1)} \ x^{(2)} \ \dots \ x^{(1000)} \mid x^{(1001)} \ \dots \ x^{(2000)} \mid \dots \mid x^{(m)}]$$

$x^{(1)}$                      $x^{(2)}$                      $\dots$                      $x^{(1000)}$   
 $x^{(1001)}$              $x^{(1002)}$              $\dots$              $x^{(2000)}$   
 $y^{(1)}$                      $y^{(2)}$                      $\dots$                      $y^{(1000)}$

Batch: process entire training set @ the same time.

MiniBatch: process a mini batch @ one time.

### Mini Batch

for  $t=1, \dots, 5000$

Forward Prop on  $X^{(t)}$

$$\tilde{Z}^{(t)} = W^{(t)} X^{(t)} + b^{(t)}$$

$$A^{(t)} = g^{(t)}(\tilde{Z}^{(t)})$$

:

$$A^{(t)} = g^{(t)}(\tilde{Z}^{(t)})$$

$$\text{Compute Cost } J^{(t)} = \frac{1}{1000} \sum_{i=1}^L h(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|w^{(l)}\|_F^2$$

Back prop to wrt  $J^{(t)}$  (using  $(X^{(t)}, y^{(t)})$ )

$$W^{(t)} := W^{(t)} - dW^{(t)}, \quad b^{(t)} := b^{(t)} - db^{(t)}$$

3 Mini Batch gradient Descent runs much faster than the Batch gradient Descent when size of training set is large.

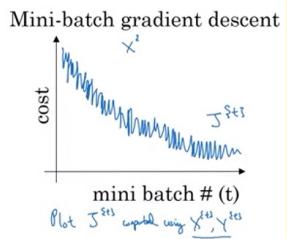
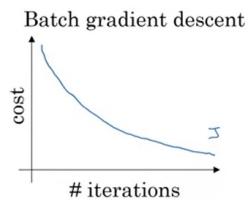
### Choosing Batches:

If small train set use Batch gradient descent ( $m \leq 2000$ )

Typical mini Batch sizes:      64      128      256      512  
 $2^6$        $2^7$        $2^8$        $2^9$

## Choosing your mini-batch size

Training with mini batch gradient descent



→ If mini-batch size =  $m$  : Batch gradient descent.  $(X^{(1)}, Y^{(1)}) = (X, Y)$

→ If mini-batch size = 1 : Stochastic gradient descent. Every example is its own mini-batch.  $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$

In practice: Somewhere in-between 1 and  $m$

Stochastic  
gradient  
descent

Use smaller  
step size

In-between  
(mini-batch size  
not too big/small)

Faster learning:  
• Vectorization.  
( $\sim 1000$ )  
• Make passes without  
finishing entire thing yet.

Batch  
gradient descent  
(mini-batch size =  $m$ )

Too long  
per iteration

Andrew Ng

\* Exponentially weighted averages.