

$c1 = 0.1$   
 $c2 = 0.4$   
 $c3 = 0.6$

use them as features and input them to another classifier

- nn = learn non-linear boundary

- we use  $\sigma$  in training

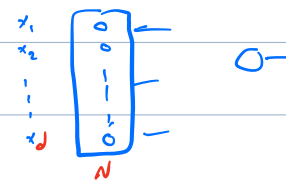
- if you don't put non-linearity in hidden layers = you will have another Logistic Regression

Standard nn (FFNN, FCNN):-  $\# \text{ of neurons} = JN + W$

-  $\# \text{ of features} = \# \text{ of neurons}$

- many features = less  $\# \text{ of neurons}$

- don't make any hidden layer  $<$  output layer



$\#$  How to make prediction?

first neuron  $\left\{ \begin{aligned} z_1^{(1)} &= w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + \dots + b_1^{(1)} \\ a_1^{(1)} &= \sigma(z_1^{(1)}) \end{aligned} \right.$

first hidden layer  $\Rightarrow \vec{a}^{(2)} = \sigma(WX + \vec{b})$

$$z^{(2)} = w_1^{(2)} \cdot a_1^{(1)} + w_2^{(2)} \cdot a_2^{(1)} + \dots + b^{(2)}$$

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

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

✱ How to training ?

to update the parameters

backpropagation

for any nn you will not get same accuracy every time

$$\frac{dJ}{dw_j^{(1)}} = \frac{1}{n} \sum (a_j^{(1)} - y) \cdot a_j^{(1)}$$

$$\frac{dJ}{db^{(1)}} = \frac{1}{n} \sum (a_j^{(1)} - y)$$

to get different gradients

initialize w and b with Random [not zeros]

$$\frac{dJ}{dw_j^{(1)}} = \frac{dJ}{dz^{(2)}} \times \frac{dz^{(2)}}{da^{(1)}} \times \frac{da^{(1)}}{dz^{(1)}} \times \frac{dz^{(1)}}{dw_j^{(1)}} \\ (a_j^{(1)} - y) \times 1 \times \sigma'(z_j^{(1)}) \times (1 - \sigma(z_j^{(1)})) \times$$

Loss function is not convex!

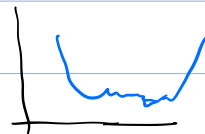
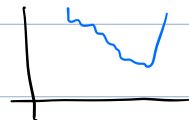
but

will not happened

in highly dimension feature space :-

→ many local minimums

→ local minimums are close to global one



what happened

GD for NN :-

1 → initial w, b with random

2 → forward pass

$z^1, a^1, z^2, a^2$

3 → compute loss, J

4 → Backpropagation  $dw^1, db^1, dw^2, db^2$

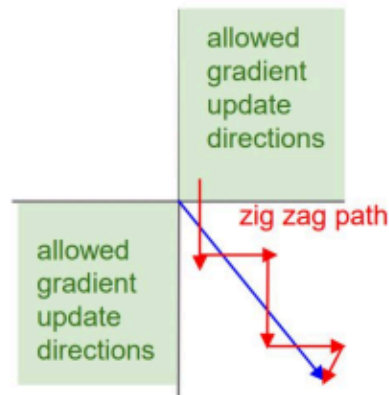
5 → update  $w^1, w^2, b^1, b^2$

repeat 2 → 5 until converge  $|J^k - J^{k-1}| < \epsilon$

## Three Problems with activation functions:-

- **Non zero-centered:** (Slower in converge)

the problem of not having zero-centered activation function is slow convergence, the gradient will move in a zig-zag path. To illustrate the problem more let us see this picture:



The blue line is the optimal path, but the non zero-center will move like the red path ( zig-zag, always move vertical or horizontal ) which lead to slow convergence.

- **Saturation:** (no learning)

the problem with saturation is that the gradient will be zero and this kill the gradient since if the gradient is zero ( $dW = 0$ ) then the update will be:

$$W = W - \alpha * dw = W - \alpha * 0 = W$$

We can see that the value of W does not change, this lead to make the network stop learning.

- **Computational Cost:**

the computational time for some activation function is more than other. For example, calculating the value of ReLU function is faster than sigmoid and tanh since both of them contain exponential calculation ( $e^x$ ).

problem with sigmoid:-

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

→ not zero centered values (0, 1) all positive (0.5 centered)

→ gradient saturate easily

→ exponential computation

$$\text{derivative} = \sigma(x)(1 - \sigma(x))$$

ReLU:-  $R(x) = \max(0, x)$

→ less computation

→ not zero centered

→ gradient not all saturate just negatives →

derivative (gradient)

$$\begin{cases} 0 & \text{if } z \leq 0 \\ 1 & \text{other} \end{cases}$$

Tanh:-

→ zero centered

→ gradient saturated easily

gradient  
→  $1 - \tanh^2$

	ReLU	Tanh	Sigmoid
zero-centered	Not zero-centered	Zero-centered	Not zero-centered
Saturation	Dose not saturate	saturated	saturated
Computational Cost	efficient	slow	slow

Why we go deep

to avoid overfitting

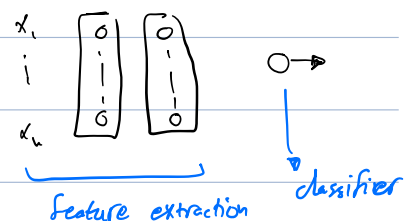
being more efficient (less # of neurons in each layer)

DL vs ML:-

having more data  $\Rightarrow$  DL

having more feature  $\Rightarrow$  DL

there is no specific features (images..)  $\Rightarrow$  DL



Structured Data  $\Rightarrow$  ML

Prosperity	Machine Learning	Deep learning
Data structure	Structure data	Unstructured data
Size of the data	Could works on small or medium data	Required huge amount of data
Training time	Less training time	Large training time
Feature Extraction	Need to understand the features	Don't have to understand the feature, it could learn the features by itself
Accuracy	Less accuracy than deep learning	Higher accuracy than machine learning

\* binary loss cross-entropy for binary (Logistics)

\* Multi class

Data:

convert  $y \Rightarrow$  one-hot-encoding

Architectures-

wider layers = overfitting

Deeper layers less overfitting

$\rightarrow$  \* neurons in output layer = \* of classes

$\rightarrow$  apply Softmax

$\rightarrow$  loss  $\Rightarrow$  categorical cross-entropy  $(\sum_{i=1}^K -y \log(\hat{y}))$

TARGET

PREDICTION

for  $i=1$

1	0.5
0	0.3
0	0.2

for multi-label classification:

outputs \* of classes

sigmoid

What we change

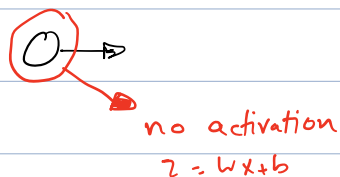
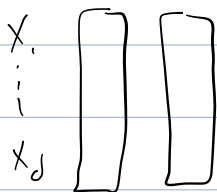
{ \* of output layer  
activation of output layer  
loss

$$\text{Softmax}(X) = \frac{e^x}{\sum e^i}$$

all output neurons

# ~~XX~~ Regression by nn

Predict the age of image



$$\text{loss} = (\hat{y} - y)^2$$

$$\text{Cost} = \frac{1}{2m} \sum (\hat{y} - y)^2$$

*MSE*

What we change

{ activation in the output layer  
loss

# ~~XX~~ big data $\Rightarrow$ 1 step is costly and may out of memory

Solutions

Gradient Descent:-

you have to do

FP  
bp  
update

for all samples

to do just one step

Costly +  
out of memory

Stochastic GD:-

you have to do

FP  
bp  
update

for just one sample

to do just one step

noisy

mini-batch GD:-

you have to do

FP  
bp  
update

for just one batch

to do just one decent

the best

batch size : 16, 32, 64

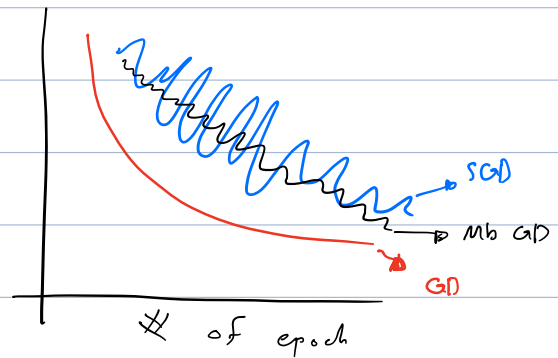
# of samples

Example:-  $m = 10000$  and we loop on  $m$

do GD  $\Rightarrow$  1 step

do SGD  $\Rightarrow$  10000 steps

do Mb GD  $\Rightarrow$  100 steps  
batch size = 100



epoch  $\Rightarrow$  cover all samples

~~Optimizers~~ to optimize the update step  
instead of  $(w_i = w_i - \alpha dw_i)$

GD with momentum :- to have a capability to escape local minimal (Ball)

GD with RMSprop :- to have faster GD

Adam combine GDM and RMSprop

