

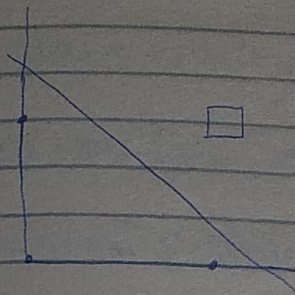
18-210-2021

L

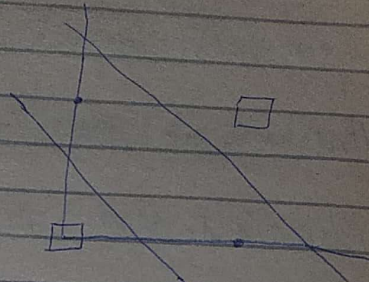
Monday

week #8 Lecture #11

Neural Networks



AND



XNOR

Each logistic unit = Neuron

— ○

— ○

⋮

— ○

○

○

○

○

○ —

○ —

○ —

Input
layer

Hidden layer

Output
layer

2 layered Neural Network

we don't count Input layer in no.'s.

Deep Learning

If more than 1 hidden layer, we say it is deep learning.

2^{four} = Neurons in each layer usually.

Question:-

If we increase No. of Neurons in one hidden layer, what does it change?

Ans: we can now learn more complicated ~~gene~~ Polynomials. More accuracy. Capacity is Increased.

Question:-

If we increase No. of layers in ~~NN~~ NN, what does it change?

Ans) we can now learn more complicated Polynomials. More accuracy, Capacity is Increased.

Problem:-

If we have only 2 classes & many layers & Neurons, then its problem is that it overfits & memorizes the data, tries to fit every example Individually, also covers noise in the data.

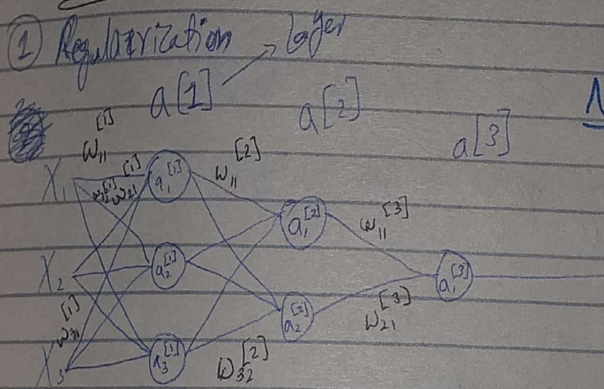
Solution:-

W

If no. of classes, is small, less, then don't make deep learning NN or make more Neurons & layers.
 or if Dataset is small, then also don't make deep learning NN or make more Neurons & layers.

Techniques to Solve this:

① Regularization



Naming Convention:-

Fully
connected
NN

Where all inputs are connected to input layer Neurons

$$Z_1^{[1]} = w_{11}^{[1]} X_1 + w_{21}^{[1]} X_2 + w_{31}^{[1]} X_3$$

$$Z_2^{[1]} = w_{12}^{[1]} X_1 + w_{22}^{[1]} X_2 + w_{32}^{[1]} X_3$$

$$Z_3^{[1]} = w_{13}^{[1]} X_1 + w_{23}^{[1]} X_2 + w_{33}^{[1]} X_3$$

Every Neuron calculates weighted sums, then it applies activation function.

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$a_2^{[1]} = \sigma(z_2^{[1]})$$

$$a_3^{[1]} = \sigma(z_3^{[1]})$$

Question:

If we don't apply this activation function in hidden layers & only apply it on output layer, what does it imply,

Solution: our learning slows down.

$$f = w_3(w_2(w_1(x)))$$

$$w = w_1 * w_2 * w_3$$

$$f = wx \quad // \text{it becomes linear classifier}$$

~~Activation function~~ Activation function acts as a linear classifier. It adds non-linearity.

Question:

If we use Sigmoid as an Activation function,
what does it imply?

Solution:

Local Gradient

$$\frac{dL}{dw_{11}^{[1]}} = \sigma(z_1^{[1]}) (1 - \sigma(z_1^{[1]})) * X_1 * \text{upstream gradients}$$

if we have $z_1^{[1]} = +10$, then

$$\sigma(z_1^{[1]}) = 1$$

$$\frac{dL}{dw_{11}^{[1]}} = 1(1-1) = 0$$

So then our weights won't be updated ever.

$$w_{11}^{[1]} = w_{11}^{[1]} - \underbrace{\alpha \frac{dL}{dw_{11}^{[1]}}}_0$$

So weights won't be update

Sigmoid Problem -

It kills gradients of weights.

If weighted sum value is -4 or $+4$
Grad > 0 or < 0

Vanishing Gradient Problem:

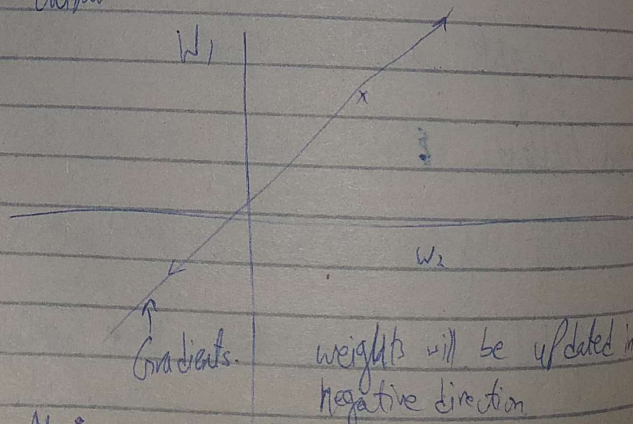
If we have more layers meaning more & bigger NN , then weights multiplied together becomes really small.

② Problem of sigmoid: ~~Vanishing Gradient Problem~~

Weights are positive, & gradients are positive, then Sigmoid output is also positive so ~~the~~ our weights are updated in either positive or negative so our learning is ^{not} improved & gets affected.

It can get negative bcz we back propagate from upstream gradients.

It saturates the output.



we only explore weights in one direction. $w_1 = w_1 - \alpha \frac{\partial L}{\partial w_1}$
 w_1 decreases. $+$