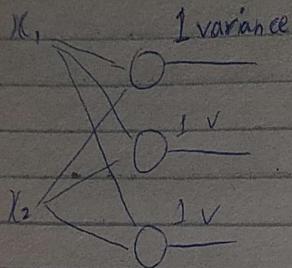


15-11-2021

Week #12 Lecture #15

Monday

Batch Normalization:-



for 1-dense (inp-shape, {3, "relu"})
fully connected layer.

no. inputs no. of outputs

↓ ↓
3 neurons

Activation function

Above is in Xavier Initialization.

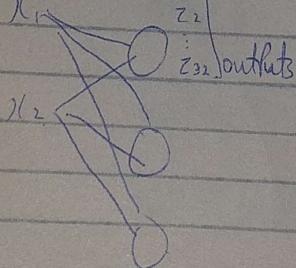
Batch Normalization:-

called Batch Norm.

We add a layer after each layer in our NN.

So we send
 \downarrow batches

\bar{x}_1 \bar{x}_2 \bar{z}_1 \bar{z}_2
 $\bar{x}_1 = \text{mean}$
 $\bar{z}_1 = \text{standardized output}$



apply standardization
on each neuron

& the input on each neuron
is 32 images

their mean is zero
& variance is 1.

32-images = batch

If we use this, then our weights initialization does affect our output that much.

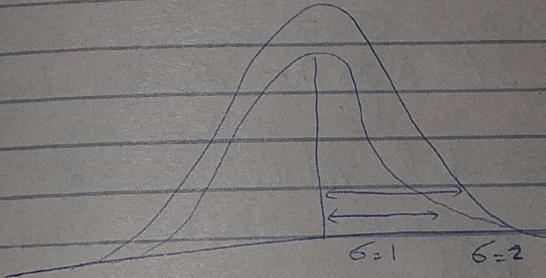
We also update μ_2 , σ_2

$$Z_i = \frac{Z_i - \mu_2}{\sigma_2} \quad \begin{array}{l} \text{standard deviation = 1} \\ \text{Mean = 0} \end{array}$$

$$Z_i = \frac{\alpha Z_i + \beta}{\text{new standard deviation}} \quad \begin{array}{l} \text{new mean} \\ \text{standard deviation} \end{array}$$

instead of having $\sigma=0.1$, $\mu=0$

We multiply it with another standard deviation



So some input requires $\sigma=1$, & some require $\sigma=2$ or some other value.

We learn these α & β .

α^{new}
= Standard deviation

β^{new}
= Mean

Reason:

some input comes from different distribution
so these images have different outputs at neuron
so it affects our learning.

So what we do is send 32 images at each neuron
so apply standardization at each neuron. Now, then

Standard deviation = 1 Mean = 0 ~~so our assumption~~
was these $\sigma=1$, $\mu=0$ will be good.

But Now, we let algorithm decide which input
requires ~~what~~ what values for ~~so that~~ α & β
for current input.

Each Neuron will have its own α & β
Value which it will be learning.

Example. $w_1=0.5$ $w_2=0.5$		Z	$\frac{Z-14}{8}$	$Z_i = \frac{Z_i - 14}{8}$	$Z_i = \alpha Z_i + \beta$
20	10	15	$(15-14)^2 = 1$	$1/6.633$	
30	20	25	$(25-14)^2 = 11^2 = 121$	$121/6.633$	
20	10	15	$(15-14)^2 = 1^2 = 1$	$1/6.633$	
10	10	10	$(10-14)^2 = 16$	$16/6.633$	
5	5	5	$(5-14)^2 = 81$	$81/6.633$	
x_1	x_2				

$$\begin{array}{l} x_1 \xrightarrow{0.5} \\ x_2 \xrightarrow{0.5} \end{array}$$

$M=14$

Loss Functions:

$$L_{\text{sum}} = \sum_{j=0}^M (0, s_j - S_y + 1)$$

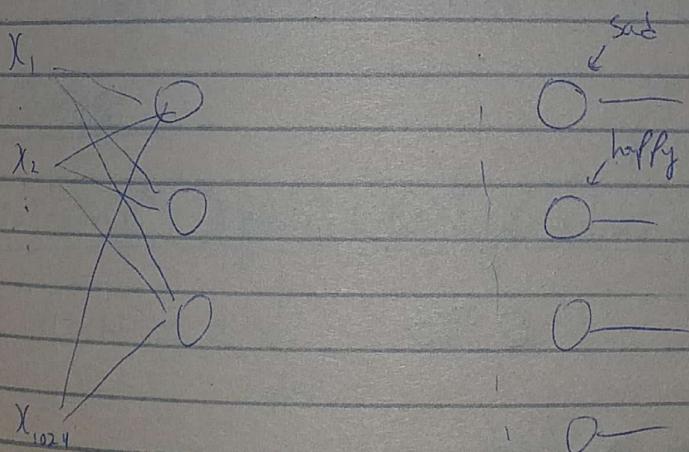
$$L_{\text{softmax}} = -\log \left(\frac{e^{s_y}}{\sum_{j=0}^M e^{s_j}} \right)$$

But Now, if we have multi-label classification, then

our output layers would be total no. of unique labels.

Now, every output has probability 0-1. So we use

Sigmoid Activation function



output layer

We can use softmax loss function,

$$L = \sum_{j \neq y} -y * \log(p_j)$$

1	belongs to this class
0	Doesn't belong to this class.
1	
0	
0	

5 labels

e.g.-

0.8
0.1
0.2
0.75
0.1

After applying sigmoid at output layer.

Now we need loss of this output so we do

$$\begin{array}{|c|c|} \hline 0.8 & 1 \\ \hline 0.1 & 0 \\ \hline 0.2 & 0 \\ \hline 0.75 & 1 \\ \hline 0.1 & 0 \\ \hline \end{array} \rightarrow \text{So} \quad = -1 * \log(0.8) - 1 * \log(0.75) \\ = 0.735$$

This ones the labels
y not Predicted labels.

Regression Problem is more difficult to solve
than classification Problem.

[3 abs]

Next Friday
deadline

Assignment:

DL learning / keras

Take images dataset 1 image = 1 label MLP
fully connected NN

8. Do image classification using sklearn

Tensorflow

Multi-label classification

③ by torch

Regression

input images

we try to convert Regression Problem to classification.

nestro]
Adam Optimizer.

[0.79]

0.0

0.1

0.2

0.3

0.4

0.5

0.6

0.7

0.8

0.9

1.0

2.0

2.1

Optimizer:

① Gradient Descent / Vanilla Gradient Descent

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

In this case, throughout our learning rates doesn't change.

for each weight, learning rate is same.

So those weights whose gradient is large & those weights

whose gradient is small has different learning rate

in each step so small gradient is learning slowly

large gradients are learning quickly.

② AdGrad:

Adaptive Gradient

$$\alpha_{w_i} = \frac{\alpha}{\sqrt{\text{sum}}} \quad | \quad \alpha = 0.01$$

$\beta w_1, \beta w_2, \dots, \beta w_n$

e.g.-

$$\begin{bmatrix} 32 \end{bmatrix} \rightarrow \frac{dL^{(1)}}{dw_1} \rightarrow \text{batch}_1$$

$$\text{next batch} \begin{bmatrix} 32 \end{bmatrix} \rightarrow \frac{dL^{(2)}}{dw_2} \rightarrow \text{batch}_2$$

$$\text{Sum. } \frac{dL}{dW_1}^{(1)} + \frac{dL}{dW_1}^{(2)}$$

we take squared sum of these gradients.

We are taking square bcz gradients could be -ve. Then we take square root to bring magnitude down.

1st problem.

a) throughout learning Sum is increasing in magnitude

b) so learning rate gets closer to zero, so our weights don't get updated as we go through our training. bcz d gets closer to zero.

b) Sum if closer to zero, our α gets increased,

or if sum gets = 0, so division by zero

error could occur