

23-11-21

Week #13 Lecture #16

Tuesday

Question

what will be $\mu \otimes \sigma$ in testing time?

we will save $\mu \otimes \sigma$ for all trains set & use it

in test time we will compute $\mu \otimes \sigma$ for all neurons

If we are not sending test images in batches.

we either send whole train set in NN & complete

$\mu \otimes \sigma$ for all neurons or send a major
portion of Train set in NN & complete $\mu \otimes \sigma$

for all neurons & use it in test time.

Optimizer:-

① Gradient Descent / Vanilla Descent :-

$$w = w - \alpha dw \quad , \quad \frac{dw}{dL}$$

Problem

learning rate

- a) Same value for all gradients.

Some Δw are large

Some Δw are small

So this leads to quick learning of first Δw 's & not
for small Δw 's

② Adam:-

different
 α value for each Δw .

$$\alpha_i = \frac{\alpha \Delta s}{\sqrt{s}} ; \alpha = 0.0001$$

usually in initial

$$S_0 + = \Delta w_i^2 \quad \text{for } w_i \text{ for 32 batch input}$$

So we get 32 Δw we take mean
of these 32 Δw &
this is Δw .

for next 32 images batch we again have

32 Δw Now we take mean of these Δw 's

Then we have Δw_i again.

Now we add it to S

$$S_1 + = \Delta w_i^2$$

whose learning rate was bigger, it will get smaller;

whose learning rate is smaller, it will get bigger.

problem.

\sqrt{S} term gets bigger as α gets smaller so our learning gets slow.

③ RMSProp:-

Instead of dividing by \sqrt{S} , we take avg of gradients & divide by that

$$d_i = \frac{\alpha}{\sqrt{A}}$$

A = Cumulative Moving Avg
time increases

e.g. $\overrightarrow{50 \ 65 \ 60 \ 70 \ 80}$

$$\begin{array}{r} 50 \\ \downarrow \\ 50 \end{array}$$

$$\begin{array}{r} 50+5 \\ \hline 2 \end{array}$$

$$\begin{array}{r} 50+65+60 \\ \hline 23 \end{array}$$

$$\begin{array}{r} 50+65+60+70 \\ \hline 24 \end{array}$$

Problem:

New entry ~~Data Point~~ isn't given much weightage in our avg.

So new ~~derivatives~~ aren't ~~given~~ given ^{much} weightage

Solution:

Simple Moving Avg.

$$50 \quad 65 \quad 60 \quad 70 \quad 80$$

, , , ,

$$\frac{50+65+60}{3} \text{ then}$$

$$\frac{65+60+70}{3} \text{ then}$$

$$\frac{60+70+80}{3}$$

Window size = 3

$$A = \frac{dw_1^{2[t_0]} + dw_1^{2[t_1]} + dw_1^{2[t_2]}}{3}$$

t_0 = batch 1

t_1 = batch 2

t_2 = batch 3

Now if batch 4 comes.

$$A = \frac{dw_1^{2[t_3]} + dw_1^{2[t_2]} + dw_1^{2[t_3]}}{3}$$

Now taking Avg of last
3 derivatives of weights.

In RMS Prof, we use Decaying Moving Average

b/cz Problem in simple Moving Avg is that we are ignoring previous history of the derivatives

Decaying Moving Avg.

$$A_{t+1} = \beta A_t + (1-\beta) d\omega_t^2$$

$$\beta = 0.9 \text{ usually}$$

A_t = Previous Avg

initially it could be 30 OR $A_0 = d\omega_1^2 [t_0]$

$$A_3 = \beta A_2 + (1-\beta) x_3$$

$$A_3 = \beta (\beta A_1 + (1-\beta) x_2) + (1-\beta) x_3$$

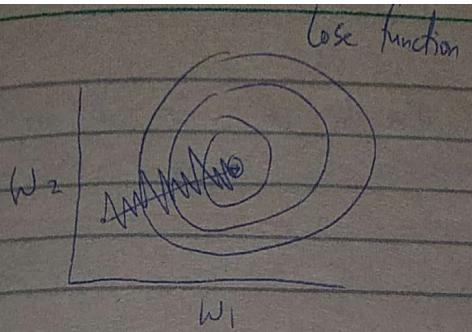
$$A_3 = \beta (\beta (\beta A_0 + (1-\beta) x_1) + (1-\beta) x_2) + (1-\beta) x_3$$

so historical Avgs are multiplied by β^n when $\beta < 1$

so β^n becomes small & historical Avgs get less

weights.

④ Gradients with Momentum:



Ideally we need to learn our weights in linear form
 in a straight line but our actual weights are
 updated in zigzag pattern
Oscillations Reason
 bcz we are taking derivative of one weight we are
 assuming that others weights are constant and only one
 weight is updated but in reality other weights are
 also updated at same time so our learning of
 weights is not right.

2nd Problem.

We update weights on batches. So one batch might
 have different gradients than others so we have oscillations
 in our loss function.

Goal:
Oscillations should be less in our improving $\mathcal{L}_{\text{loss}}$
function.

Solution-

$$w_i = w_i - \alpha (\text{change})$$

$$w_i = w_i - \alpha (c) \quad c = \text{change}$$

$$C_{t+1} = \beta C_t + (1-\beta) d\omega_i$$

↙

by doing this, we get to know in what direction batch

Avg are going so we update weights according

Initially, $C_0 = 0$ or $C_0 = d\omega_i^{[0]}$

problem usually we use 0.

In this, α value is constant.

⑤ ADAM:-

It has adaptive gradients & we also lower oscillations in this

$$w_i = w_i - \alpha_i C_t$$

$$C_{t+1} = \beta C_t + (1-\beta) d\omega_i^{[t]}$$

$$\alpha_i = \frac{\alpha}{\sqrt{A_i}}, \quad A_i = \beta_0 A_{i-1} + (1-\beta_2) d\omega_i^2$$

C_t Force β_i is called Momentum

Convolution:-

We move a filter across the image & do some process, then at the end we get a filtered image.

e.g.

100	100	100
100	0	100
100	100	100

Image Part

1	1	1
1	1	1
1	1	1

Filter

We do element wise multiplication so we get

100×1	100×1	100×1
100×1	0×1	100×1
100×1	100×1	100×1

100	100	100
100	0	100
100	100	100

8 Then we sum up all the values in that result

& also take sum of all values in filter &

divide sum of result by sum of filter

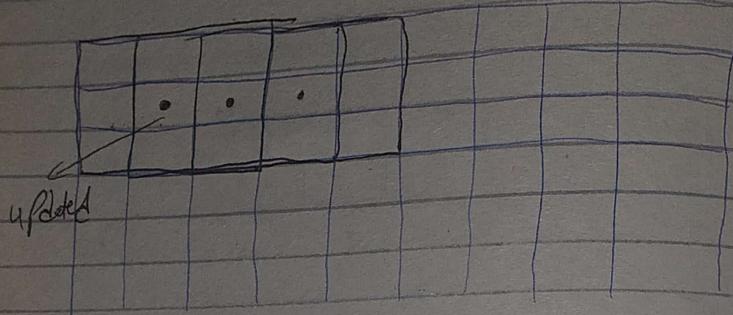
$$= \frac{800}{9}$$

= 88.889 we update 0 in original ^{Imod} part

& write this result.

100	100	100
100	88.89	100
100	100	100

This as a result reduces noises in the image.



we move Filter & process the image