

Activations function in depth 1

→ In the last class we discuss different activation function and talked about where to use which activation function.

I gave you kind of some general rules

→ Today we will talk about optimisers

① we ~~also~~ have already discuss about 1 optimiser gradient descent.

→ write the formula for gradient descent

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$$

Gradient descent is most the simplest optimisers. and even gradient descent 3 types.

① Batch gradient descent

② Mini Batch gradient descent

③ Stochastic gradient descent

→ Thank ~~you~~ fully, the formula for all this gradient descent doesn't change

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$$

① Batch gradient descent → So, what batch gradient descent does is, it calculates the differentiation that is $\frac{dL}{dw}$ for all the data points.

Let say we have 10000, data points, so, it will calculate differentiation for 10000 points.

→ Once it has calculated the differentiation for all the 10000 data points it averages it out and update the weight.

Q → So, what is the problem with batch gradient descent.

→ Since it has 10000 data points and it calculate the gradient for all this 10000 data points it becomes

kind of very slow.

→ See, we have to update the weight for every epoch. and for each epoch, if we calculate the gradient descent for 10000 data points, this becomes very very slow.

→ Now, I have just taken 10000 data points, but some time in training dataset we have million, billions. in that case it will take a very long time just to update the weight.

So, we do not want that

Hence, we have mini Batch gradient descent, Now as the name suggest it will not use entire data point. So

$$MB. GD \rightarrow \frac{dL}{dw}_{100}$$

So, MB GD will do, from 10000 data point it will randomly select the 100 data points, it will calculate the differentiation of weights of this 100 data point and it will update the weights

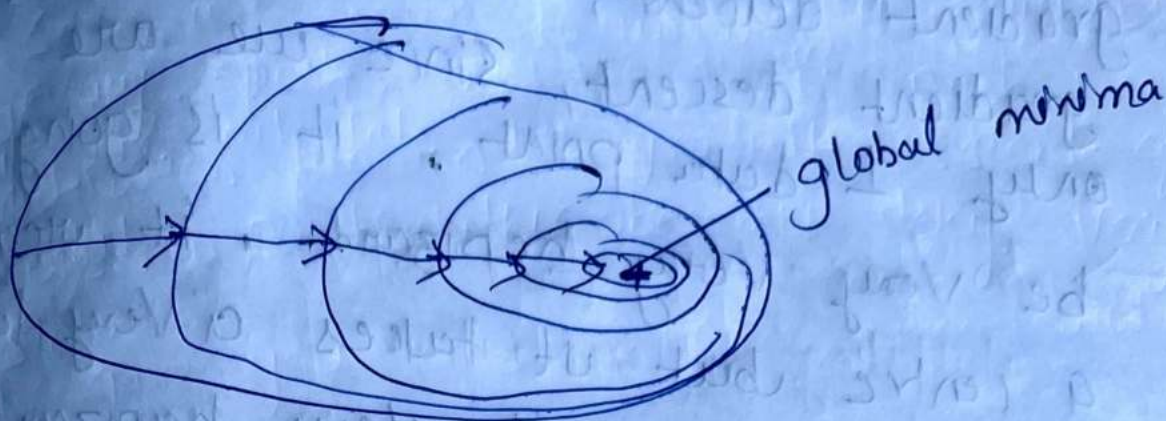
again ~~the~~ itself has some problem
we will talk about later,

then we have Stochastic Gradient
descent

SGD \rightarrow ~~in SGD~~ in SGD they
randomly select 1 data points
and update the weights based
on that.

So, the differentiation is calculated
for 1 data points and weight
are change based on the differentiation
of that 1 data points.

Batch gradient descent \rightarrow



Here you can see this indicates a loss function going inside towards the min.

So, this is Batch gradient descent since we average out for all the data points it is very very smooth.

So, it goes very smoothly & straightly towards the centre.

in MBGD and SGD since we are not using all the data points, the different.

ation $\frac{dL}{dw}$ for minibatch will be quite equal to differentiation $\frac{dL}{dw}$

for batch gradient descent but

it will not be the same

$$MB \frac{dL}{dw} \sim \frac{dL}{dw} BGD$$

and some thing happens with stochastic gradient descent, so in stochastic gradient descent, since we are using only 1 data point, it is going to be very very haphazard, it reaches a centre but it takes a very long time, and it's very haphazard.

→ So, what can we do, so, that we don't have to use some many data points, and it requires so much of time, but at the ~~same time, we~~

→ So, we use some thing known as stochastic gradient descent or minibatch gradient descent with momentum

→ So, ^{new} term is added i.e. the momentum term,

by the way I forgot to mention but mini batch gradient descent is most used gradient descent among all this 3. I think

that is pretty understandable because
that line was not that happy and
it does not take this much of time.

→ So if you see in the graph. there
are two figures, this SLD, and
SLD with momentum,

SLD → Here you can see, this very
happy. b.

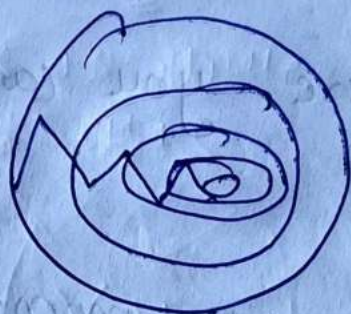
~~SLD~~ SLD with momentum → Here this
is happy but it's not as happy as
one with ~~no~~ momentum

→ So, what is the very big idea ~~behind~~
behind with momentum

→ So, what momentum says. ~~no interest~~
~~of using~~ ~~instead of using~~
when differentiating we are using
 $\frac{dL}{dw}$. instead of using $\frac{dL}{dw}$
differentiation for the last data
point. let us calculate

differentiation for all the ~~other points~~
other step as well.

→ so, for any 1 datapoint instead of
calculating differentiation for that point
let us, somehow add this data
points as well into that



So, that means it will average out
the points and won't be that jagged

→ so, let me first change the
notation

$$= w_{\text{new}} = w_{\text{old}} - W \frac{dL}{dw_{\text{old}}}$$

instead of new an old
let us just write iteration number
 t^{th} iteration, any
iteration before t , is $t-1$
iteration, the old will

replace with $(t-1)$ and anything before that will be $(t-2)(t-3)$ and so, on

$$w_t = w_{t-1} - \eta \frac{dL}{dw_{t-1}}$$

the reason why i am writing this is, because it becomes easier to write instead of writing new and old all the time. Let us just represent this with a iteration number,

and instead of writing $\frac{dL}{dw}$ all the time, let me just replace with g_t (that is gradient tag and t^{th} point)

So, new equation becomes something like this

$$w_t = w_{t-1} - \eta g_t$$

→ Now, just let's get back to the formula. So, what happens is, the formula remains the same, we just have a new

term. so, that means

$$w_t = w_{t-1} - \cancel{\gamma} \cancel{v_t} \tilde{v}_t$$

$$\rightarrow v_t = \gamma v_{t-1} + W g_t,$$

where,

$$\cancel{w_t} v_t = W g_t$$

\rightarrow I know this is a little complex, but expand this

\rightarrow

$$w_t = w_{t-1} - (\gamma v_{t-1} + \gamma g_t)$$

with the formula above so this will be at

$$= w_{t-1} - W g_t - \gamma \underbrace{v_{t-1}}_{\text{if}}$$

we replace v_{t-1} :

$$= w_{t-1} - W g_t - \gamma (v_{t-1} + W g_{t-1})$$

and this goes on an on. honestly, and this comes out - some thing like this

$$= w_{t-1} - W g_t - \gamma W g_{t-1} - \gamma^2 W g_{t-2}$$

$$\cancel{w_t = w_{t-1} - (\gamma w_{t-1} + \gamma g_{t-1})}$$

$$w_t = w_{t-1} - \gamma g_t - \gamma \gamma g_{t-1} - \gamma^2 \gamma g_{t-2} - \gamma^3 g_{t-3} - \gamma^4 g_{t-4} \dots$$

$$- \gamma^T g_1$$

→ So, what happens if's. the γ (gamma)
if's always $[0, 1]$

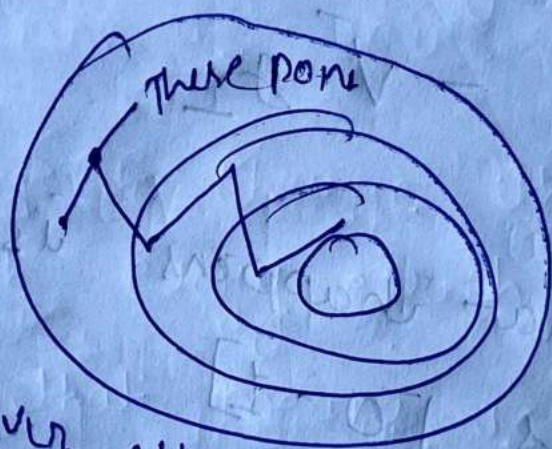
let's assume $\gamma = 0.1$

weight = 1

$$w_t = w_{t-1} - 1 - \gamma g_t - (\gamma)^2 \gamma g_{t-1} - (\gamma)^2 \gamma^2 g_{t-2}$$

At this t^{th} iteration the weight
if's one, for the iteration before
that we have the weight of 0.1
so before that we have square
(0.1)² and this goes on.

→ So, you can understand pattern, ~~what~~ happening is, you give the most weightage to the point for which it is being calculated but you also give some weightage here to these points.



→ However the weightage keeps on decreasing. So for the t th iteration you have a weightage of 1 for $t-1$ iteration you have weightage of 0.1 for gamma = 0.1

= for $t-2$ iteration you have weightage $t-2 = 0.01$

and this goes on and on.

→ So, what happens is, since you are giving weightage to the previous point it does not become it's happened



it kind averages out depending upon the position and it's happened yes, but it less happened now.

Momentum → momentum in short it's basically means giving weightage to the previous calculated gradient descent; so that the next gradient descent as happened as the previous one.

* Some other optimisers

① Adagrad if you remember the equation of gradient descent is

$$w_t = w_{t-1} - \eta \cdot g_t$$

~~the~~