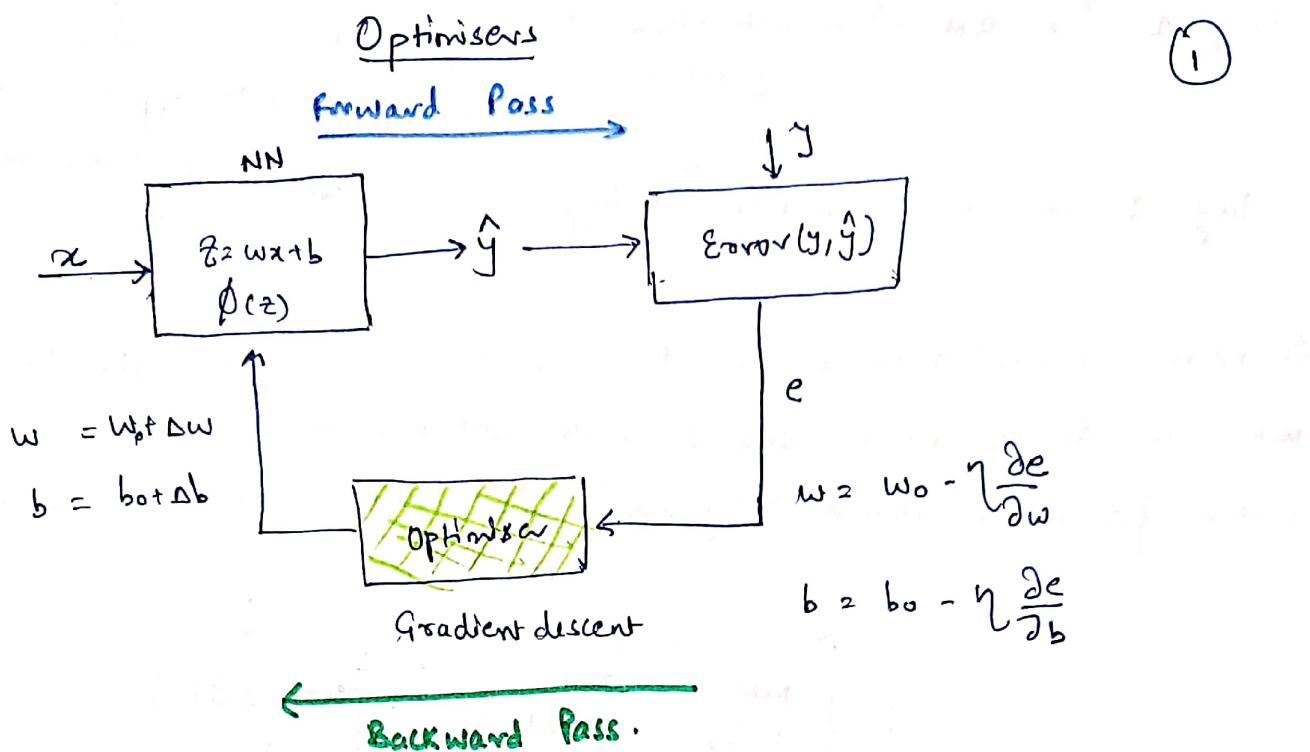


(1)



Drawback of Gradient Descent

$\Delta w \propto \frac{\partial c}{\partial w}|_{w=old}$ hence this is only dependent on current weight update

It does not get influenced by previous updates.

↓ moves very slow or may get stuck here

slow training

Observation

(i) If $\Delta w \uparrow$ if $\frac{\partial c}{\partial w} \uparrow$

If slope steep at that point weight update will be high

$\Delta w \downarrow$ if $\frac{\partial c}{\partial w} \downarrow \rightarrow$ flat surface

Simplified reprn of gradient descent

$$w_{new} = w_{old} - \eta \frac{\partial c}{\partial w} \quad |_{w=w_{old}}$$

or

$$\theta_{new} = \theta_{old} - \eta \nabla_{\theta} J(\theta) \quad |_{\theta=\theta_{old}}$$

$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ b_1 \\ b_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} w_{1,old} \\ \vdots \\ b_{1,old} \\ \vdots \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial J}{\partial w_1} \\ \vdots \\ \frac{\partial J}{\partial w_n} \\ \frac{\partial J}{\partial b_1} \\ \vdots \\ \frac{\partial J}{\partial b_n} \end{bmatrix}$$

$$\Delta w = -\eta \frac{\partial c}{\partial w}$$

$$\Delta w \propto \frac{\partial c}{\partial w} \quad |_{\theta=\theta_{old}}$$

slope of tangent

(iii) chain rule at lower layers includes multiple $\frac{\partial c}{\partial w}$ terms and hence

②

$$\frac{\partial e}{\partial w_r} = \frac{\partial e}{\partial a_r} \cdot \frac{\partial a_r}{\partial z_r} \cdot \frac{\partial z_r}{\partial w_r} \quad [\text{lower weight update}]$$

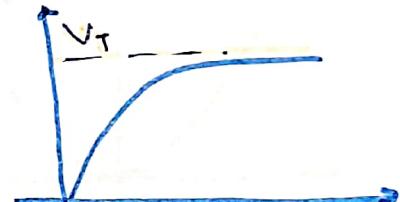
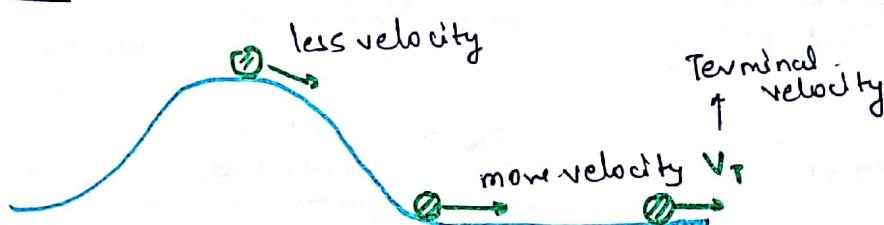
This will train very slowly at lower layers as compared to first one

Type of Optimizers

- Fast optimiser
 - 1) Momentum Optimisation
 - 2) Nesterov accelerated gradient
 - 3) Adaptive gradient
 - 4) RMS Prop
 - 5) Adam Optimisation.

Momentum Optimisation (Boris Polyak) 1964

Idea



It tries to accumulate past gradient as well so that you reach minima faster.

Algorithm

$$m_{\text{new}} \leftarrow \beta m_{\text{old}} + \gamma \nabla_{\theta} J(\theta)$$

$$\theta \leftarrow \theta - m$$

or

$$m_{\text{new}} = \beta m_{\text{old}} + \gamma \frac{\partial c}{\partial w} \quad |_{w=\text{current}}$$

$$w_{\text{new}} = w_{\text{old}} - m$$

bias update will also happen in same way

$$\theta = \theta - [\beta m + \gamma \nabla_{\theta} J(\theta)]$$

β = coefficient of momentum

analogous to coefficient of friction.

(3)

$$\beta = 0$$

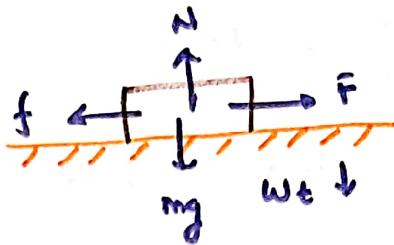
\Rightarrow frictionless surface

$$m = \beta m^0 + \eta \frac{\partial c}{\partial w} \Big|_{w=\text{current weight}}$$

$$= \eta \frac{\partial c}{\partial w} \Big|_{w=\text{current weight}}$$

$$\mu = 0$$

$$f = 0$$



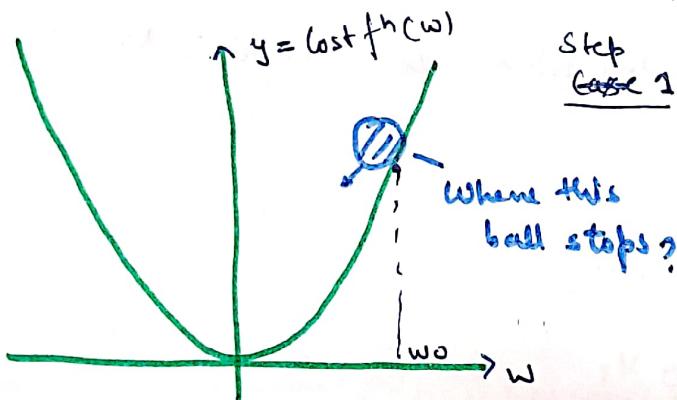
$$f = \mu N = \mu mg$$

$$w = w - m = w - \eta \frac{\partial c}{\partial w} \Big|_{w=\text{current weight}}$$

This is normal gradient descent if $\beta = 0$

What is value of β ?

It is experimentally proven in most of the cases that $\beta = 0.9$ but you can take any other value as well and not bounded by it.



Step Case 1 :-

Initially $w=w_0$, $m_0=0$, $\beta=0.9$

$$m_1 = \beta m_0^0 + \eta \frac{\partial c}{\partial w} \Big|_{w=w_0}$$

$$m_1 = \eta \frac{\partial c}{\partial w} \Big|_{w=w_0}$$

$$w_1 = w_0 - m_1 = w_0 - \eta \frac{\partial c}{\partial w} \Big|_{w=w_0}$$

like a gradient descent

Step 2 :-

$$m_1 = \eta \frac{\partial c}{\partial w} \Big|_{w=w_0}, \quad \beta = 0.9, \quad w=w_1$$

$$m_2 = \beta m_1 + \eta \frac{\partial c}{\partial w} \Big|_{w=w_1}$$

$$m_2 = \underbrace{\left(\beta \eta \frac{\partial c}{\partial w} \Big|_{w=0} \right)}_{\checkmark} + \eta \frac{\partial c}{\partial w} \Big|_{w=w_1}$$

it considers past weight updates and does not solely depends on the current weights.

(4)

Step n :

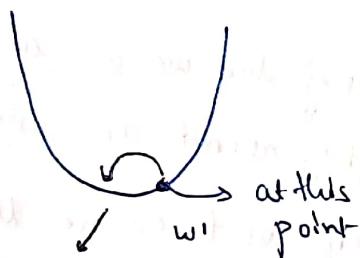
$$m_n = \eta \left[\beta^{n-1} \frac{\partial c}{\partial w} \Big|_{w=w_0} + \beta^{n-2} \frac{\partial c}{\partial w} \Big|_{w=w_1} + \dots + \frac{\partial c}{\partial w} \Big|_{w=w_{n-1}} \right]$$

It takes more fraction of latest one & less fraction of previous one.

Drawback

- ① One extra parameter β but $\beta=0.9$ works well for most of the cases.

②



oscillation

Issue in backprop

It oscillates closer to

global or local minima due

to accumulation of fast momentum

$$\eta \left[\beta^{n-1} \frac{\partial c}{\partial w} \Big|_{w=w_0} + \dots + \frac{\partial c}{\partial w} \Big|_{w=w_1} \right]$$

This value is
already very less

Hence in this case past
gradient will dominate over
the current gradient

Advantages

- ② Oscillation can also help to come out of local minima

① Fast Convergence

Keras

optimizer = tf.keras.optimizers.SGD(lr=0.01, momentum=0.9)

AdaGradAdaptive Gradient

6

Nesterov Accelerated Gradient

5

(NAG) - Yuri Nesterov, 1983

aka Nesterov Momentum Optimisation
as it's faster than momentum optimisation

What's new?
 It calculates gradient slightly ahead in the direction of momentum.
 Instead of calculating at θ but $\theta - \beta m$ or $w - \beta m$.

Algorithm

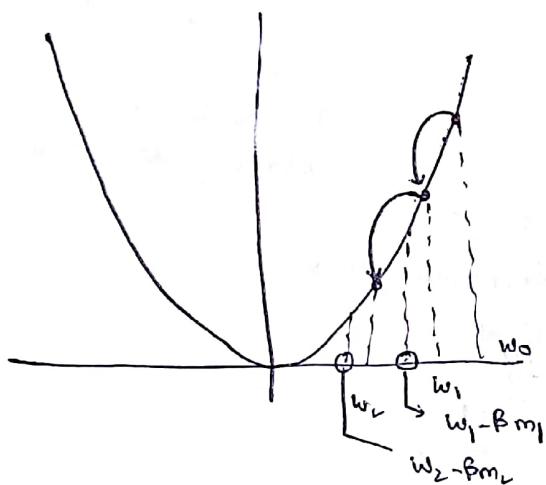
$$m = \beta m + \eta \frac{\partial C}{\partial w} \quad |_{w=(w-\beta m)}$$

$$w = w - m$$

$$m \leftarrow \beta m + \eta \nabla_C(\theta - \beta m)$$

or

$$\theta \leftarrow \theta - m$$



It is also accumulating but we are also looking for terms slightly ahead and knowing what exists. This helps us to ensure that we are actually converging faster and also know what exists ahead.

Advantages

- ① faster than momentum optimisation, it gives us less oscillations. Oscillations are significantly reduced at the minimum.
- ② Reaches close to minima because of less oscillations.

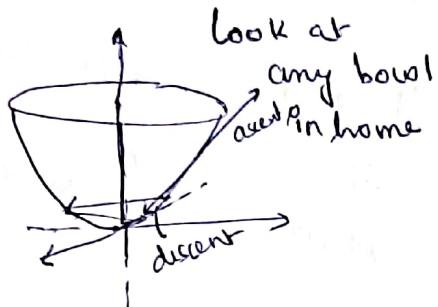
Disadvantages

- ① Extra parameter to tune i.e. β but as you know $\beta = 0.9$ works well in most of the cases.
- * optm... = Et. Keras. optimizers.SGD(lr= , momentum= 0.9, nesterov= True)

(6)

AdaGradAdaptive Gradient

It solves one of the problems of gradient descent or the fast optimisers that we have seen. \rightarrow The problem is known as **Elongated Bowl Problem**



$$w = w - \eta \frac{\partial c}{\partial w}$$

$$w = w - \nabla C$$

\hookrightarrow This negative sign leads to gradient descent

Sometimes this negative descent in an elongated bowl may lead us to a divergence track rather than convergence and this particular problem is solved by ada grad.

~~"Put around - others to MATH TEST App"~~

The $-\nabla C$ does not always focus towards the direction of global minimum.

Dot product :- $\vec{v} = x^i \hat{i} + y^j \hat{j}$ $\vec{p} = a^i \hat{i} + b^j \hat{j}$

Scalar Product $\vec{v} \cdot \vec{p} = xa + yb = |\vec{v}| |\vec{p}| \cos \theta$

Algorithm (Adaptive Gradient)

* It corrects the direction in the initial steps. \rightarrow scalar term

$$s \leftarrow s + \nabla_{\theta} J(\theta) \otimes \nabla_{\theta} J(\theta)$$

$$s = s + \left(\frac{\partial c}{\partial w} | w \right)^2$$

$$\theta \leftarrow \theta - \eta \nabla_{\theta} J(\theta) \oslash \sqrt{s + \epsilon}$$

$$\theta = \theta - \eta \left(\frac{\partial c}{\partial w} \right) w / \sqrt{s + \epsilon}$$

s = scaling down factor

ϵ = smoothing term to avoid zero division error = 10^{-7}

\otimes = element wise multiplication + dot product

\oslash = element wise division

(7)

$$s \leftarrow s + \frac{\nabla_{\theta} J(\theta)}{\vec{v}} \otimes \frac{\nabla_{\theta} J(\theta)}{\vec{v}}$$

$$\theta \leftarrow \theta - \eta \frac{\frac{\nabla_{\theta} J(\theta)}{\sqrt{s + |\vec{v}|^2 + \epsilon}}}{\sqrt{s + |\vec{v}|^2 + \epsilon}}$$

* if $s_0 = 0$, $\epsilon \approx 0$ (1st iteration)

* After some iteration

$$\frac{\vec{v}}{\sqrt{|\vec{v}|^2}} = \hat{\vec{v}} \text{ unit vector}$$

$$\sqrt{s + |\vec{v}|^2 + \epsilon}$$

$$|\vec{v}'| < |\vec{v}|$$

Step by step

Step 1: Initially $s_0 = 0$, $w = w_0$

$$s_1 = \vec{v}_0^T + \left(\frac{\partial C}{\partial w} \Big|_{w=w_0} \right)^2 - ①$$

$$w_1 = w_0 - \eta \frac{\left(\frac{\partial C}{\partial w} \Big|_{w=w_0} \right)}{\sqrt{\left(\frac{\partial C}{\partial w} \Big|_{w=w_0} \right)^2 + \epsilon}} - ②$$

$\epsilon \approx 0$

$$w_1 = w_0 - \eta \frac{\left(\frac{\partial C}{\partial w} \Big|_{w=w_0} \right)}{\left| \frac{\partial C}{\partial w} \Big|_{w=w_0} \right|}$$

$$w_1 = w_0 - \eta \hat{\vec{v}}$$

$\therefore \eta \in (0, 1) \Rightarrow$ overall weight update will be very less.

Step 2:

(Q)

$$S_2 = S_1 + \left(\frac{\partial c}{\partial w} \Big|_{w=w_0} \right)^2$$

$$w_2 = w_1 - \eta \frac{\frac{\partial c}{\partial w}}{\sqrt{S_2 + \epsilon}}$$

$$S_2 = \underbrace{\left(\frac{\partial c}{\partial w} \Big|_{w=w_0} \right)^2}_{\text{past}} + \underbrace{\left(\frac{\partial c}{\partial w} \Big|_{w=w_1} \right)^2}_{\text{present}}$$

Now we are taking 100%.

of it unlike the previous cases.

larger term than the previous one

Step n

$$S_n = \sum_{i=0}^n \left(\frac{\partial c}{\partial w} \Big|_{w=w_i} \right)^2$$

S_n is very large due to accumulation compared to S_1

$$w_{n+1} = w_n - \eta \frac{\frac{\partial c}{\partial w} \Big|_{w=w_n}}{\sqrt{S_n + \epsilon}}$$

$$\Delta w \propto \frac{\frac{\partial c}{\partial w} \Big|_{w=w_n}}{\sqrt{\sum_{i=0}^n \frac{\partial c}{\partial w} \Big|_{w=w_i}} + \epsilon} \Rightarrow \Delta w \propto \frac{\frac{\partial c}{\partial w} \Big|_{w=w_n}}{K}$$

Assumption $K = \left\{ \sqrt{\sum_{i=0}^n \frac{\partial c}{\partial w} \Big|_{w=w_i}} + \epsilon \right\}$

$$w = w - \eta \frac{\frac{\partial c}{\partial w} \Big|_{w=w_n}}{K}$$

$$w = w - \eta \frac{\frac{\partial c}{\partial w} \Big|_{w=w_n}}{K}$$

varying term

overall it is similar to gradient descent with learning rate

$$\frac{\eta}{K} \rightarrow \text{decaying learning rate}$$

K will keep increasing

This is the property of Adagrad

⑨

Hence due to ada grad the jump will be lower and zig-zag fashion can be avoided. It is trying to self correct the direction.

Gradient Descent

- Moves quickly in direction of deepest descent without considering where global minima is located.

Ada Grad

- It corrects its direction by scaling down the gradient vector along the steepest direction.

Advantages of Ada grad

- 1) Corrects the direction ~~Fully~~ of gradient vector partially unlike the gradient descent.
- 2) Less tuning of learning rate is required because of decaying factor.

Disadvantages

- 1) It oftenly stops early before reaching global minima
- 2) It takes longer time to converge due to decaying learning rate.

~~xx~~ Even though Ada grad has direction correction advantage but due to point ① in disadvantage it is not recommended to use.

RMS Prop

↓ ↳ Geoffrey Hinton et al.

10

Tackles the drawback of AdaGrad - early stopping problem by accumulating gradient from recent iteration and using exponential decay

Algorithm

$$\begin{aligned} s &\leftarrow \beta s + (1-\beta) \nabla_{\theta} J(\theta) \otimes \nabla_{\theta} J(\theta) \\ \theta &\leftarrow \theta - \eta \nabla_{\theta} J(\theta) \otimes \sqrt{s + \epsilon} \end{aligned}$$

$$s = \beta s + (1-\beta) \left(\frac{\partial C}{\partial w} \right)_w^2$$

$$w = w - \eta \left(\frac{\partial C}{\partial w} \right)_w \Big| \sqrt{s + \epsilon}$$

β = coefficient value = 0.9

experimentally proven not
any theoretical basis.

Smoothly
term

RMS Prop

$$w_1 = w_0 - \eta \cdot \frac{\frac{\partial C}{\partial w}|_{w=w_0}}{\sqrt{\left(\frac{\partial C}{\partial w}|_{w=w_0} \right)^2 (1-\beta) + \epsilon}}$$

Smaller
denominator

$\Delta w_{RMS} > \Delta w_{Adagrad}$

AdaGrad

$$w_1 = w_0 - \eta \frac{\frac{\partial C}{\partial w}|_{w=w_0}}{\sqrt{\left(\frac{\partial C}{\partial w}|_{w=w_0} \right)^2} + \epsilon}$$

large
denominator

* decay in learning rate is also low.

Keras

tf. keras.optimizer.RMSProp (learning_rate = 0.01, rho = 0.9)

* Adam ↗ Adaptive momentum estimation. ↗ combines the idea of

Momentum optimisation
RMS prop
exponential decay

Algorithm

1. $m \leftarrow \beta_1 m + (1-\beta_1) \nabla_{\theta} J(\theta)$
2. $s \leftarrow \beta_2 s + (1-\beta_2) \nabla_{\theta} J(\theta) \otimes \nabla_{\theta} J(\theta)$
3. $\hat{m} \leftarrow \frac{m}{1-\beta_1 t}$ ↗ Iteration step.

4. $\hat{s} \leftarrow \frac{s}{1-\beta_2 t}$

5. $\theta \leftarrow \theta - \eta \hat{m} \otimes \sqrt{\hat{s} + \epsilon}$

- or —
1. $m = \beta_1 m + (1-\beta_1) \frac{\partial c}{\partial w} \Big|_{w=w}$ → Inspired from momentum optimisation
Exponential decay of gradients.
 2. $s = \beta_2 s + (1-\beta_2) \left(\frac{\partial c}{\partial w} \Big|_{w=w} \right)^2$ — $\beta^3 = 0.729$ $\beta^2 = 0.81$ $\beta = 0.9$
 3. $\hat{m} = \frac{m}{1-\beta_1^t}$
 4. $\hat{s} = \frac{s}{1-\beta_2^t}$
 5. $w = w - \eta \frac{\hat{m}}{\sqrt{\hat{s} + \epsilon}}$
- 1 1 1 1 1
- Inspired from RMS prop - Exponential decay of squared gradients
- Bias correction step:-- zero bias
Correction step
→ Inspired from Time series smoothing operation

Step ③ & Step ④

Initial condition: $m_0 = 0$, $s_0 = 0$, $t_0 = 0$

$$\beta_1 = 0.9, \beta_2 = 0.999 \quad \epsilon = 10^{-8}$$

$$m_1 = \beta_1 m_0 + (1-\beta_1) \frac{\partial c}{\partial w} \Big|_{w=w_0} \quad s_1 = \beta_2 s_0 + (1-\beta_2) \left(\frac{\partial c}{\partial w} \Big|_{w=w_0} \right)^2$$

$$m_1 = (1-\beta_1) \frac{\partial c}{\partial w} \Big|_{w=w_0} \quad s_1 = (1-\beta_2) \left(\frac{\partial c}{\partial w} \Big|_{w=w_0} \right)^2$$

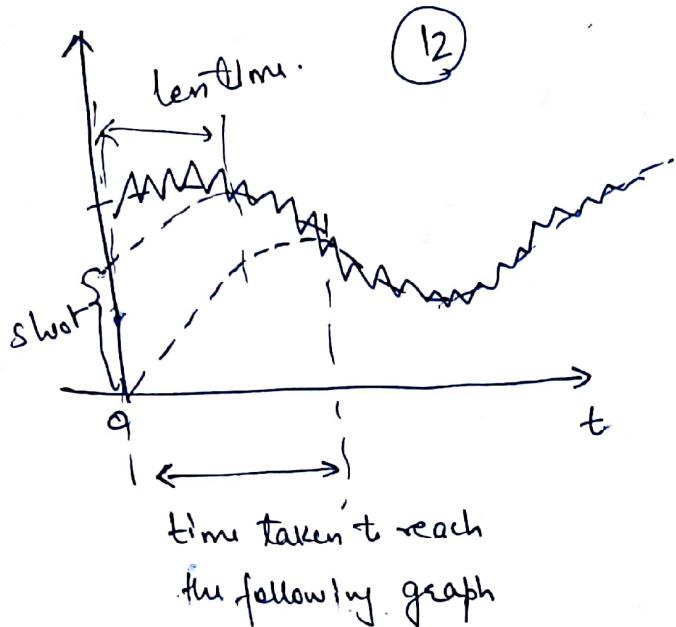
Let's assume $\frac{\partial c}{\partial w} \Big|_{w=w_0} = 0$

$$m_1 = (1-0.9) \times 10 = 1 \quad s_1 = (0.001) \times 100 \\ = 0.1$$

m_1 and s_1 both are close to zero

$$\hat{m}_1 = \frac{m_1}{1-\beta_1^t} = \frac{m_1}{1-(0.9)^t} = \frac{1}{0.1} = 10$$

$$\hat{s}_1 = \frac{s_1}{1-\beta_2^t} = \frac{0.1}{0.001} = 100$$



In this example of time series, we need to create a smoothing function. \hat{m}_t & \hat{s}_t shoot it far away from zero so that \hat{q}_t is able to follow along the graph q_t in very few iterations.

- * This is why we use a zero bias correction step.

Observation:- Since, we are starting m_0 and s_0 from zero due to this in initial few steps they will be biased towards zero. So to avoid that we are shooting q_t away from zero using zero bias correction step.

Expanding Adam Algorithm.

$$\hat{m}_t = \frac{m_t}{1 - \beta_1} = \frac{(1 - \beta_1) \frac{\partial C}{\partial w} \Big|_{w=w_0}}{(1 - \beta_1)} = \frac{\frac{\partial C}{\partial w} \Big|_{w=w_0}}{1 - \beta_1}$$

$$\hat{s}_t = \frac{s_t}{1 - \beta_2} = \frac{(1 - \beta_2) \left(\frac{\partial C}{\partial w} \Big|_{w=w_0} \right)^2}{(1 - \beta_2)} = \left(\left(\frac{\partial C}{\partial w} \Big|_{w=w_0} \right)^2 \right)$$

$$w_t = w_0 - \eta \frac{\frac{\partial C}{\partial w} \Big|_{w=w_0}}{\sqrt{\left(\frac{\partial C}{\partial w} \Big|_{w=w_0} \right)^2 + \epsilon}} \quad] \text{ like Adagrad.}$$

Step 2

$$t_2 = t_1 + 1 = 2 \quad w_2 = w_1, \beta_1 = 0.9, R_2 = 0.991$$

13

$$w_2 = w_1 - \eta \frac{\hat{m}_2}{\hat{s}_2 + \epsilon}$$

$$\Delta w = \frac{\eta}{K} \hat{m}_2 \approx \underbrace{\eta}_{\text{varying number}} \underbrace{\hat{m}_2}_{\text{varying learning rate}}$$

varying number varying learning rate.

Q. When to use what algorithm?

Class of Algorithm

Convergence speed

Convergent quality

Good for simple ML data - set	$\begin{cases} \text{SGD} \\ \text{Momentum Opt} \\ \text{NAG} \end{cases}$	2	3
		2	3
		2	3

Not recommended

AdaGrad	2	1 (stop early)
---------	---	----------------

to use

RMS Prop	3	2-3
(Adam)	3	2-3

good for complex datasets.
like images, video, etc.

Best optimiser in terms of convergence speed and quality. Most commonly used optimiser.