

# Optimizers

67

## Introduction

### Performance of NN

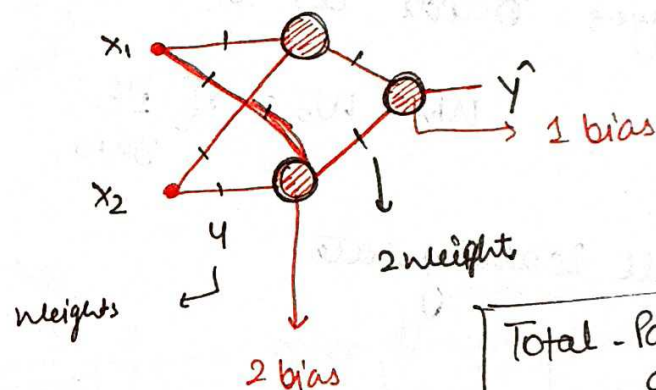
↳ How to speed the training

#### Technique

- Weight init
- Batch Norm
- Activation function
- optimizers

### Role of optimizer

Gpa	iq	placed
8	80	1
9	90	1
7	70	0



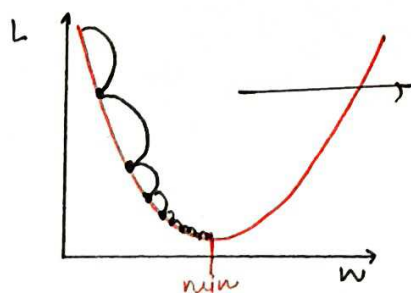
Total - Parameter is 9

Hum 9 parameters ki aisi value dhundi hai ki jab hum data NN mai dale to output close to real output se.

$$\text{Loss} \rightarrow \hat{y} - y$$

final value  $\rightarrow$  final value  
 $w, b$  ki value ki optimum value hai.

Start with Random Value ( $w, b$ )  $\rightarrow$  end with Optimize Value ( $w, b$ )



Optimize value  
 Optimize technique  $\rightarrow$  (gradient descent)

Gradient Descent Update rule:

$$W_n = W_0 - \eta \frac{\partial L}{\partial W}$$

Types of Gradient descent

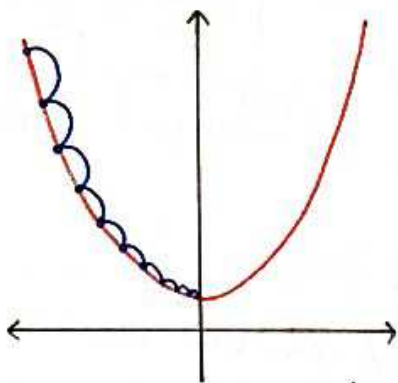
- Batch GD
- Stochastic GD
- Mini-batch GD

Challenge No. 1

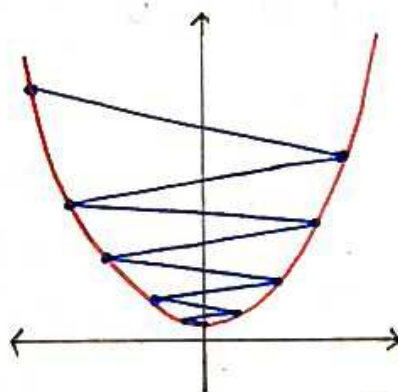
1. challenge -> Decide correct ~~value~~ learning rate value.

$$W_n = W_0 - \eta \frac{\partial L}{\partial W_0}$$

Small learning rate



Large learning rate



2) Solution of this problem is

→

Learning rate scheduling

Not perfect with every dataset

Pre-define

threshold → some range

means before training define

[kis time pe LR ghatega aur kabhi badhaye]

3) Learning rate same  $\rightarrow$  every direction.

68

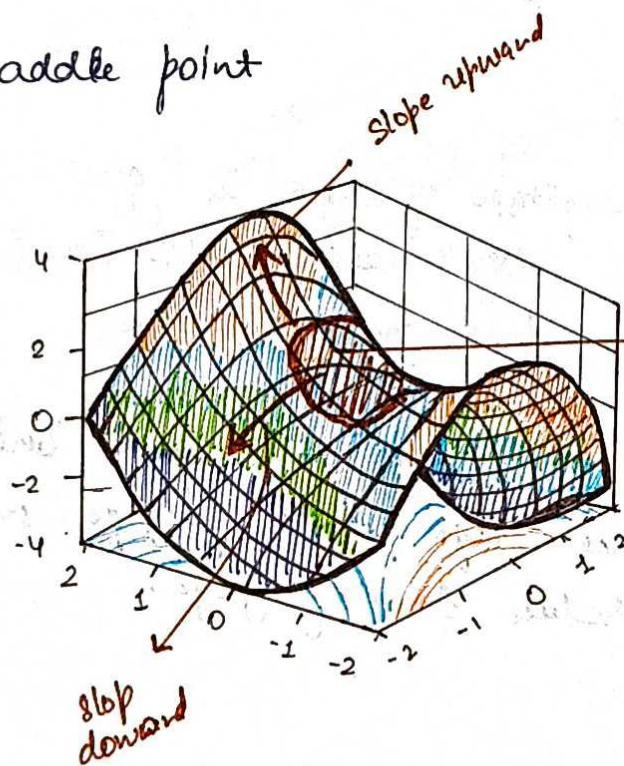
for eg:- 9 weights and bias  
 $\rightarrow$  graph  $\rightarrow$  10-D

If one weight want to decrease the value with high speed and another weight want low speed for reach best solution. This statement is true if learning for weight one is different to learning rate of another weight. This is not possible.

$\hookrightarrow$  LR Separate  
Learning rate is same for every weight.

4. Local minima

5. Saddle point



every direction slope is same  
(not change in slope)

which means  $\frac{\partial L}{\partial W} = 0$

$\hookrightarrow$  differentiation

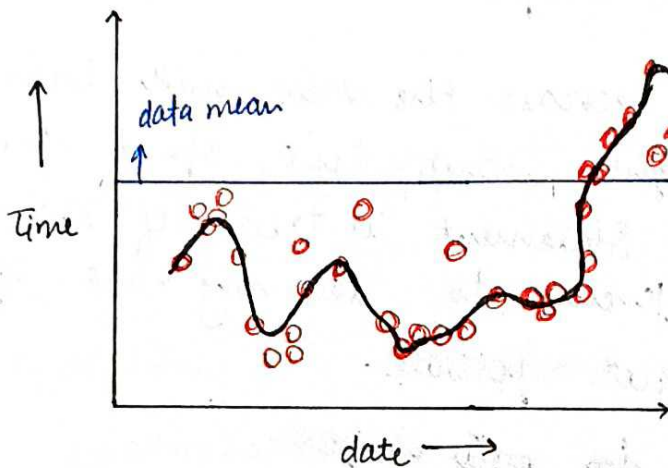
$$W_n = W_0 - \eta \left[ \frac{\partial L}{\partial W} \right] \rightarrow 0$$

old and new weight is same.



# Exponential Weight moving Average (EWMA)

Time Series Data



EWMA

uses

- time series
- financial
- Signal process
- Deep learning
- optimization

①

EWMA →

Day 1 → • → weightage

Day 2 → •

Day 3 → •

Day 4 → •

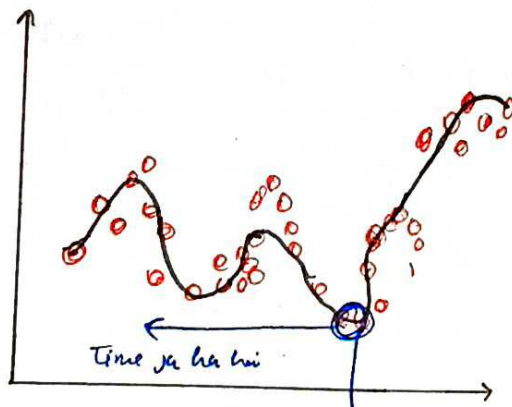
Day 5 → •

Day 3 Weightage > Day 2 Weightage

weightage of Day 5 is greater than Day 4 weightage. (Cz Day 5 weightage bad mai aaya hai)

②

But time ke sath weightage reduce bhi hota jayega. Time age badhta jayega weightage reduce hota jayega.



Weightage reduce ho rha

# Mathematical Formula

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

EWMA  
at particular  
time  $t$ .

pervious time  
of EWMA

time instance per  
data hai.

$$\beta \Rightarrow 0 < \beta < 1$$

ex:  $\rightarrow$

Index	Temp( $\theta$ )
D <sub>1</sub>	25
D <sub>2</sub>	13
D <sub>3</sub>	17
D <sub>4</sub>	31
D <sub>5</sub>	43

Every time stamp we will calculate  
EWMA.

$V_0 = 0$  or  $V_0 = \theta_0$  some test book  $\rightarrow$  assume this one

let  $\beta = 0.9$

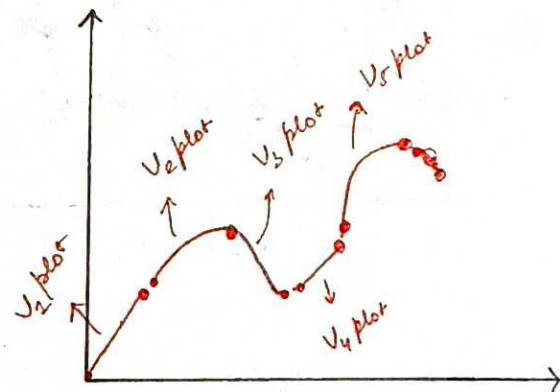
$$V_1 = 0.9 \times V_0 + 0.1 \times 13 \rightarrow \text{②}$$

$$V_1 = 1.3$$

$$V_3 = 0.9 \times 1.3 + 0.1 \times 17 =$$

⋮

$$V_5 = \dots$$



$$\beta \Rightarrow \frac{1}{1-\beta}$$

$$\beta = 0.9$$

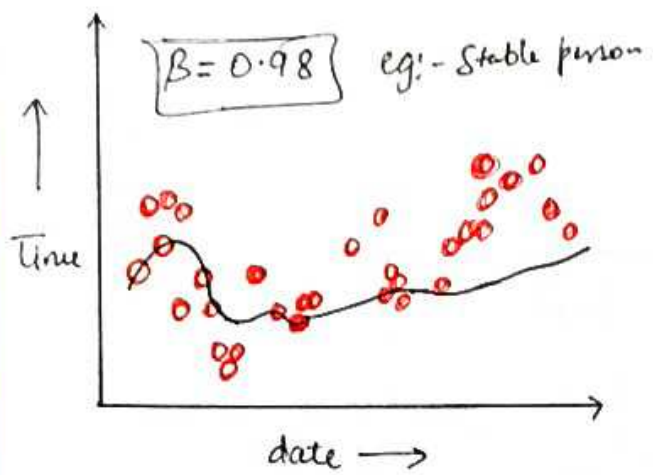
$$\frac{1}{1-0.9} = 10 \rightarrow$$

So,  $\beta$  <sup>act</sup> behave as the  
average of last 10 days data

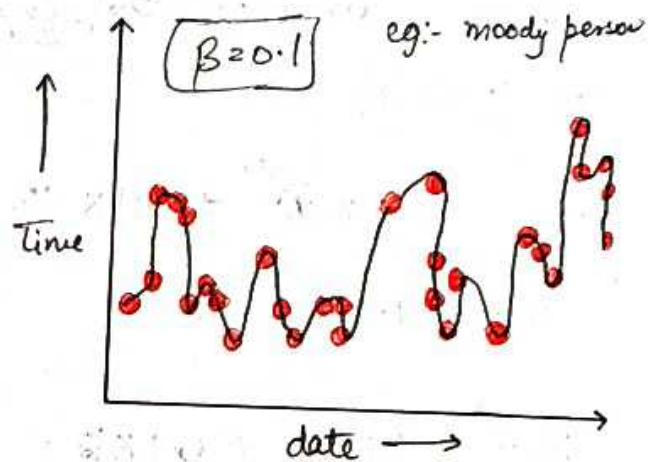
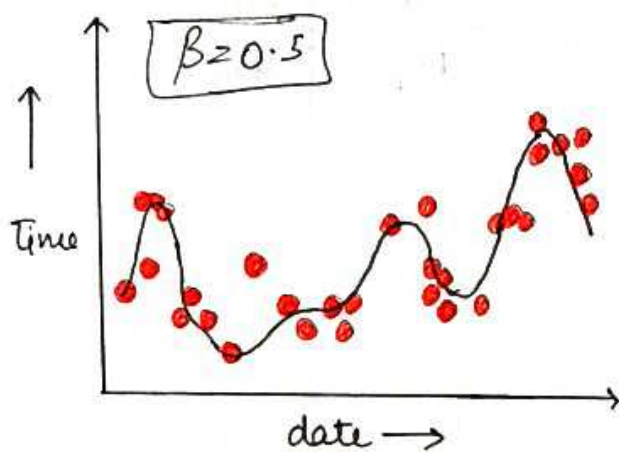
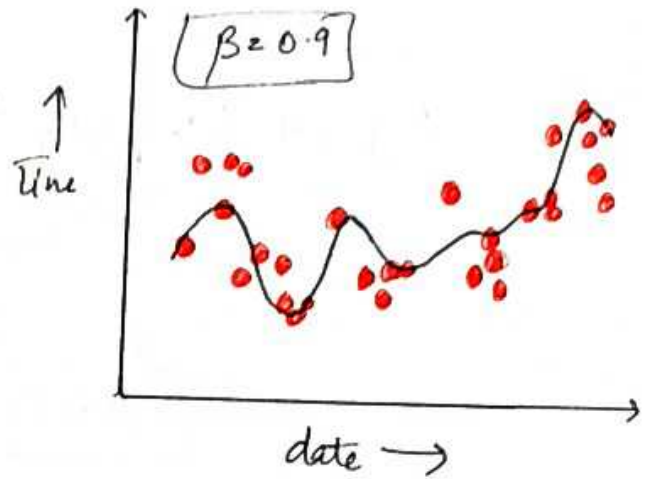
let  $\beta = 0.5$

$$\frac{1}{1-0.5} = 2 \rightarrow$$

$\beta$  <sup>act</sup> behave as the avg of  
last 2 days data



$\beta$  is very high which means past ke points ko weightage de rahi hai.



When  $\beta$  value is very low and graph is very spiky. Because Modelate graph on the base of current point.

Moody person  $\rightarrow$  affect<sup>on</sup> moody on the basis of today's day

Stable person  $\rightarrow$  Person not focus on today. Person focus on past and decide mood.



Sweet spot is  $\rightarrow 0.9$   
 Most of ~~the~~ time you see 0.9

(70)

## Mathematical Intuition

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$$V_0 = 0$$

$$V_1 = (1-\beta) \theta_1$$

$$V_2 = \beta V_1 + (1-\beta) \theta_2$$

$$= \beta(1-\beta) \theta_1 + (1-\beta) \theta_2$$

$$V_3 = \beta V_2 + (1-\beta) \theta_3 \Rightarrow \beta^2(1-\beta) \theta_1 + \beta(1-\beta) \theta_2 + (1-\beta) \theta_3$$

$$V_4 = \beta V_3 + (1-\beta) \theta_4 \Rightarrow \beta^3(1-\beta) \theta_1 + \beta^2(1-\beta) \theta_2 + \beta(1-\beta) \theta_3 + (1-\beta) \theta_4$$

$$= (1-\beta) \left[ \underbrace{\beta^3 \theta_1}_{\substack{\uparrow \\ \text{oldest} \\ \text{term}}} + \underbrace{\beta^2 \theta_2}_{\substack{\uparrow \\ \text{new} \\ \text{the } \theta_1}} + \beta \theta_3 + \underbrace{\theta_4}_{\substack{\uparrow \\ \text{newest} \\ \text{seen}}} \right]$$

$$\beta^3 < \beta^2 < \beta^1 \xrightarrow{\text{lies betw}} 0 < \beta < 1$$

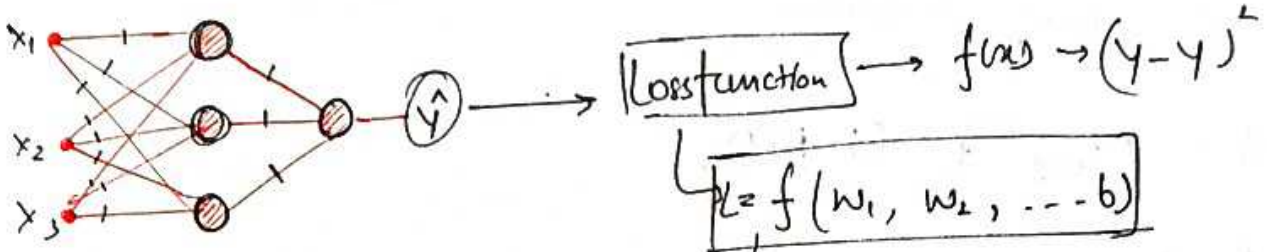
point jitna dikhta hota ja raha hum aur chhoti value se multiply kar rahi hai.

This means that "purani value ko kam weightage de rahi aur new value ko jayada weightage"

[code]  
[EWMA file]

# Understanding Graphs

revise some NN concept



$$y = f(w)$$

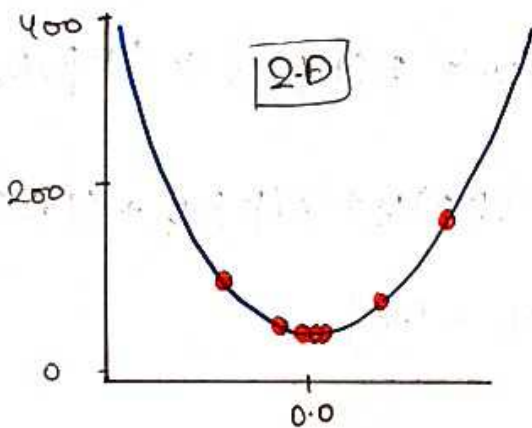
→ x mai  
Change hone pr  
y kitna change  
ho rha.

Let's say weights are 24  
So, graph → 25D  
but human → max 3-D

Let assume, single neuron and then cal.

$$w + 0 \rightarrow \hat{y} \Rightarrow L = f(w) \rightarrow \text{This make 2-D graph}$$

w mai change hone se  
L mai kitna change ho  
rha.

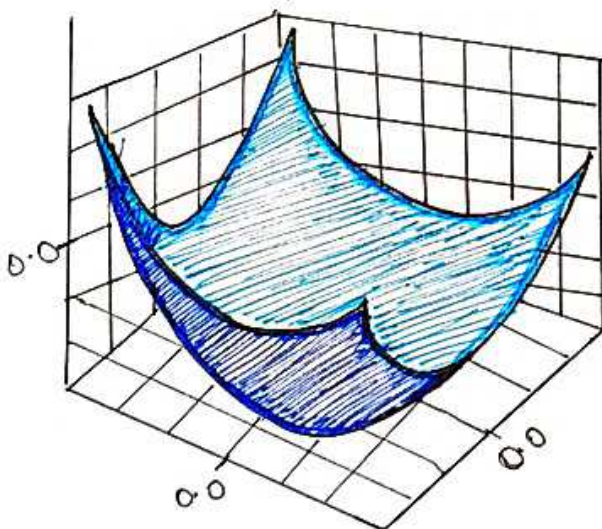


Let assume, single Neuron

$$w + b \rightarrow \hat{y}$$

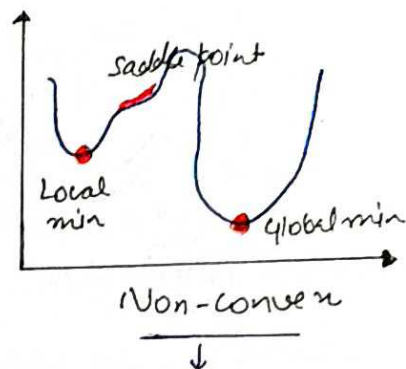
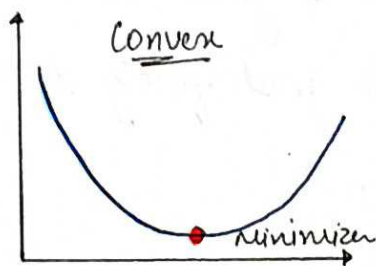
Now we make

3-D Graph





# Convex Vs Non-Convex optimization



1) Local minima

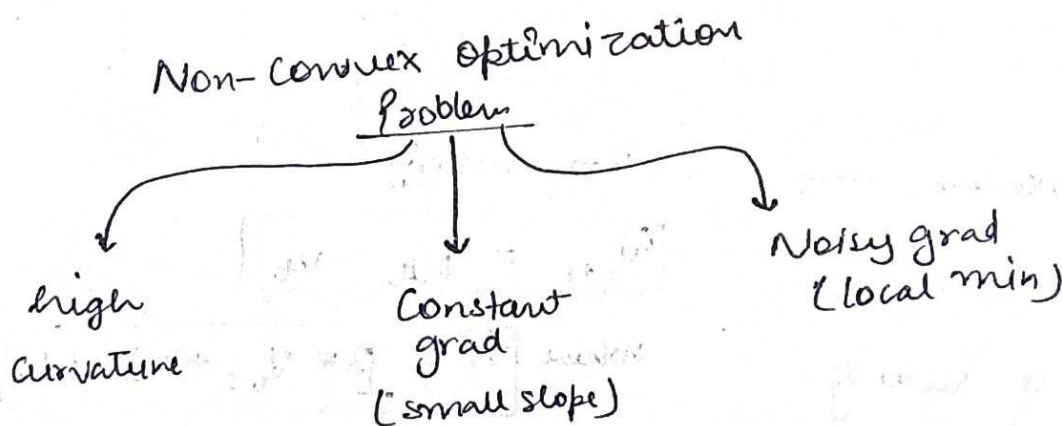
2) Saddle point



slope is small. So, it takes time to reach global min.

3) Higher curvature.

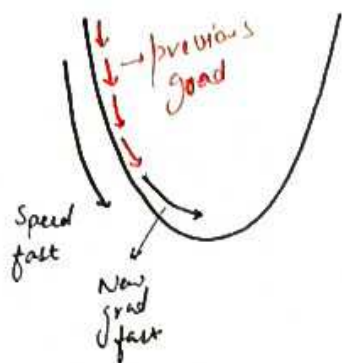
## Momentum Optimization - The why?



→ Momentum can solve all 3 problems

## Momentum

↳



→

If previous 4 grad downward hai to new grad fastly downward jayega.

Same in Physics →



Momentum =  $mv$   
mass      velocity

## Momentum Optimization - Mathematics (The How?)

Normal grad descent algo for update weight →

$$W_{t+1} = W_t - \eta \Delta W_t$$

Momentum →

$V \rightarrow$  velocity

$$W_{t+1} = W_t - V_t$$

history of velocity using for momentum

where  $V_t = \beta * V_{t-1} + \eta \Delta W_t$   
↳  $0 < \beta < 1$



SGD



SGD Momentum

## Effect of beta ( $\beta$ )

$$W_{t+1} = W_t - V_t$$

$$V_t = \beta V_{t-1} + \eta \nabla W_t$$

let say  $\beta = 0$

$$V_t = \eta \nabla W_t$$

$$W_{t+1} = W_t - \eta \nabla W_t$$

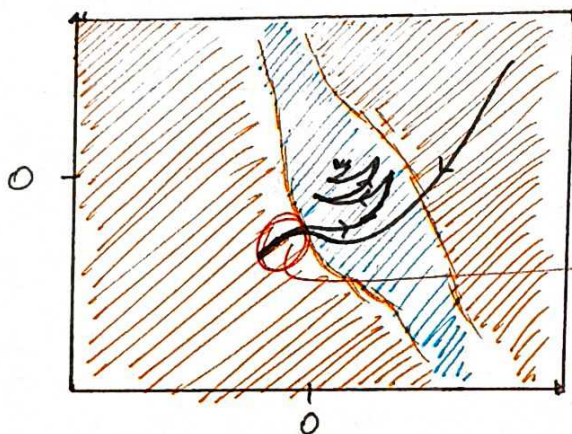
→ Normal SGD

So, if  $\beta = 0$  then momentum act like Normal SGD.

If  $\beta = 0.9 \rightarrow \frac{1}{1-\beta} \rightarrow \frac{1}{1-0.9} = 10$  which means

current velocity = avg of past 10 velocity

## Problem with Momentum Optimization



Time waste in this Oscillation  
That's why momentum is  
not fastest technique



# Nesterov Accelerated Gradient (NAG)

## Mathematical Intuition

$$W_{t+1} = W_t - V_t$$

Where

$$V_t = \beta V_{t-1} + \eta \nabla W_t$$

$\beta$  = decay factor

$$W_{t+1} = W_t - (\beta V_{t-1} + \eta \nabla W_t)$$

for new weight depend  $\rightarrow$  Past velocity + gradient at that point

But in Simple GD

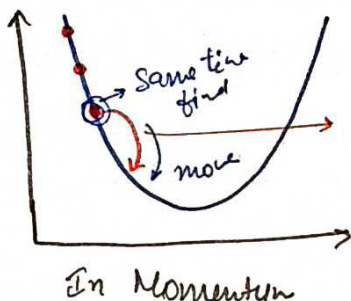
$$W_{t+1} = W_t - \eta \nabla W_t$$

New weight depend  $\rightarrow$  gradient at that point

## Nesterov Accelerated Gradient (NAG)

$$\hookrightarrow W_{t+1} = W_t - (\beta V_{t-1} + \eta \nabla W_t)$$

ye kis position pe calculate hoga (NAG)

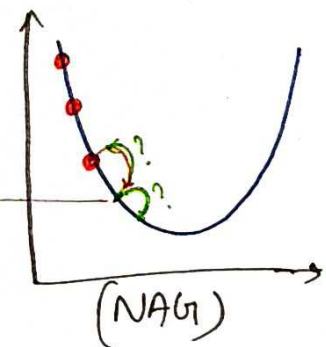


Calculate  
History of velocity + gradient at that point

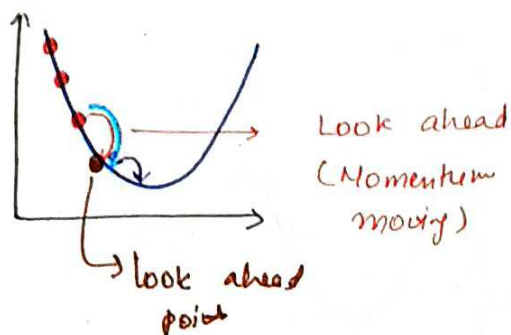
find both  $\rightarrow$  Same time then move.

Where the Momentum move from that point, we calculate "gradient at the point". And then again move.

At this point we first calculate Momentum and move.



After gradient point  $\rightarrow$  decide to move upward? or downward?



formula  $\rightarrow W_{la} = W_t - \beta V_{t-1} - 0$

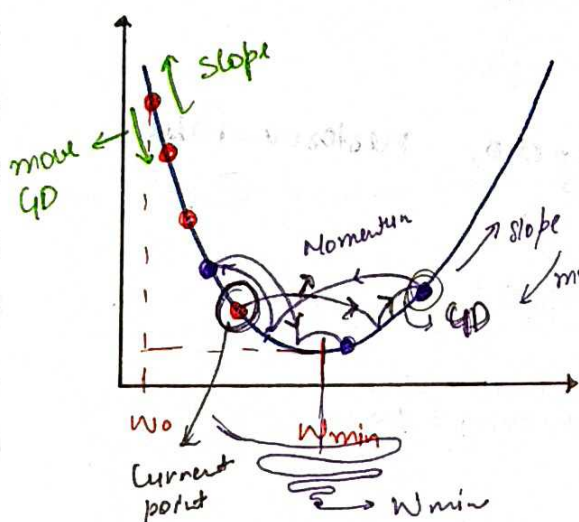
$$V_t = \beta V_{t-1} + \eta \nabla W_{la} \quad (2)$$

eq ① put in eq ②

$$V_t = (W_t - W_{la}) + \eta \nabla W_{la}$$

$$W_{t+1} = W_t - V_t$$

### Geometric Intuition

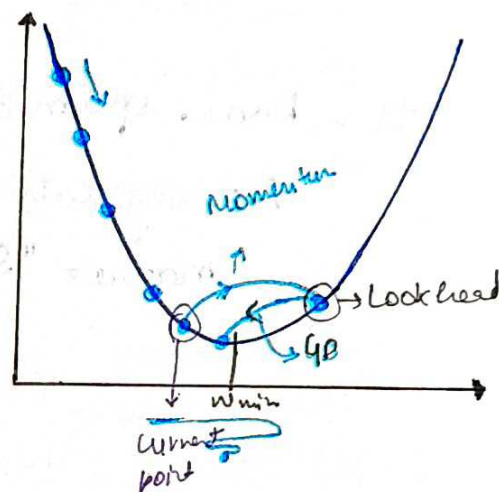


Slope ke opposite gradient descent move krta

first move kiya  $\rightarrow$  Momentum  
cause of previous point

Gradient Descent move ~~from~~  
cause of ~~previous~~ <sup>current</sup> point not  
momentum point.

finding Momentum and  
Gradient descent  
at a same time.



Here's we find  
first Momentum and  
move the line then  
find Gradient descent  
and again move the  
line.

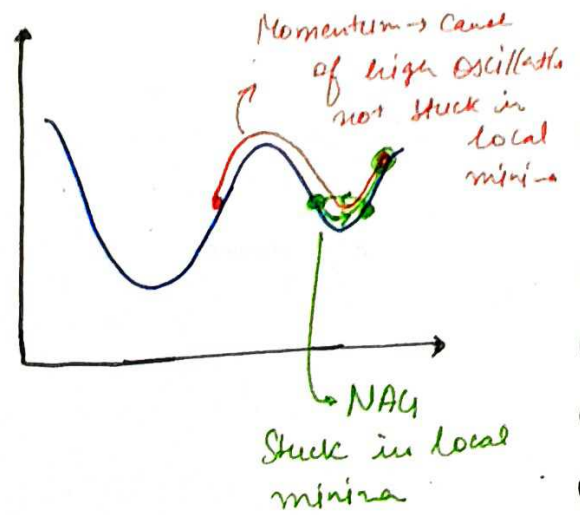
Gradient descent on the  
basis of look ahead point

but the Momentum  
Gradient descent find on  
the basis of ~~previous~~ <sup>current</sup> point

## Disadvantage

NAG  $\rightarrow$  Oscillation  
 $\downarrow$   
dampen

May be stuck in local minima



## Keras Code $\rightarrow$ Momentum

```
tf.keras.optimizers.SGD (
```

```
    learning_rate = 0.01, Momentum = 0.0, nesterov = False,
```

```
    name = "SGD", **kwargs
```

```
)
```

Simple SGD  $\rightarrow$  Momentum = 0 , Nesterov = False

Momentum  $\rightarrow$  Momentum = 0.9 , Nesterov = False

NAG  $\rightarrow$  Momentum = 0.9 , nesterov = True



# AdaGrad (Adaptive Gradient)

74

AdaGrad  $\rightarrow$  learning rate is not fix

change or adapt according to the situation.

where we use:

$\rightarrow$  Input features  $\rightarrow$  Scale is different

eg:- 

cgpa	Salary	iq
1-10	0-100k	0-100k

 $\rightarrow$  Scales are different

$\uparrow$  range

$\rightarrow$  In this Normalize or Scaling can be use.

$\rightarrow$  Input feature  $\rightarrow$  Sparse  $\rightarrow$  most of values are zero.

eg:- 

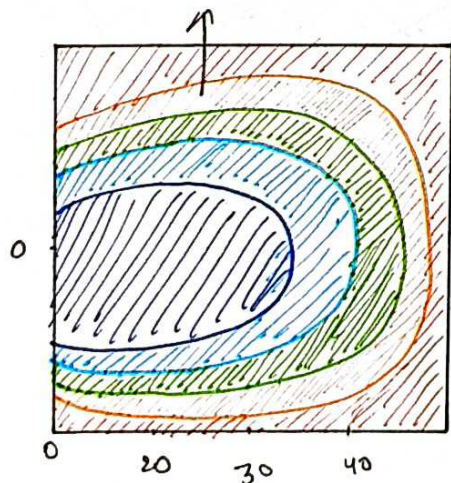
iq	cgpa	from lit	package
1	1	0	1

 $\rightarrow$  if most of the value is 0.

## Problem

Elongated bowl problem

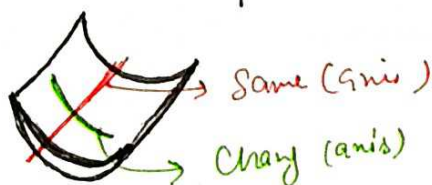
Sparse



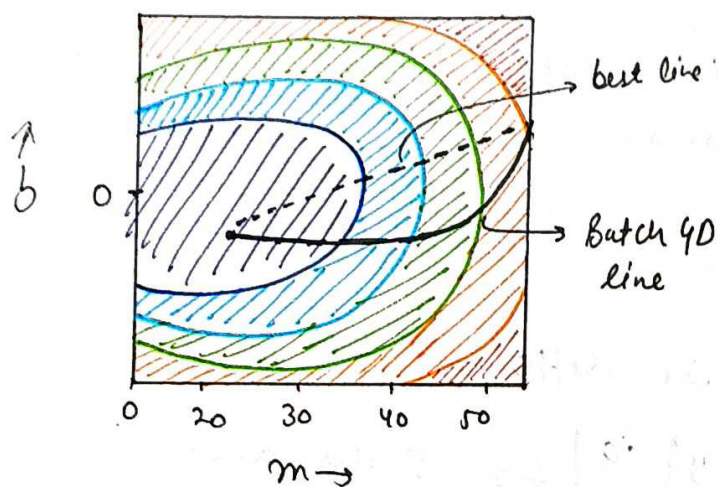
Half Contour plot

Cause of maximum zeroes in data Half Contour plot made.

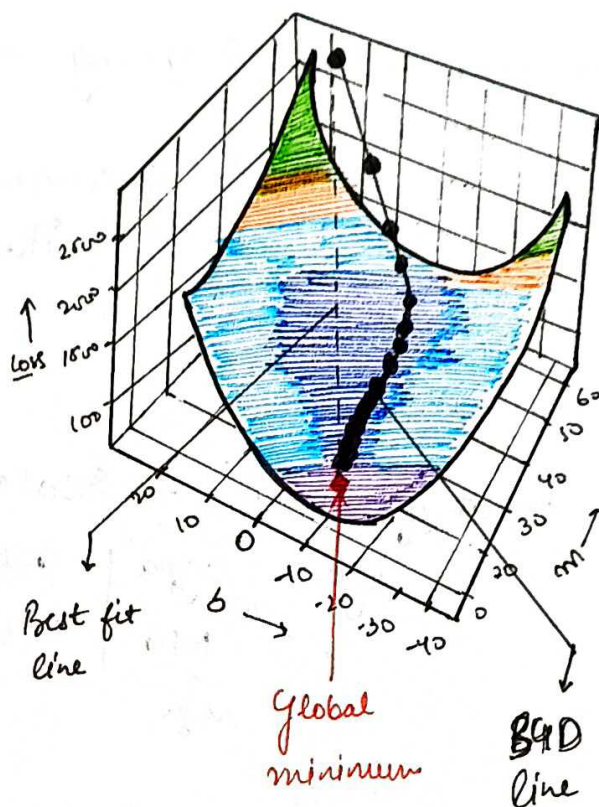
Elongate  $\rightarrow$  slope change along 1-axis  
slope same at 2nd axis



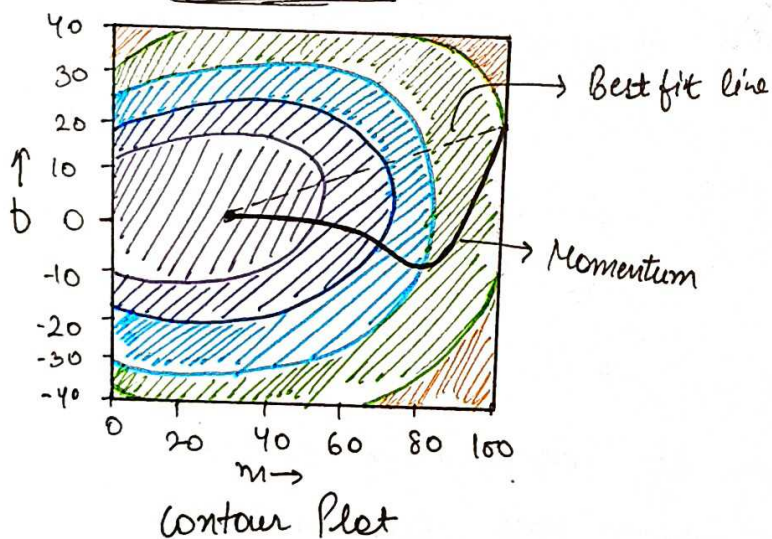
Batch Gradient Descent and Momentum is not good for sparse data.



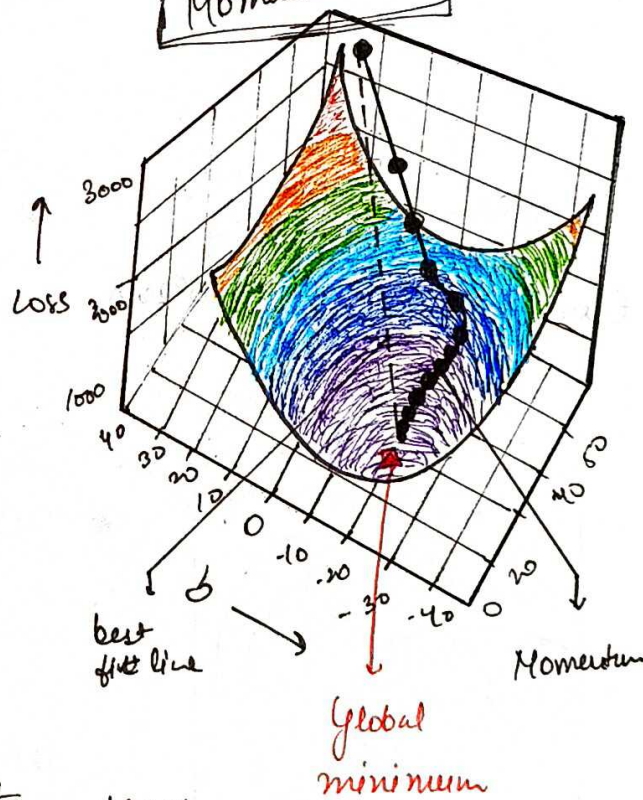
Cause of curve  $\rightarrow$  Batch Gradient Descent take more time to reach global minima.



Momentum



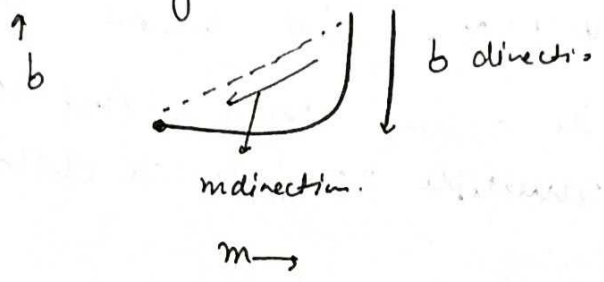
Momentum 3-D



Cause of overshoot  $\rightarrow$  Momentum take more time to reach global minimum.



Why we moving toward b direction not m?



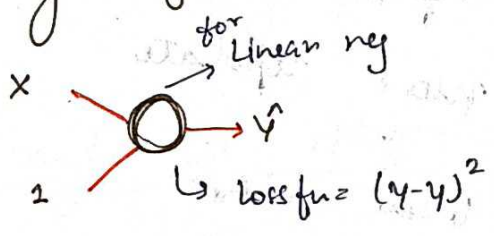
because sparse data present in m  
 correlated with b  
 m ← correlated

Let assume b col input

	X	Y
1	0.72	-
1	0.81	-
1	0	-
1	0	-
1	0	-
1	0	-
1	0.2	-

Sparse value present in ~~m~~ and X correlated with m. That's why line not move to m.

Explaining Why?



BGD → Sparse  
 pseudocode?

for i in epochs:

$$W = W - \eta \frac{\partial L}{\partial W}$$

$$b = b - \eta \frac{\partial L}{\partial b}$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial W}$$

$$\frac{\partial L}{\partial W} = -2(y - \hat{y})x$$

$$\frac{\partial L}{\partial b} = -2(y - \hat{y}) \text{ (1)}$$

∵ input is 1

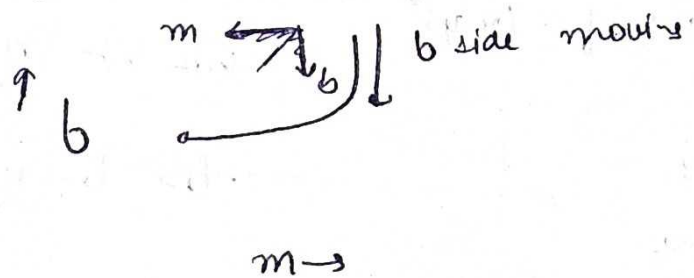
Let assume 100 rows  
 so, we calculate 100 times  $\frac{\partial L}{\partial W}$  and add all 100 repeated weights  $\frac{\partial L}{\partial W}$  term. we will get big number for  $\eta \left( \frac{\partial L}{\partial W} \right)$



→ but if lots of 0 value present in  $x$  and when we calculate  $\frac{\partial L}{\partial w} = -2(y - \hat{y})x$  and whole term will be 0 because  $x$  is 0 in data. And there are multiple zero  $\frac{\partial L}{\partial w}$   $\hookrightarrow$  multiple zero in  $x$  data. After

that we add whole  $\frac{\partial L}{\partial w}$  values  $\rightarrow$  Small values cause of zero present in data. Now update also small if  $\frac{\partial L}{\partial w}$  is small value. Small update in every epoch

→ In  $b$ ,  $\frac{\partial L}{\partial b} = -2(y - \hat{y}) * 1$  we will get big values because whole term multiply with 1 not  $x$ . So,  $\frac{\partial L}{\partial b}$  never be 0. So after addition, value will be bigger and <sup>see</sup> more update



## Adagrad Mathematics + Intuition

(76)

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

$$b = b - \eta \frac{\partial L}{\partial b}$$

→ Check whether  $w$  get higher update or  $b$  get higher update.

→ According to previous discussion  $b$  are getting higher updates. And  $w$  getting lower updates.

So maintain the both update → will reduce learning rate of  $b$  (higher updates).

low learning rate  $\times$  higher update =  $\overset{\text{low update}}{\uparrow}$  low value

Through this we can match with  $w$  update.

⇒ "for every parameter learning rate is different"

### Formula for Adagrad

$$V_t = V_{t-1} + (\nabla w_t)^2$$

$$w_{t+1} = w_t - \frac{\eta \nabla w_t}{\sqrt{V_t + \epsilon}}$$

$$b_{t+1} = b_t - \frac{\eta \nabla b_t}{\sqrt{V_t + \epsilon}}$$

if  $V_t = 0$  then whole not to be 0.

Normal Grad Descent

$$w_{t+1} = w_t - \eta \nabla w_t$$

Let's discuss about  $V_t = \underbrace{V_{t-1}}_{\text{Past } V_t} + \underbrace{(\nabla W_t)^2}_{\text{gradient ka square}}$

$\left(\frac{\partial L}{\partial w}\right)^2$  or  $\left(\frac{\partial L}{\partial b}\right)^2$

→ Why gradient square?

Because gradient may be negative or positive.

That's why we squared gradient for negative to positive

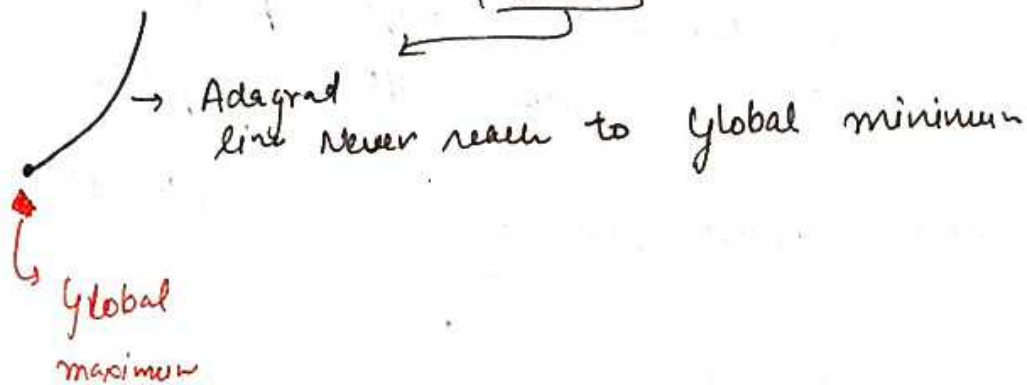
→ Why not use mod?

because differentiation nhi hoga mod ka.

### Disadvantage

(Adagrad → not use in complex NN)  
→ use in linear regression

Reason



→ Why Adagrad never reach global minimum?

Because learning rate divide with past updates

$$\frac{\eta}{\sqrt{V_t}} \rightarrow \text{past updates}$$

After sometime or many epochs  $V_t$  getting bigger and learning rate smaller. So, after divide learning rate become more smaller and no more updates. So never reach to global minimum.



# RMSProp

(72)

RMSProp → Root Mean Square Propagation

Improvement

Adagrad ← Sparse data

Disadvantage → reduce learning rate



## Mathematical Formulation

### Adagrad

$$V_t = V_{t-1} + (\nabla w_t)^2$$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{V_{t+1}}} \nabla w_t$$

past gradient → giving high value

$w_t$  multiply with smaller  
So, whole term will be more  
smaller and not getting any  
updates

because of past gradient this term will be  
bigger.

Learning rate divide with bigger term  
so value is smaller

### RMSProp

Old value ko  
kann value de  
sake

expo decay any

Recent value ko  
jaye value de sake

$$V_t = \beta V_{t-1} + (1-\beta) (\nabla w_t)^2$$

Where  $\beta = 0.95$   
↳ Hyperparam

any time change

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{V_{t+1}}} \nabla w_t$$

so,  $V_t$  will not be much bigger

$$V_0 = 0$$

$$V_1 = 0.95 \times 0 + 0.05 (\nabla W_1)^2$$

$$V_2 = 0.95 \times 0.05 (\nabla W_1)^2 + 0.05 (\nabla W_2)^2$$

$$V_3 = \underbrace{0.95 \times 0.95 \times 0.5 (\nabla W_1)^2}_{\substack{\uparrow \text{Smallest} \\ \text{epoch 1} \\ \downarrow \text{old}}} + \underbrace{0.95 \times 0.05 (\nabla W_2)^2}_{\text{epoch 2}} + \underbrace{0.5 (\nabla W_3)^2}_{\substack{\uparrow \text{epoch 3} \\ \text{recent}}} \quad \nearrow \text{Largest}$$

→ RMSprop →  $V_t$  shoot →  $\infty$

$\eta \rightarrow$  very small  $\times$

recent → accordingly move → updates are relatively  
global minimum → reach



Adagrad → good with → convex / optimization problem as regression

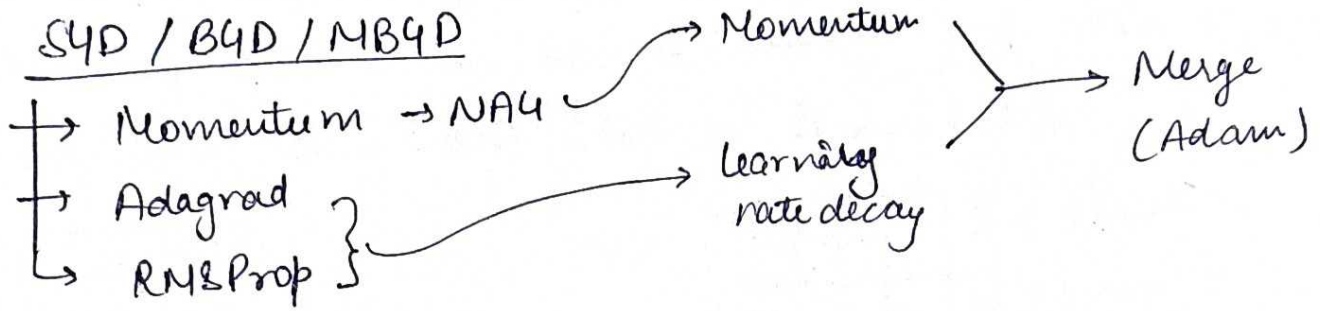
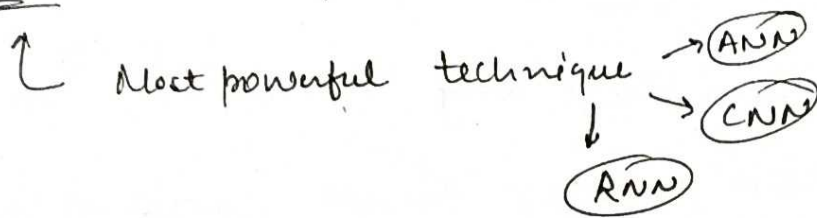
↳ bad with → non-convex optimization and complex NN

In RMSprop → No disadvantage

# ADAM Optimizer

78

Adam → Adaptive Moment Estimation



## Mathematical Formulation

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{V_t + \epsilon}} \times m_t$$

where,

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla W_t \rightarrow \text{Momentum}$$

$$V_t = \beta_2 V_{t-1} + (1 - \beta_2) (\nabla W_t)^2 \rightarrow \text{Adagrad}$$

Bias Correction → find bias correction after  $m_t$  and  $V_t$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

epoch Number

$\beta_1 = 0.9$

$$\hat{V}_t = \frac{V_t}{1 - \beta_2^t}$$

epoch Number

$\beta_2 = 0.99$  according to keras



Why bias correction?

→ Because <sup>during</sup> starting →  $m_t$  start from 0 and  $V_t$  also start from 0. So, biased can be in the  $m_t$  and  $V_t$ . That's why we find bias correction.