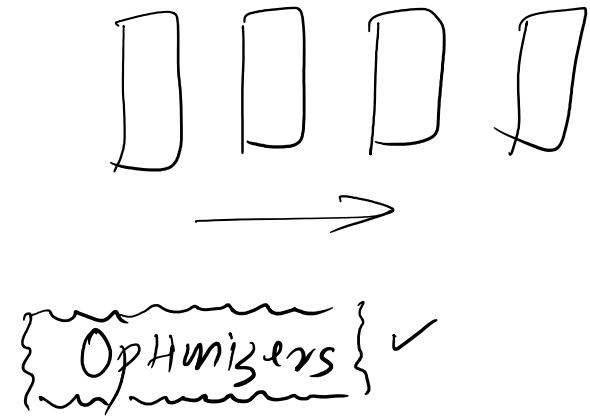


Introduction

05 July 2022 09:57

Performance of NN

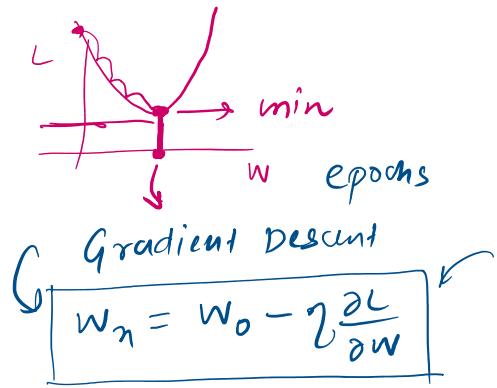
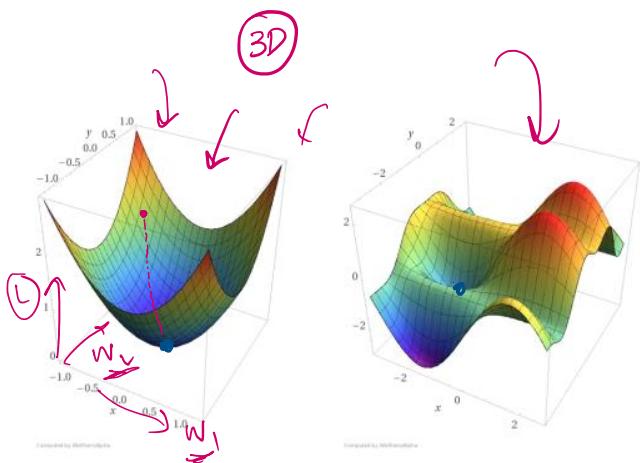
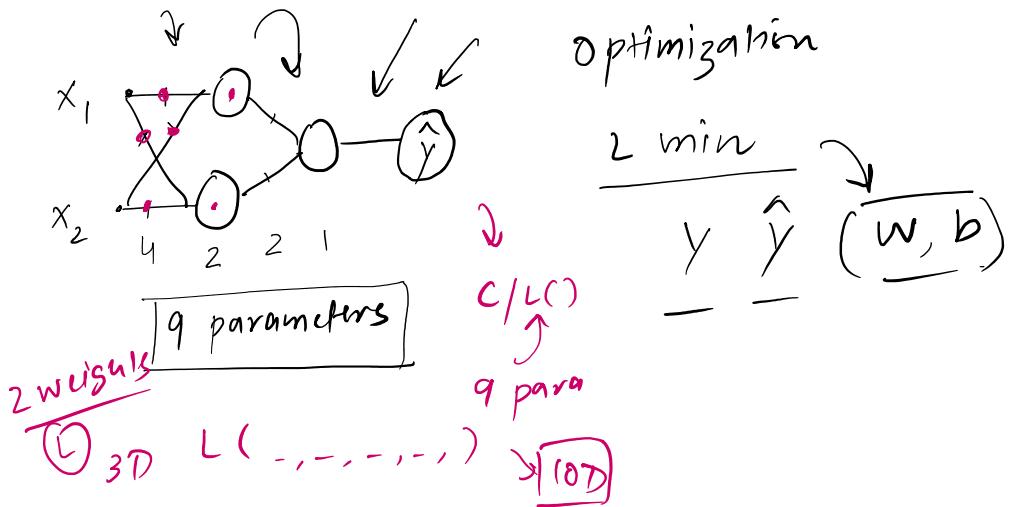
↳ how the speed the training



- Weigh init
- Batch Norm
- Activation function

Role of Optimizer

05 July 2022 10:01



Types of Optimizers

05 July 2022 10:02

- Batch GD
- Stochastic GD
- Mini batch GD

10

epochs = 10

$$w_n = w_0 - \eta \frac{\partial L}{\partial w}$$

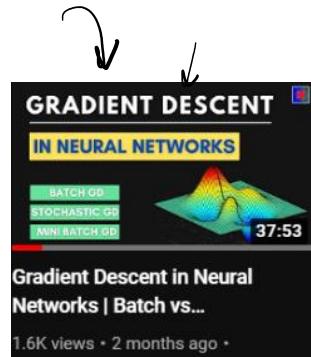
500 rows → pred

↓
loss → weight update

$[10 \times 500] \rightarrow 5000$

batch size = 100

5 batch → 10×5 times

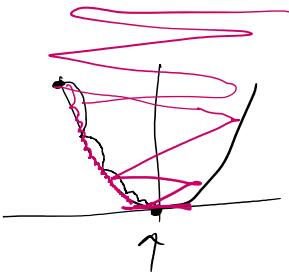


Challenges

05 July 2022 10:02

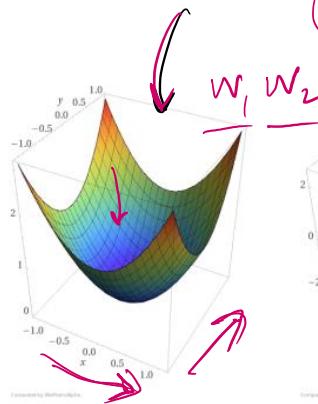
1) learning rate

$$w_n = w_0 - \eta \frac{\partial L}{\partial w_0}$$

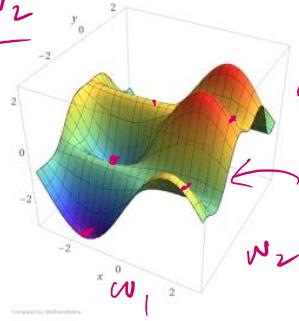


2) learning rate scheduling → pre define

3)



(η) same both the directions



⑨ weight's and bias' ⑨ sep learn

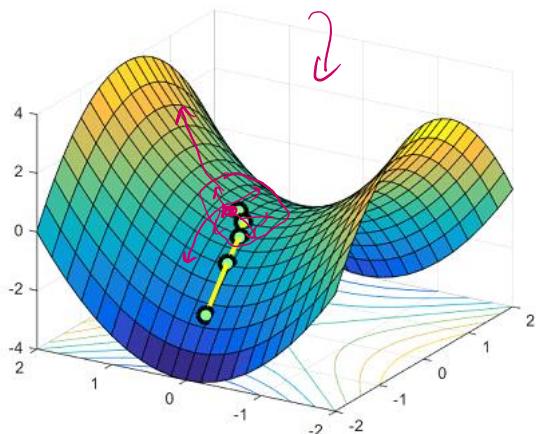
10-D ~ 10 dimensions

minima

mm ↗
sub optimum

4) local minima

5) saddle point



$$\frac{\partial L}{\partial w} = 0$$

$$w_n = w_0 - \eta \left[\frac{\partial L}{\partial w} \right]$$

What next?

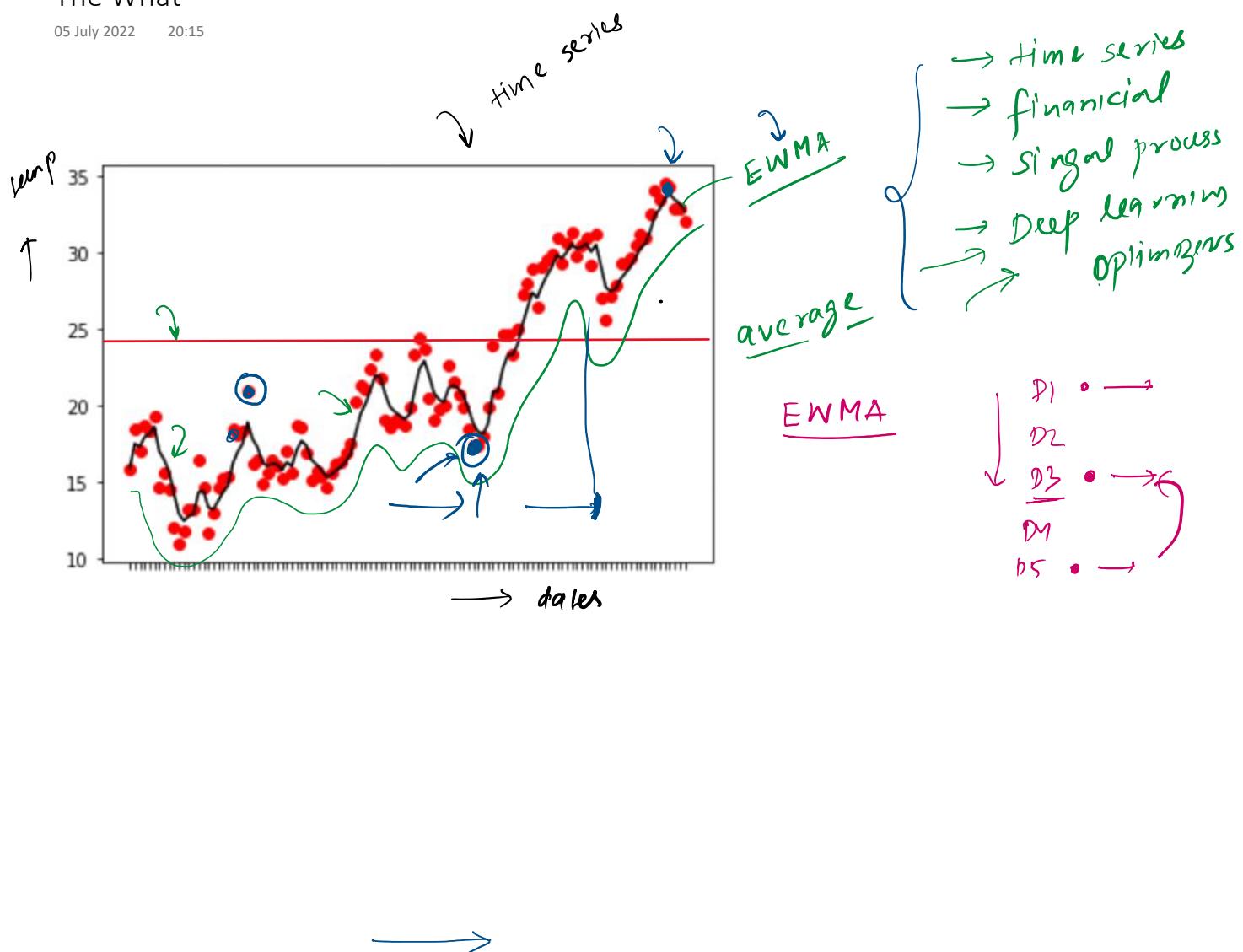
05 July 2022 10:02

- 1) Momentum
- 2) Adagrad
- 3) NAG
- 4) RMSprop
- 5) Adam

→ Exponentially
Weighted
Moving
Average

The What

05 July 2022 20:15



Mathematical Formulation

05 July 2022 20:16

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

$\beta \rightarrow 0 < \beta < 1$

$$V_0 = 0$$

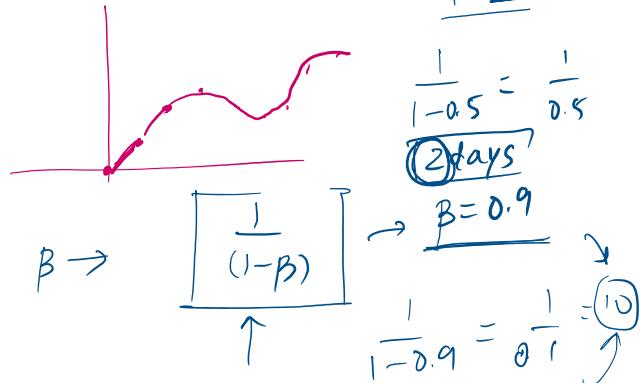
$$V_0 = \theta_0$$

$$\begin{aligned} V_1 &= 0.9 \times V_0 + 0.1 \times 13 \\ &= 1.3 \end{aligned}$$

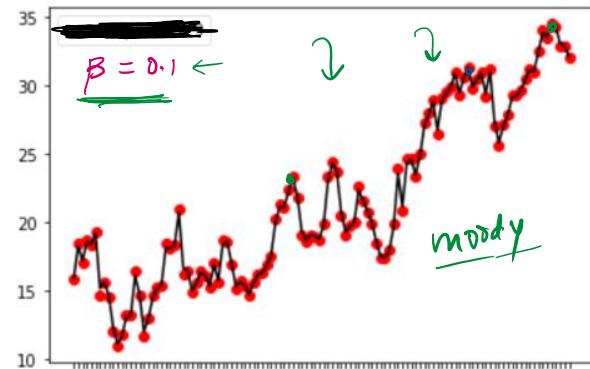
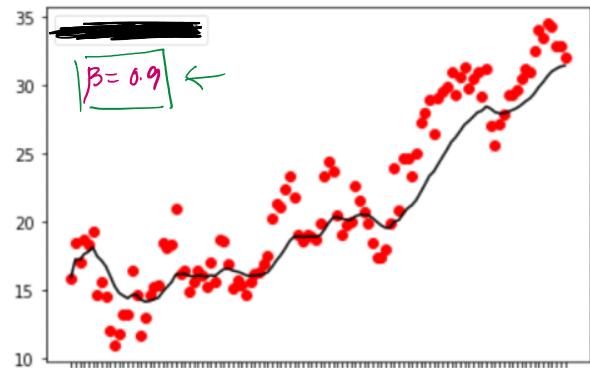
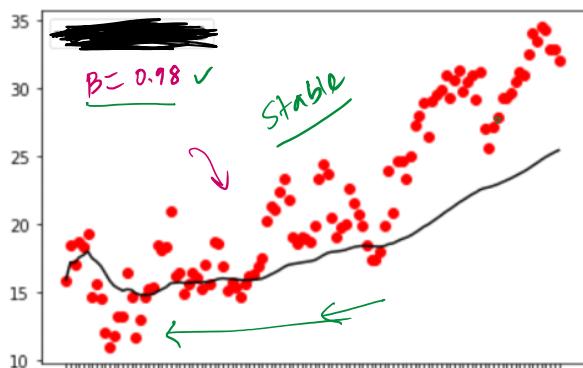
$$V_2 = 0.9 \times 1.3 + 0.1 \times 17 =$$

Index	$temp(\theta)$
D1	25
D2	13
D3	17
D4	31
D5	43

$$\beta = 0.5$$



Effect of β



Mathematical Intuition

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

$$\left. \begin{aligned} V_4 &= \beta V_3 + (1-\beta) \theta_4 \\ V_4 &= \beta^3 (1-\beta) \theta_1 + \beta^2 (1-\beta) \theta_2 + \beta (1-\beta) \theta_3 + (1-\beta) \theta_4 \end{aligned} \right\}$$

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

$$V_0 = 0$$

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

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

$$= \underline{\beta (1-\beta) \theta_1} + \underline{(1-\beta) \theta_2}$$

$$V_3 = \underline{\beta V_2} + \underline{(1-\beta) \theta_3}$$

$$= \underline{\beta^2 (1-\beta) \theta_1} + \underline{\beta (1-\beta) \theta_2} + \underline{(1-\beta) \theta_3}$$

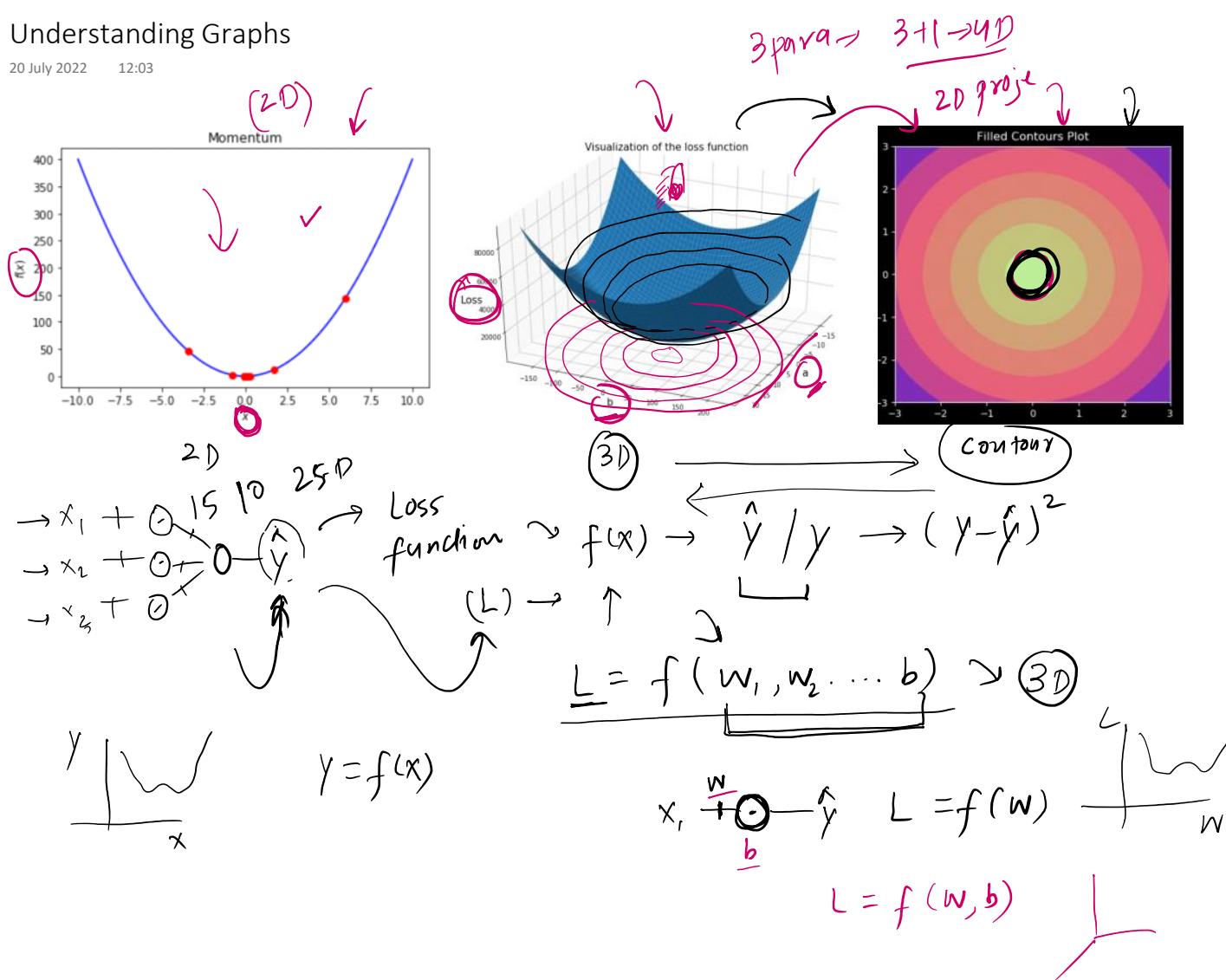
$$\left. \begin{aligned} V_4 &= \beta^3 (1-\beta) \theta_1 + \beta^2 (1-\beta) \theta_2 + \beta (1-\beta) \theta_3 + (1-\beta) \theta_4 \\ &= (1-\beta) [\beta^3 \theta_1 + \beta^2 \theta_2 + \beta \theta_3 + \theta_4] \end{aligned} \right\}$$

$\beta^3 < \beta^2 < \beta < 1$

—————>

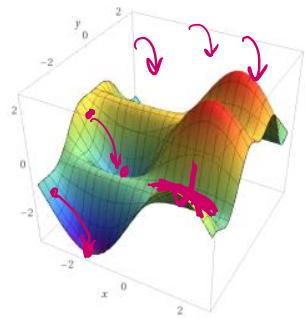
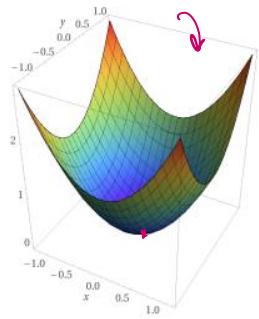
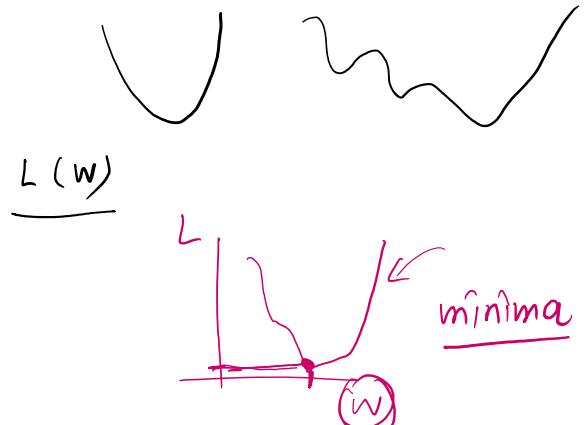
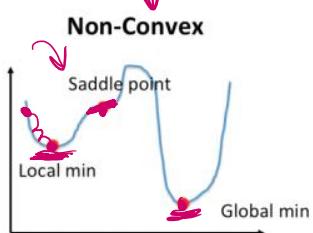
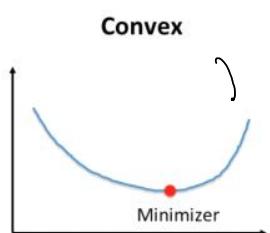
Understanding Graphs

20 July 2022 12:03

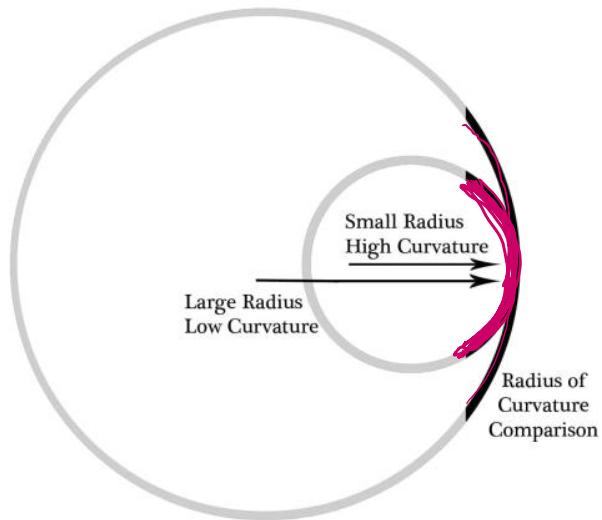


Convex Vs Non-Convex Optimization

20 July 2022 13:06

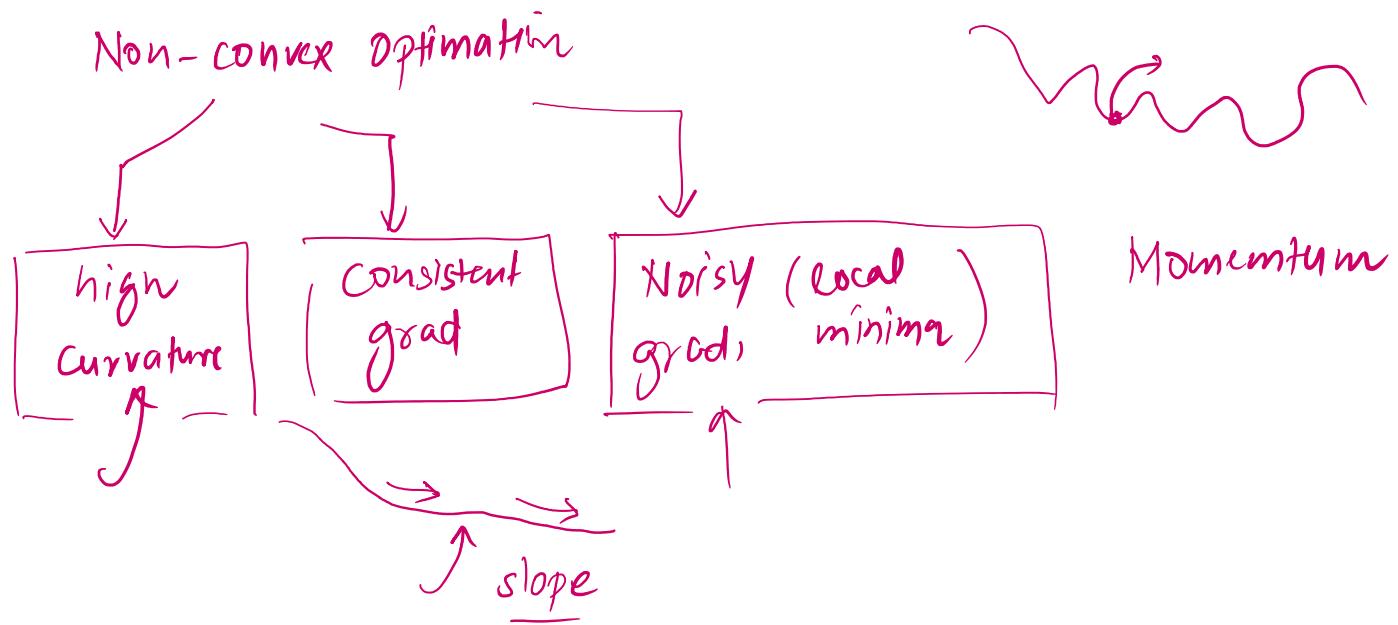


- 1) local minimum
- 2) saddle point
- 3) High curvature



Momentum Optimization - The Why?

20 July 2022 14:28

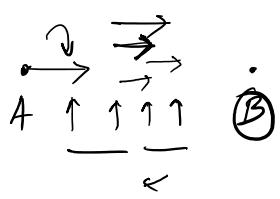


Momentum Optimization - The What?

20 July 2022 14:33

Speed

Common



prev grad \rightarrow

$$\boxed{mV}$$

$m \rightarrow$ unit mass
velocity \rightarrow previous history update



Momentum Optimization - Mathematics(The How?)

20 July 2022 14:56

$$\underline{w_{t+1}} = \overbrace{\underline{w_t}}^{\text{acceleration}} - \boxed{\eta \nabla w_t}$$

$\} \quad v \rightarrow \text{velocity}$

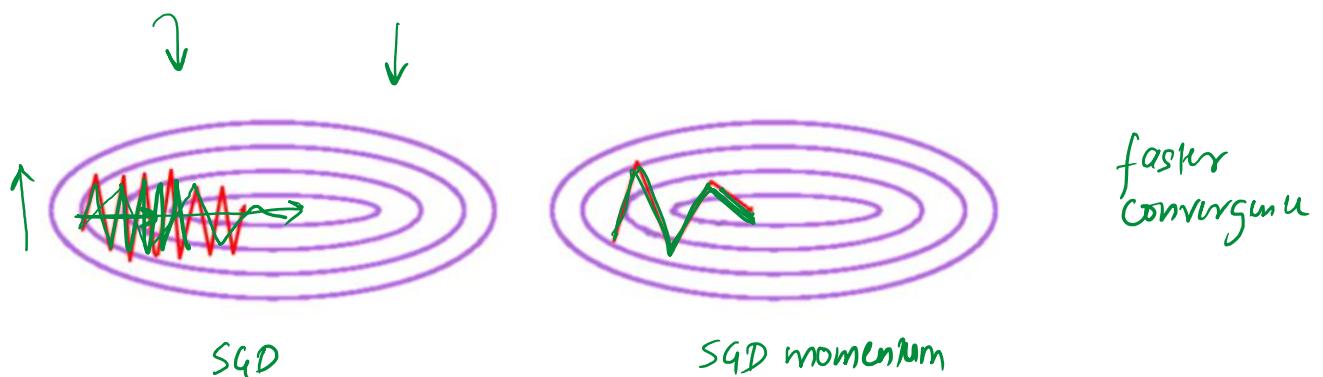
$$w_{t+1} = w_t - \underline{v_t} \quad \rightarrow$$

$$v_t = \boxed{\beta * v_{t-1}} + \boxed{\eta \nabla w_t}$$

$\boxed{0 < \beta < 1}$ $\boxed{v_{t-2} \quad v_{t-3}}$

momentum $\text{history of post velo-} \quad \boxed{\text{use}}$

prev velocity



Effect of beta

20 July 2022 15:09

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

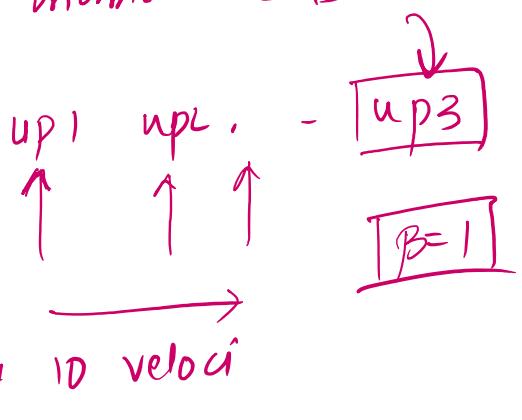
$$w_{t+1} = w_t - v_t$$

SGD

$$\underline{v_t} = \underline{\beta v_{t-1}} + \underline{\eta \nabla w_t}$$

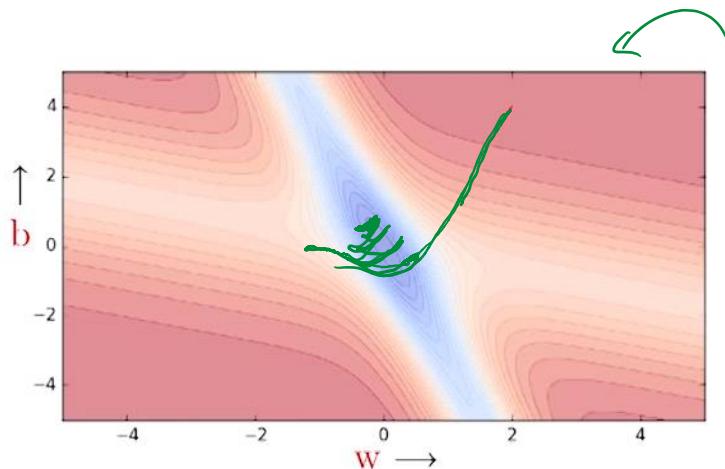
$\beta = 0$ momen \rightarrow SGD

$$\frac{1}{1-\beta} = \frac{1=10}{0.1} \frac{1}{1-\beta} \quad \begin{matrix} \text{decaying factor} \\ \text{current} \\ \text{avg past} \end{matrix}$$

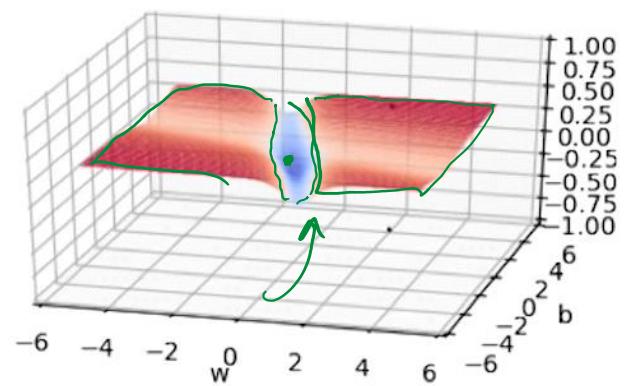


Problem with momentum optimization

20 July 2022 15:16

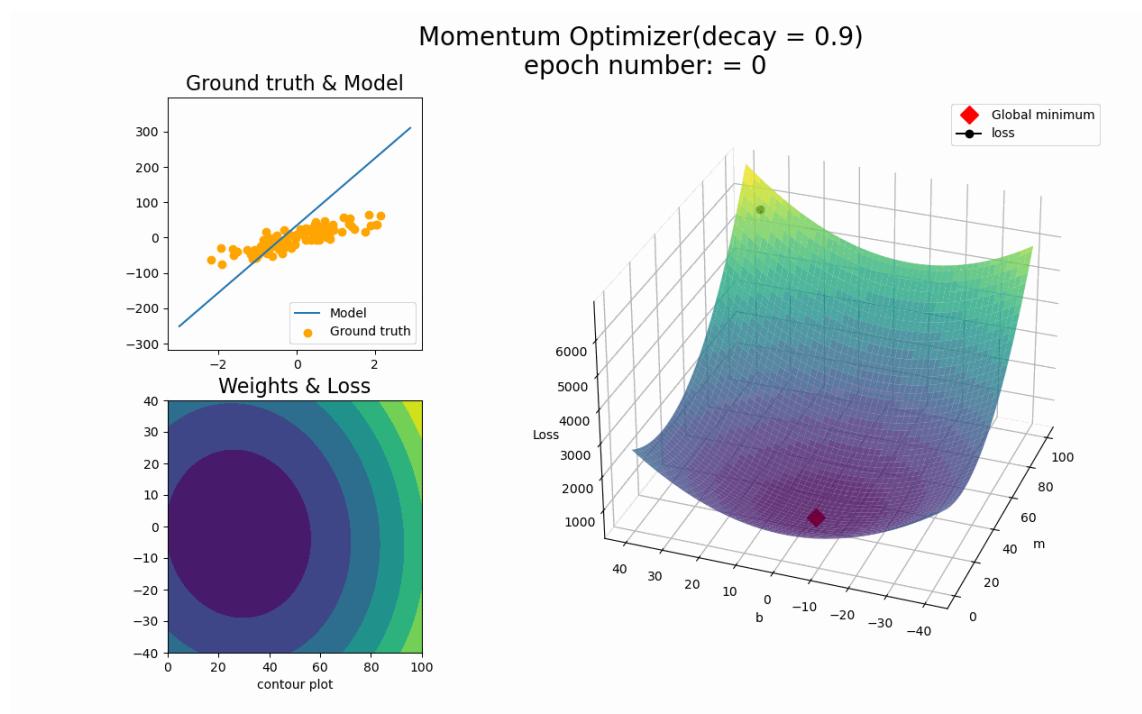


washle



decay 0.9

Momentum Optimizer(decay = 0.9)
epoch number: = 0



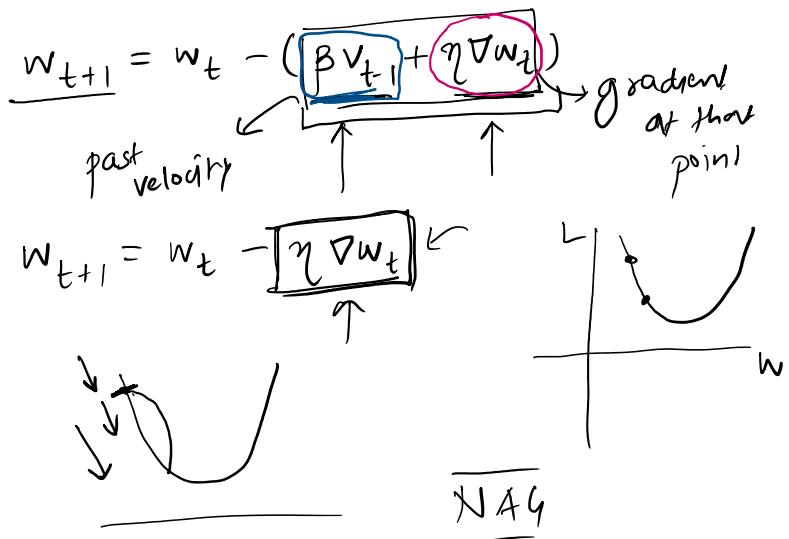
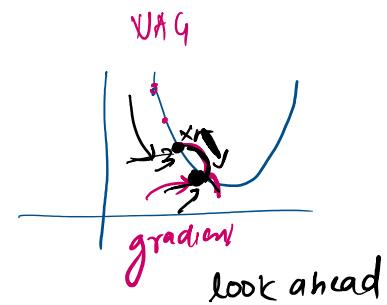
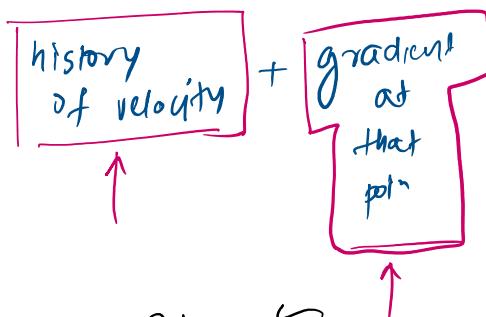
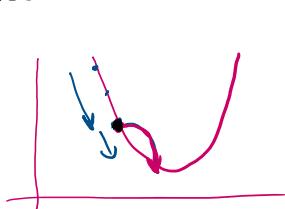
Momentum

$$w_{t+1} = w_t - v_t$$

where

$$v_t = \beta v_{t-1} + \eta \nabla w_t$$

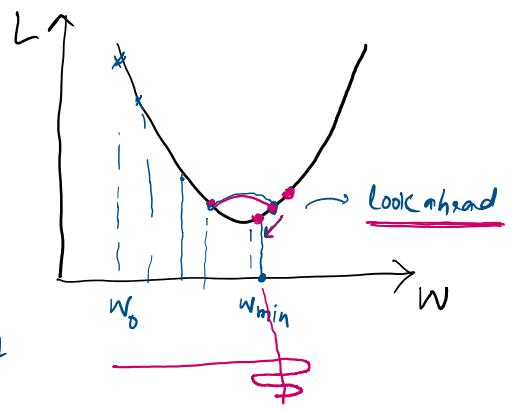
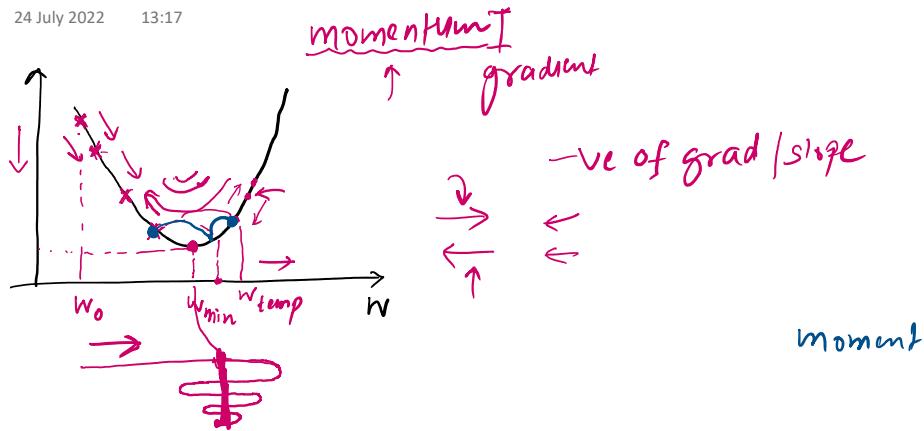
$\beta \rightarrow$ decay factor

Nesterov Accelerated Gradient *better*

$$\left\{ \begin{array}{l} w_{\text{la}} = w_t - \beta v_{t-1} \\ v_t = \beta v_{t-1} + \eta \nabla w_{\text{la}} \\ w_{t+1} = w_t - v_t \end{array} \right. \quad \begin{aligned} v_t &= (w_t - w_{\text{la}}) + \eta \nabla w_{\text{la}} \end{aligned}$$

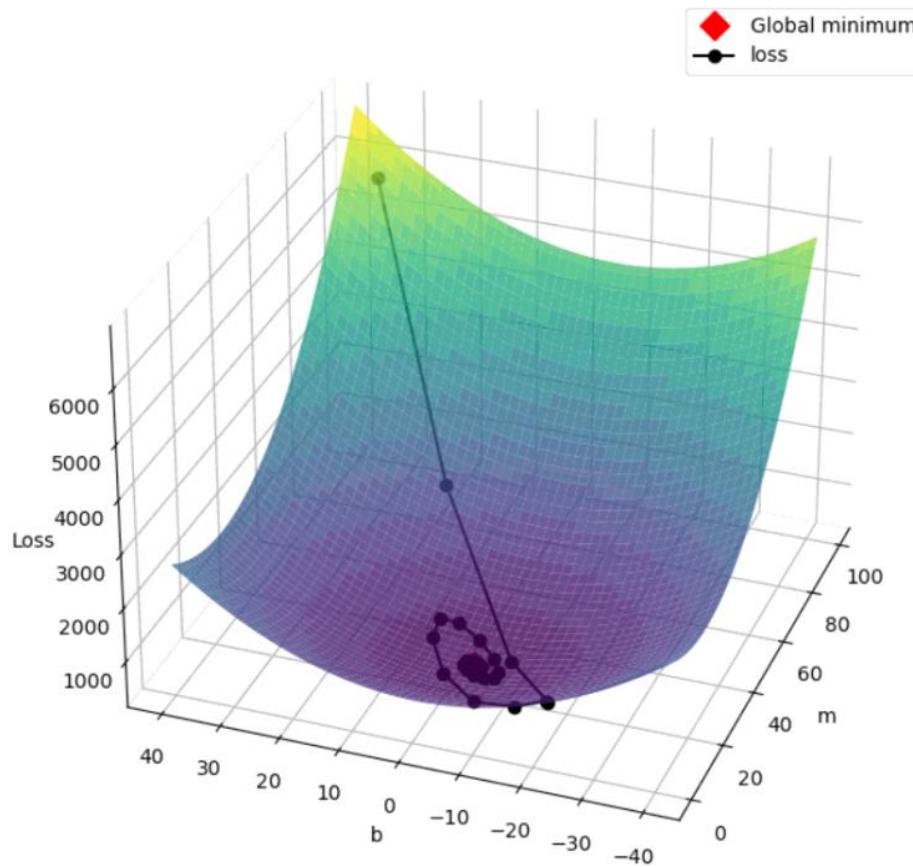
Geometric Intuition

24 July 2022 13:17



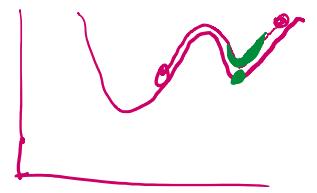
Disadvantage

24 July 2022 13:18



Disadvantage

NAG \rightsquigarrow oscillation
 \downarrow
dampen
local minima



Keras Code

24 July 2022 13:18



```
tf.keras.optimizers.SGD(  
    learning_rate=0.01, momentum=0.0, nesterov=False, name="SGD", **kwargs  
)
```

SGD



Momentum

momentum = 0.9

nesterov = False

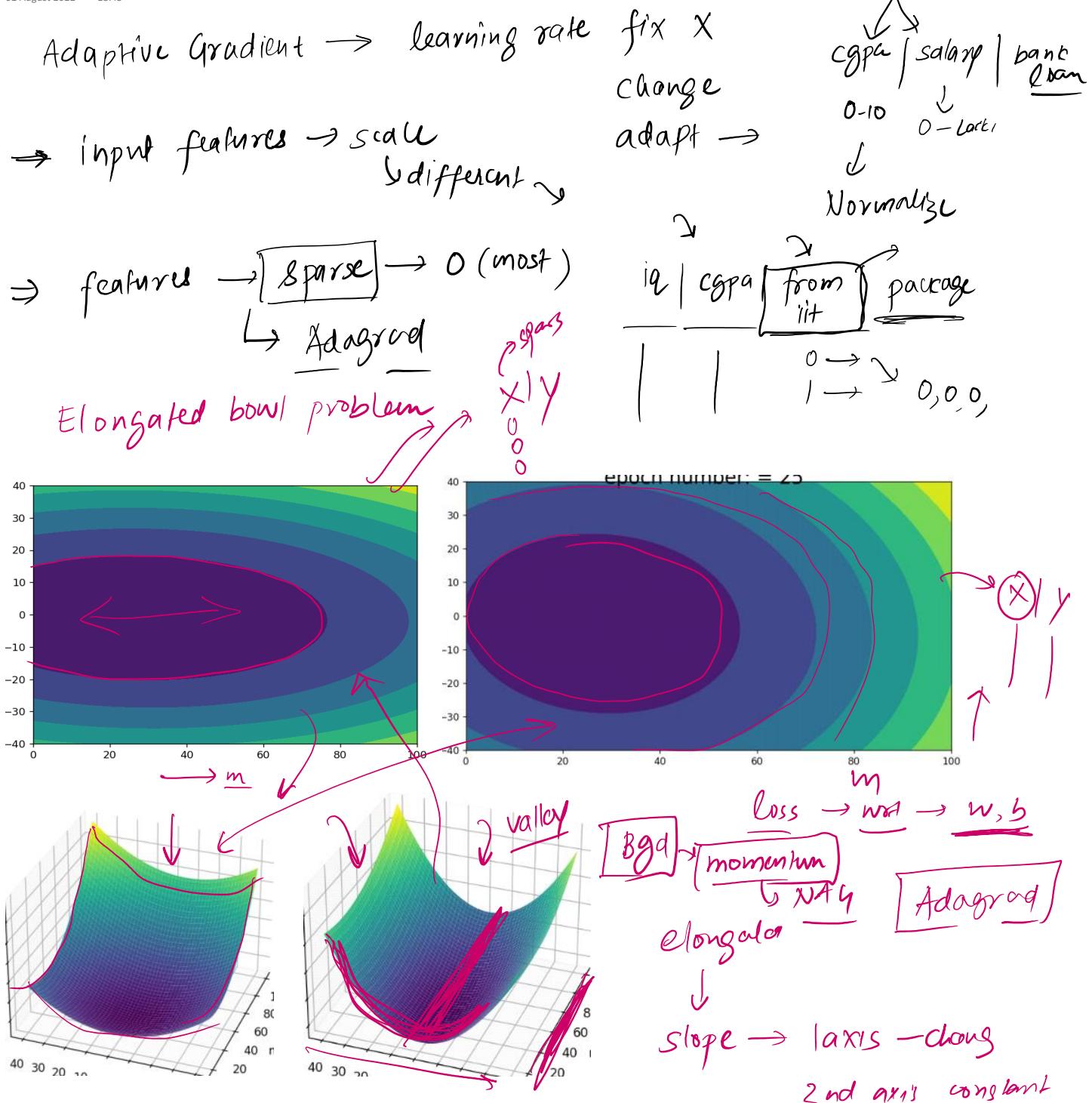


NAG

momentum = 0.9

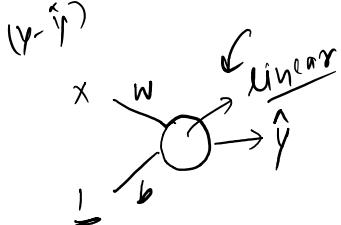
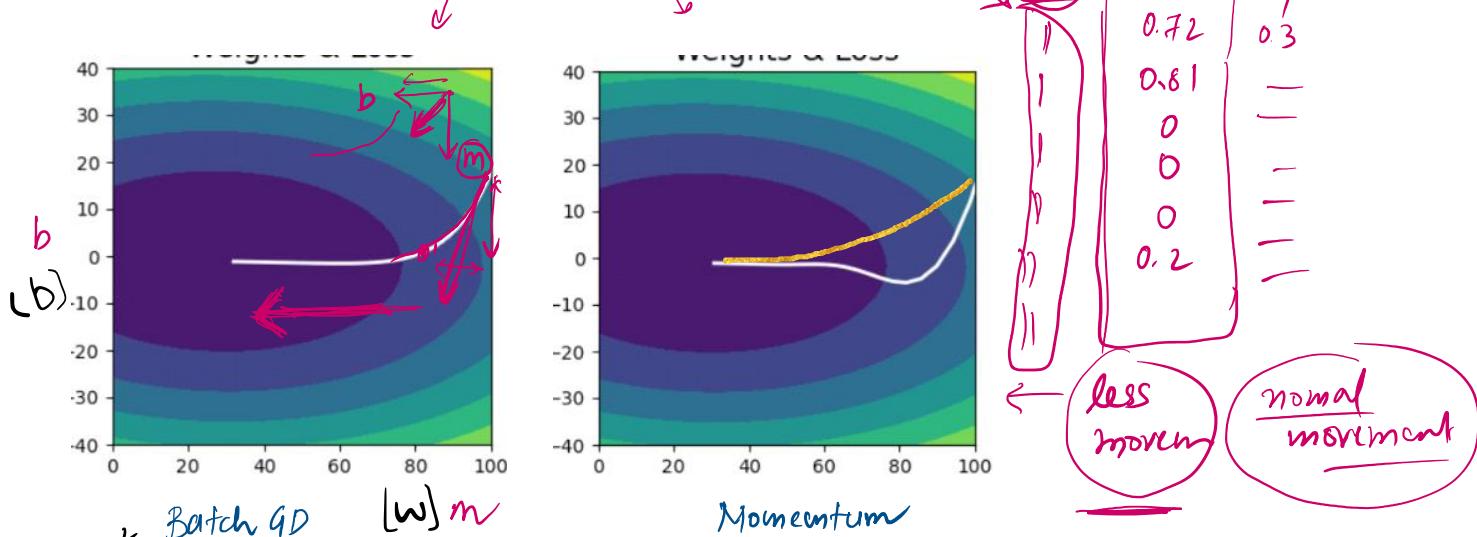
nesterov = True





How optimizers behave (Why?)

02 August 2022 18:44



BGD → sparse

for i in epochs:

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

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

100 rows 0 sparse

small in every epoch

big update in every epoch

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

$$= -2(y - \hat{y}) \underline{x}$$

add

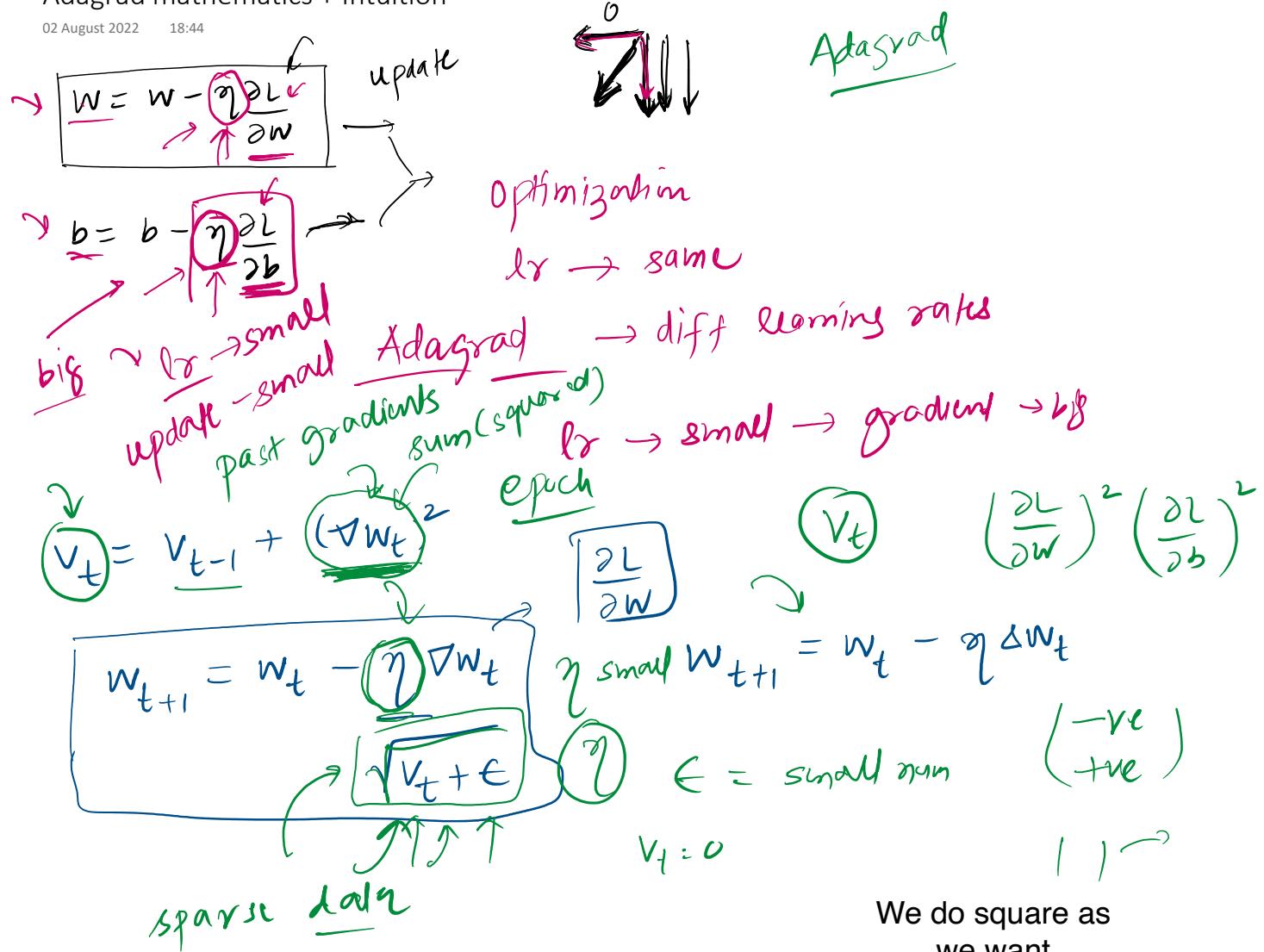
$$\frac{\partial L}{\partial b} = \underline{-2(y - \hat{y}) \times 1}$$

bignum

sparse col →
small non sparse
normal / big

Adagrad mathematics + Intuition

02 August 2022 18:44

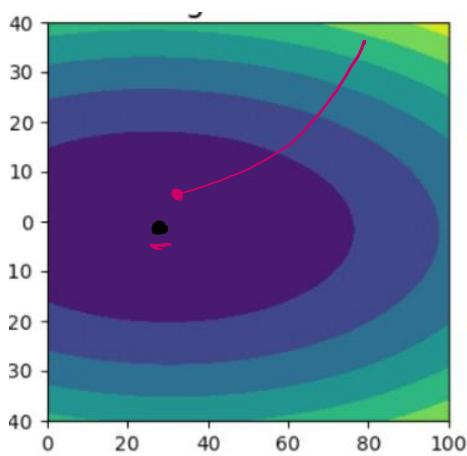


We do square as
we want
magnitude, not
do mod as we
have to do
backpropagation
also.

Disadvantage

02 August 2022 18:43

NN \rightarrow use $\times \propto$ linear degrees



Initially it is fast to reach near , but later this update become small , and very hard to converge towards global minima

The Why?

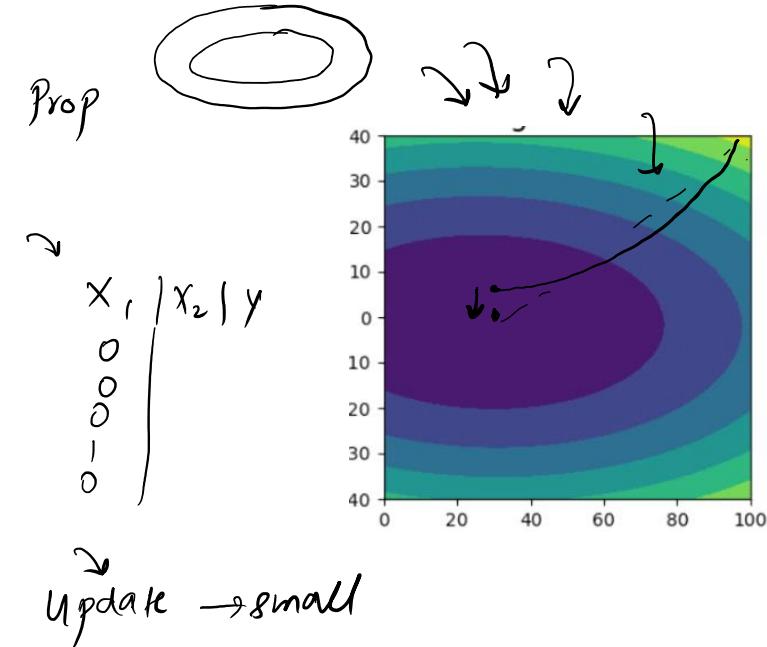
03 August 2022 14:03

RMSprop → Root Mean square Prop
↑ improvement

→ Adagrad ← sparse data

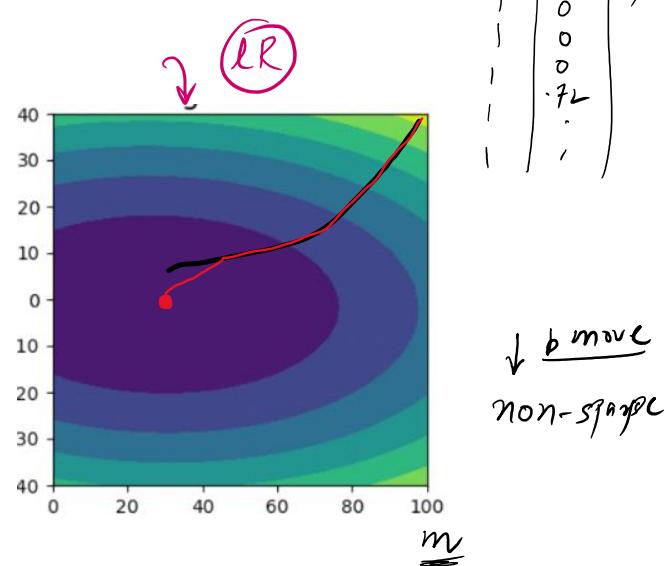
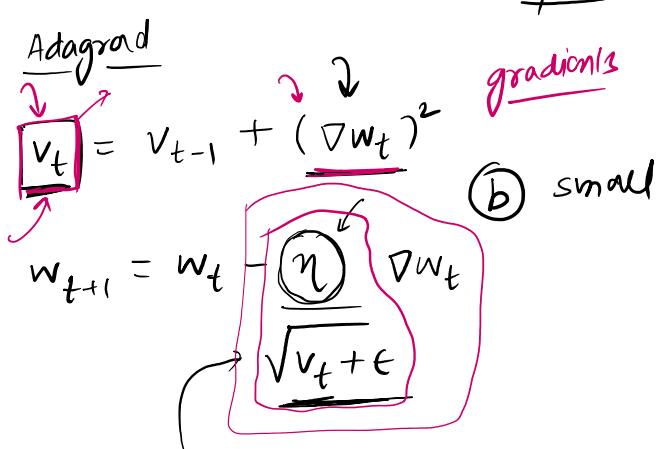
BGD → momentum

Disadvantage → learning rate



Mathematical Formulation

03 August 2022 14:04



RMSprop

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

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

exp decaying avg
old epoch

$$\beta = 0.95$$

$$\rightarrow v_0 = 0$$

$$\rightarrow v_1 = 0.95 \times 0 + 0.05(\nabla w_1)^2$$

$$\rightarrow v_2 = 0.95 \times 0.05(\nabla w_1)^2 + 0.05(\nabla w_2)^2$$

$$\rightarrow v_3 = \frac{0.95 \times 0.95 \times 0.05(\nabla w_1)^2}{\text{epoch 1}} + \frac{0.95 \times 0.05(\nabla w_2)^2}{\text{epoch 2}} + \frac{0.05(\nabla w_3)^2}{\text{epoch 3}}$$

↑ smaller ↑ Epoch 1 ↑

v_t shoot $\rightarrow \eta$ small

Adagrad \rightarrow NN \rightarrow non-convex optimization

convex optimis \rightarrow linear $\xrightarrow{\sigma \text{ eq}}$

Not performs good in non-convex optimization.

Performs good in convex optimization.

✓ convex \rightarrow linear
optim's seg

global
minima RMSprop

Disadvantage

03 August 2022 14:54

Before adam is there , we used RMSProp only .

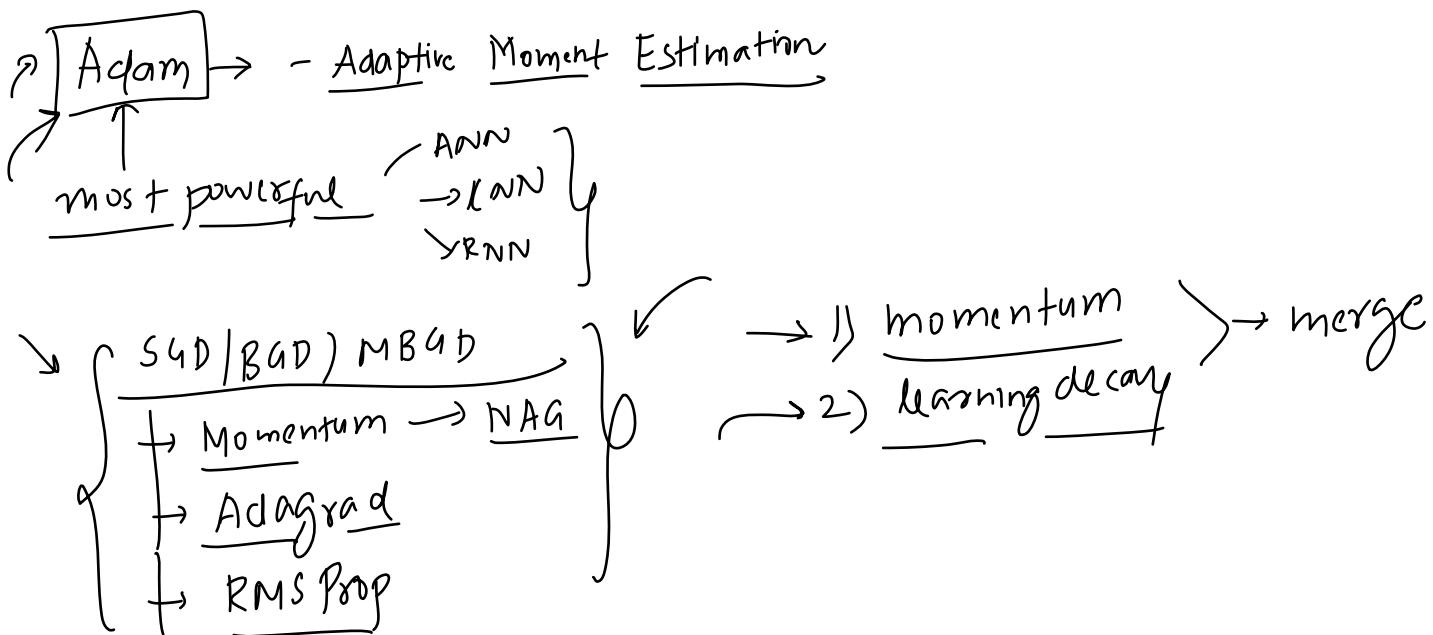
Infact today if adam not works , we consider RMSProp

No,
Optimization

↓
Adam

Introduction

04 August 2022 10:45



Momentum : It introduces the concept of momentum i.e using history values also. (ewma on gradients)

NAG : It uses same momentum concept with little modification of the point where gradients are calculated.

Adagrad : It considers change learning rate.

RMSProp : In overtime , adagrad fails to converge towards global minima due to very low learning rate, it adds exp weighted moving avg to modify learning rate.

Adam : It merges concept of both momentum and learning rate. By using concept of momentum for gradients and learning rate change .

Mathematical Formulation

04 August 2022 10:45

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} * m_t \quad \xrightarrow{\text{Keras}} \boxed{\text{Bias correction}} \quad \text{epoch \#}$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^{[t]}} \rightarrow \hat{v}_t = \frac{v_t}{1 - \beta_2^{[t]}}$$

$$\beta_1 = 0.9 \quad \beta_2 = 0.99 \rightarrow \underline{\text{Keras}}$$

where

$$\rightarrow \boxed{m_t = \beta_1 m_{t-1} + (1-\beta_1) \nabla w_t} \rightarrow \text{momentum}$$

$$\rightarrow \boxed{v_t = \beta_2 v_{t-1} + (1-\beta_2) (\nabla w_t)^2} \rightarrow \text{Adagrad}$$

$$m_0 = 0 \quad m_t = 0 \quad v_t = 0 \quad 0 \rightarrow 0$$

$$\eta = 0.1 \quad 0.01$$

This bias correction is there , as in initial weights as zero so kick off we used it.

Verdict

04 August 2022 10:45

