

Optimizer

1

$$G.D. \Rightarrow W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial w} \Rightarrow \text{entire Data}$$

1

Epoch \Rightarrow Complete cycle with entire Data

[FP + BP]

2

SGD \Rightarrow Single Point at a time

2

Iteration — SGD, Mini Batch GD

3

Mini-Batch GD \Rightarrow Batch.

3

Batch size

4

Momentum (SGD with momentum)

5

NAG (Nesterov Accelerated gradient)

- Reduce the Loss
- Train a Variable Parameter in Best Possible way.
- Fast convergence to global minima

6

Adgarde (Adaptive gradient)

7

RMS Prop (root-mean Square Propogation)

8

ADAM



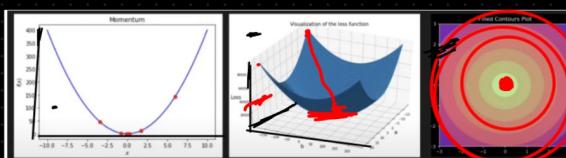
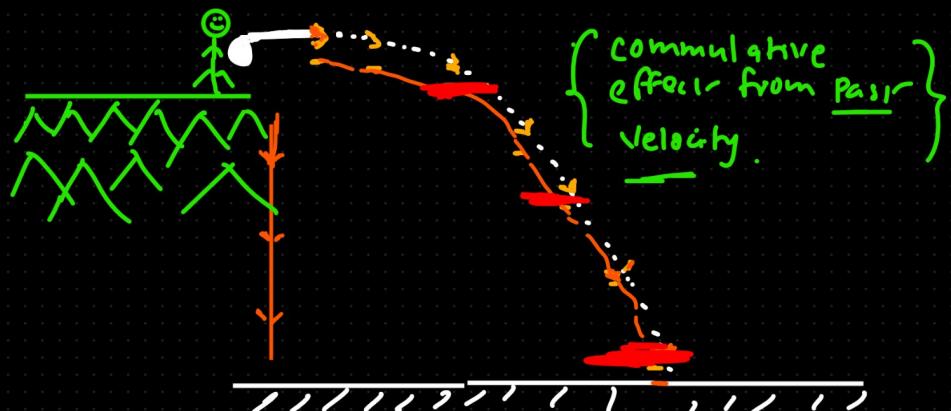
Minute change in Gradient - Decreasing itself²

1

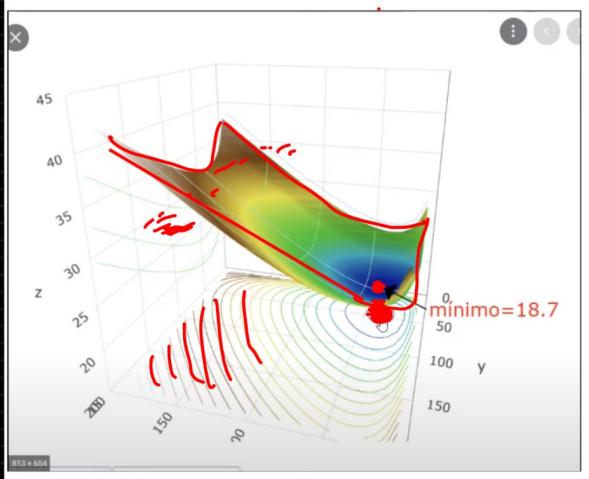
Momentum (SGD with momentum)

Phy (Accerlation) change into the Pushon

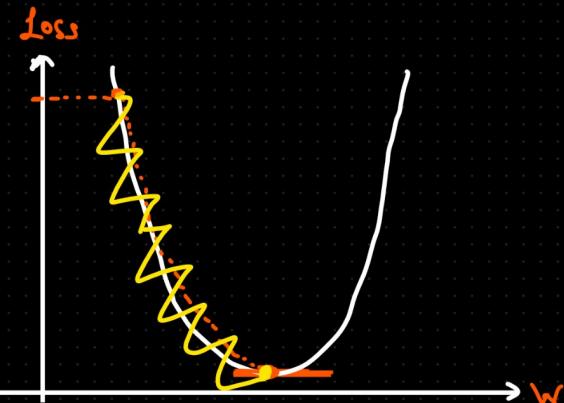
$$M \times V_f =$$



← counter



$\zeta_D, S\zeta_D, M\zeta_D$

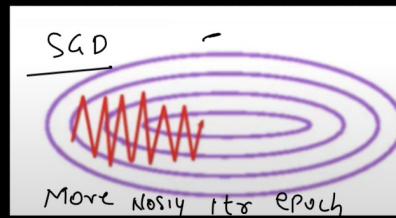


$\zeta_D \Rightarrow$ More time to converge

More memory

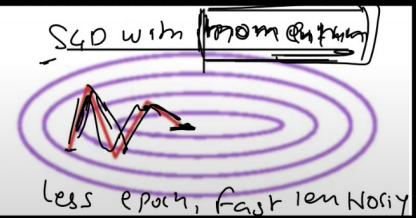
$S\zeta_D$

epoch, iteration, Noise



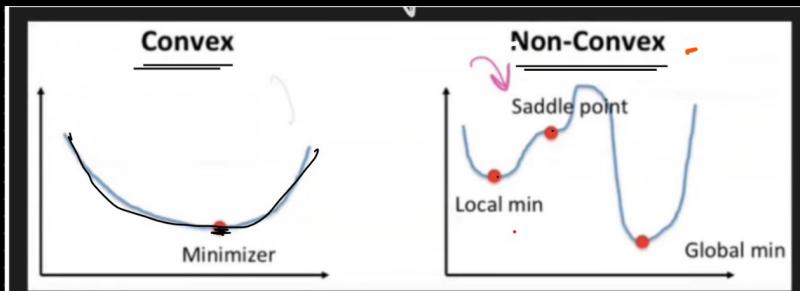
SGD

More noisy it's epoch

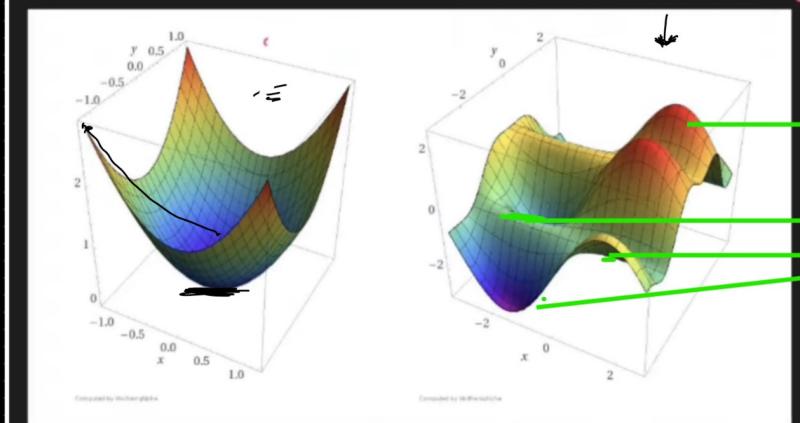


SGD with momentum

Less epoch, faster in noisy



Flat Surface \Rightarrow Slope = 0 (would be w/ update)

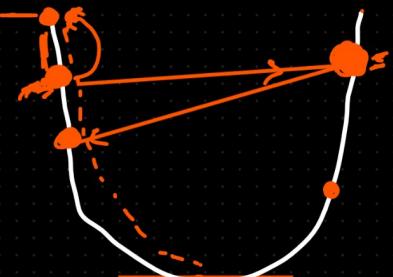


Curvature
Saddle Point
Local minima
Noisy

Global minima

G.D.

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial w}$$



Right dir, faster

{ SGD
Momentum

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

$$\text{MOM} = M \times \nabla$$



$$V_t = \beta V_{t-1} + \eta \frac{\partial L}{\partial w}$$

β = decay factor

$$0 < \beta < 1 \Rightarrow \text{0.9}$$

$$w_{t+1} = w_t - \left[\beta v_{t-1} + \eta \frac{\partial L}{\partial w} \right]$$

\Rightarrow EWMA

Adding the momentum

Considering a history of gradient



D_1	25	$t=0$
D_2	13	$t=1$
D_3	14	$t=2$
D_4	31	$t=3$
D_5	43	$t=4$

SMA Window = 3, 5, 7, 8, 9, 10, ... n

$= \frac{\text{MA}}{\text{---}}$

trend, pattern, smoothing

EWMA

$$= \underline{\text{EWMA}} =$$

DL, Optimizer

$$\boxed{Y_t = \beta X_{t-1} + (1-\beta) \theta_t}$$

at t time at t-1 time Value at t
 EWMA EWMA time

$$\begin{cases} Y_0 = \theta_0 = 2 \\ Y_0 = \theta_0 = 0 \end{cases}$$

$$\boxed{0 < \beta < 1}$$

$\beta = 0.9$

$$Y_1 = 0.9 \times 0 + (0.1) \times (1.3)$$

$$Y_1 = 0 + 1.3$$

$$\boxed{Y_1 = 1.3} =$$

$$\begin{aligned} Y_2 &= 0.9 \times 1.3 + (0.1) \times (1.7) \\ &= \underline{1.17 + 1.7} \end{aligned}$$

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

$$V_0 = 0$$

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

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

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

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

$$= \beta [\beta (1-\beta) \theta_2 + (1-\beta) \theta_2] + (1-\beta) \theta_3$$

$$= \beta^2 (1-\beta) \theta_2 + \beta (1-\beta) \theta_2 + (1-\beta) \theta_3$$

$$= (1-\beta) \left[\underline{\beta^2 \theta_2} + \beta \underline{\theta_2} + \underline{\theta_3} \right]$$

↑ ↑ ↑
 V_t line value 3
 $t+1$ $t+2$

$$\frac{0.5 \times 0.5}{0.81}$$

?

\downarrow

$$\frac{0 < \beta < 1}{\underline{\beta^2 \theta_2} + \underline{\beta \theta_2} + \underline{\theta_3}}$$

$\cancel{\beta^2 \theta_2}$ $\cancel{\beta \theta_2}$ $\cancel{\theta_3}$

NAG \Rightarrow Nesterov Accelerated Gradient (Updated Version of Momentum)

$$\text{C.D.} \Rightarrow w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w}$$

Momentum in UD

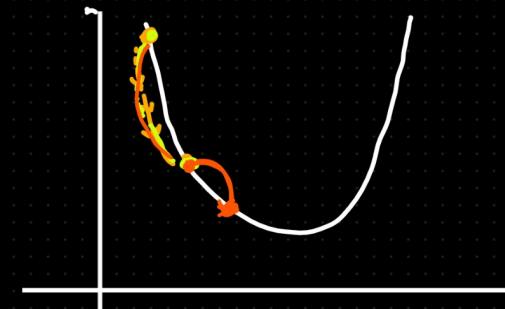
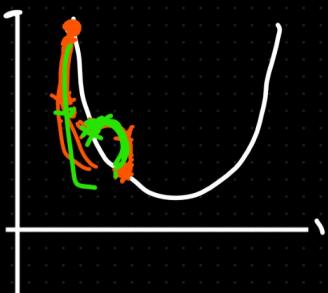
$$w_{t+1} = w_t - v_t$$

$$v_t = \beta v_{t-1} + \eta \frac{\partial L}{\partial w}$$

$$w_{t+1} = w_t - \left[\beta v_{t-1} + \eta \frac{\partial L}{\partial w} \right]$$

history of Gradient + gradient at point r

NAG \Rightarrow Look ahead factor



$$\beta v_{t-1} = w_t - w_{Lq}$$

$$w_{Lq} = w_t - \beta v_{t-1}$$

$$v_t = \beta v_{t-1} + \eta \frac{\partial L}{\partial w_{Lq}}$$

$$v_t = (w_t - w_{Lq}) + \eta \frac{\partial L}{\partial w_{Lq}}$$

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

Adagrad (Adaptive gradient)

$$GD \Rightarrow W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial w}$$

Learning rate.

Dynamically

Momentum

NAG

Updated momentum

Look ahead factor

Adding Past. history of Gradient \Rightarrow Acceleration

$$W_t = W_{t-1} - \eta \frac{\partial L}{\partial w}$$

$$\cancel{W_t = W_{t-1} - \eta \frac{\partial L}{\partial w}}$$

LR

D.G.

$$= V_t = V_{t-1} + \left(\frac{\partial L}{\partial w} \right)^2$$

Past gradient sum of square

$\epsilon = \text{small value}$

= RMS-Prop \rightarrow

Adadeltq

small \downarrow
Dynamik

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} X \frac{\partial L}{\partial w},$$

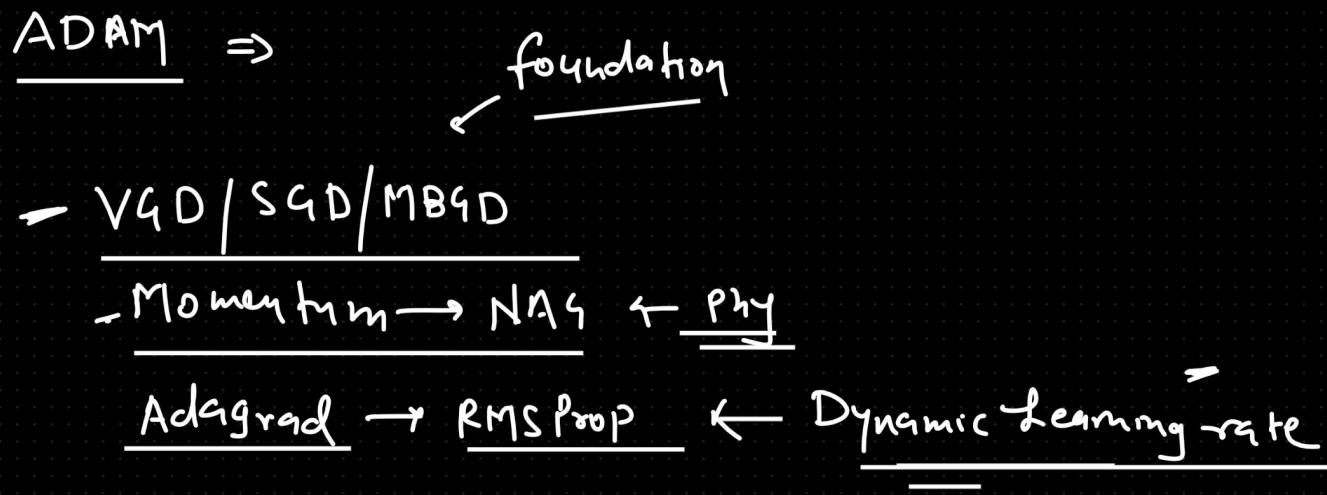
Adagrad

- Adagrad

$$v_t = v_{t-1} + \left(\frac{\partial L}{\partial w} \right)^2$$

- RMSprop

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



ADAM \Rightarrow Momentum + Dynamic Learning rate

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

$$\boxed{m_t = \beta_1 m_{t-1} + (1-\beta_1) \frac{\partial L}{\partial w}}$$

$$\boxed{v_t = \beta_2 v_{t-1} + (1-\beta_2) \frac{\partial L}{\partial w}}$$

Friday → pm
8 to 11 pm
hyperparameter tuning in ANN

Saturday →
- Vanishing / exploding grad -
→ - BN / Reg | Dropout -
Overfitting
- Weight initialization -

