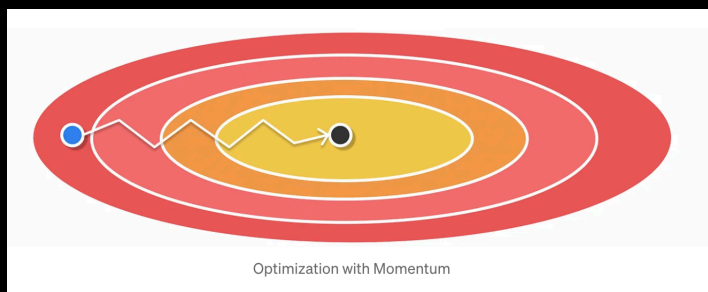
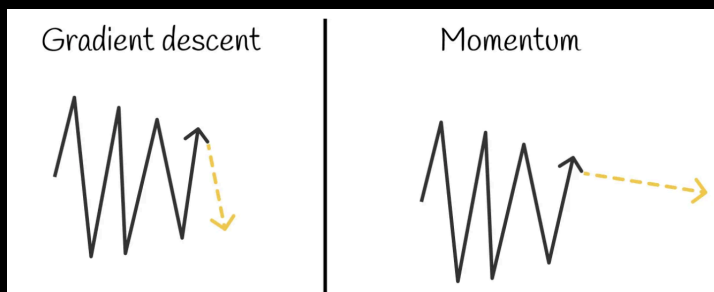


SGD with Momentum



$$\omega_{\text{new}} = \omega_{\text{old}} - \eta \frac{\partial L}{\partial \omega_{\text{old}}} \rightarrow \omega_t = \omega_{t-1} - \eta \frac{\partial L}{\partial \omega_{t-1}}$$

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

$$t \rightarrow (t-1)$$

$$t_{-2}, t_{-1}, (t), t_{+1}, t_{+2}$$

Past old / happened instance

$\omega_{t-1} \approx \omega_t$
 impact of ω_{t-1}
 in the output
 of ω_t

Time Series

Removal \rightarrow Noise / Filter / Oscillations

1) Exponential Weighted Average ✓

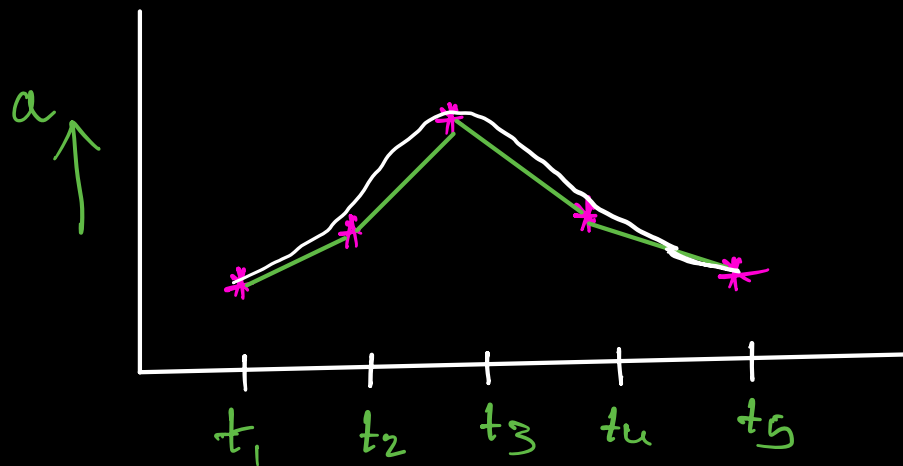
2) Rolling Mean ✗

Exponential Weighted Average (EWA) Smoothing

Time t_1 t_2 t_3 t_4 \dots t_n

Value a_1 a_2 a_3 a_4 \dots a_n

t_1 should be controlling
 t_2



$$\text{EWA}_{V_{t_2}} = \boxed{\beta * V_{t_1} + (1 - \beta) * V_{t_2}}$$

$$= (0.95 * a_1) + ((1 - 0.95) * a_2)$$

Recommended β value should be close to 1

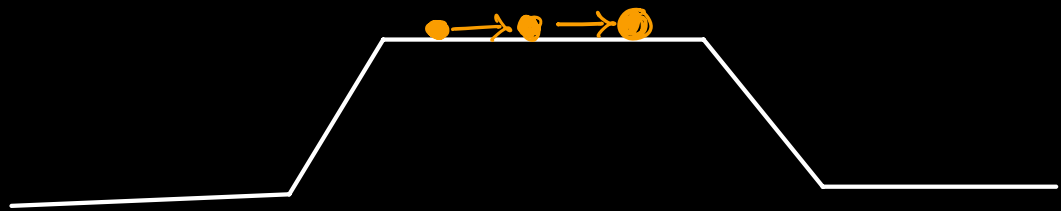
$\beta = \underline{0.95}$

$$= \left\{ \underbrace{0.95 * a_1}_{\text{High Impact}} + \underbrace{0.05 * a_2}_{\text{Low impact}} \right\}$$

$$\text{EWA}_{\text{for } V_{t_3}} = \beta * V_{t_2} + (1 - \beta) * V_{t_3}$$

$$= 0.95 \{ \quad \} + (1 - \beta) a_3$$

Momentum is the push towards global Minima

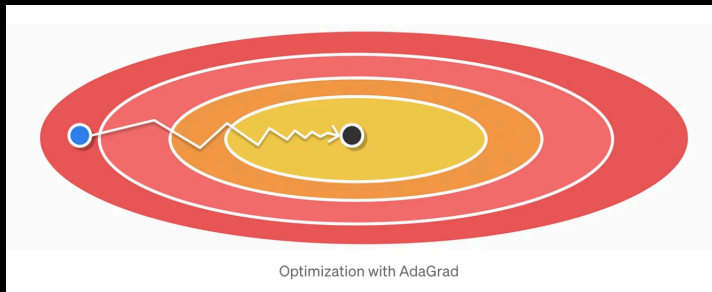


Plateau Region Problems

Advantages of SGD + Momentum

- 1) Noise Reduction
- 2) Smoothing of the noise
- 3) Smooth Convergence

Adagrad (Adaptive Gradient Descent)



Learning Rate η = fixed

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

$$\eta' = \frac{\eta}{\text{large value}}$$

$$\eta' = \text{small value}$$

$$\eta' = \frac{\eta}{\sqrt{\alpha^t + \epsilon}} \rightarrow \text{very small value}$$

$$\alpha^t = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_t} \right)^2 \rightarrow \text{large value}$$

Slope Squared

Adaptive Learning \rightarrow Dynamic Learning Rate

When we start we need high L.R

When we converge we need lower L.R

Disadvantage

1) Slow Convergence only in the end.

\hookrightarrow very small update

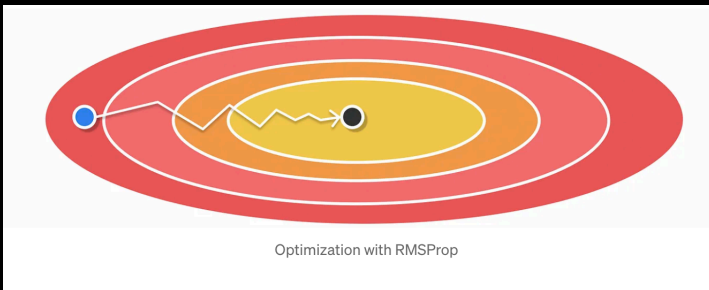
2) $\eta' =$ Possibility to come α will be very small
 $\alpha^+ \approx 0$

Note : —

✓ Solved with Momentum (Academically)

✓ Adam (Industry Applications)

RMSProp (Root Mean Squared Propagation)



(EWA) \rightarrow Exponential
Weighted
Average

$$\eta' = \frac{\eta}{\sqrt{Sdw_t} + \epsilon}$$

$$S_{dw_t} = 0$$

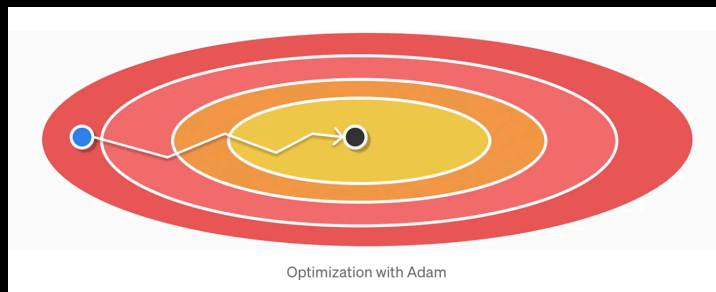
$$= \underbrace{\beta * S_{dw_{t-1}}}_{\text{green}} + \underbrace{(1-\beta) \left(\frac{\partial L}{\partial w_t} \right)^2}_{\text{blue}} *$$

① SGD with Momentum \rightarrow No DLR

② Adagrad \rightarrow DLR

③ RMSProp \rightarrow SGD with Momentum + Exponential Weighted Avg LR

Adam Optimizer (Adaptive Momentum Estimation)



\Rightarrow RMSProp + SGD with M
D.L.R + SGD with M

$$w_t = w_{t-1} - \eta' \nabla_{dw_t}$$

where, $\nabla_{dw_t} = \beta * \nabla_{dw_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial w_t} \right)$

$$\eta' = \frac{\eta}{\sqrt{S_{dw_t} + \epsilon}}$$

Bias Calculation

$$b_t = b_{t-1} - \eta \nabla_{\theta} L_t$$

where $\nabla_{\theta} L_t = \beta * \nabla_{\theta} L_{t-1} + (1-\beta) \left(\frac{\partial L}{\partial \theta} \right)_t$

$$\eta = \frac{\eta_0}{\sqrt{S_{\theta} + \epsilon}}$$

Learning Rate Scheduler

Dynamic Learning Rate

$$\text{Epochs} = \underline{100}$$

$$\text{if epochs} \leq 25;$$

$$\text{lr} = 0.01$$

$$\text{elif epochs} > 25 \text{ and } \leq 50$$

$$\text{lr} = 0.001$$

$$\text{elif epochs} > 50 \text{ and } \leq 75$$

$$\text{lr} = 0.0001$$

$$\text{else epochs} > 75$$

$$\text{lr} = 0.000001$$

CLR

Remaining Topics

1) Back propagation

2) Regularization

3) Callbacks

Task Val Accuracy = 98%.
 Loss = 0.02

①

MNIST

↳ Different Activation
Functions

mish
↳ Optimizers

Final Report

8 Activation, 8 Optimizers
= 64 combinations