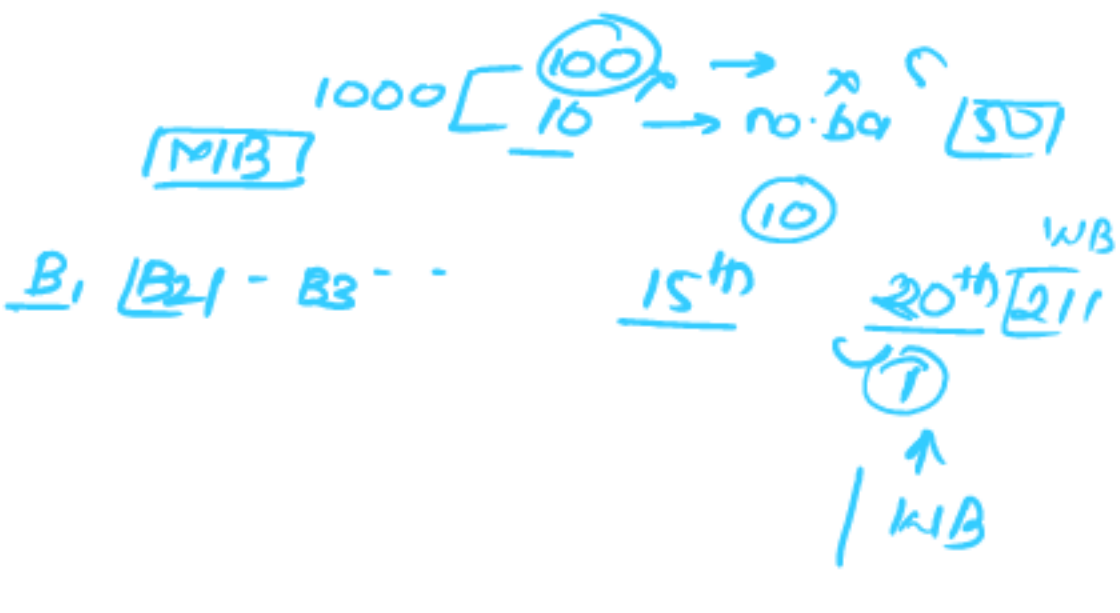


Stochastic gradient Descent with Momentum.

$$gD \rightarrow Nq \rightarrow \text{CUG SMOOTH Diff}$$

$$w_{new} = w_{old} - \alpha \frac{\partial L}{\partial w_{old}}$$

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



→ pop →

$$v_{dw,t} = \beta \times v_{dw,t-1} + (1-\beta) \times \frac{\partial L}{\partial w,t}$$

1st previous / Today.

$$\beta = \text{importance} \quad 0.8 = 5^{th} \quad 0.9 = 10^{th}$$

1 2 3 4 5 6 7 8 9 10

$$v_{dw} = 0 \quad v_{db} = 0$$

$$w_{new} = w_{old} - \alpha \frac{\partial L}{\partial w_{old}}$$

derivative of loss w.r.t weights.

$$v_{dw,t} = v_{dw,t-1} \times \beta + (1-\beta) \times \frac{\partial L}{\partial w,t}$$

$$v_{t1} = \beta \times v_{t0} + (1-\beta) \times a_1 \quad 0.9$$

$$= 0.1 \times a_1$$

$$v_{t2} = \beta \times v_{t1} + (1-\beta) \times a_2$$

$$= 0.9 \times (0.1 \times a_1) + (1-\beta) \times a_2$$

$$= 0.9 \times (0.1 \times a_1) + 0.1 \times a_2$$

$$v_{t3} = \beta \times v_{t2} + (1-\beta) \times a_3$$

$$= 0.9 \times (0.9 \times 0.1 \times a_1 + 0.1 \times a_2) + (0.1 \times a_3)$$

PIOW, B

$$\beta = 0.8 \quad 5^{th}$$

Current weight and bias.

new.

noise

the algorithm become MBSGD

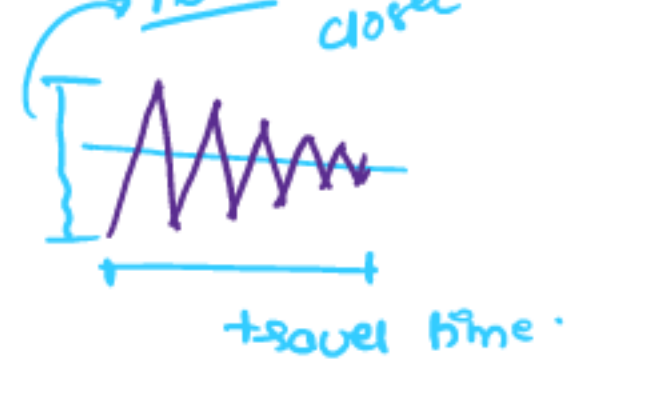
algorithm goes into dynamic equilibrium keep moving on error surface.



By using Exponential moving avg.

We have avg of the oscillations in vertical direction close to zero.

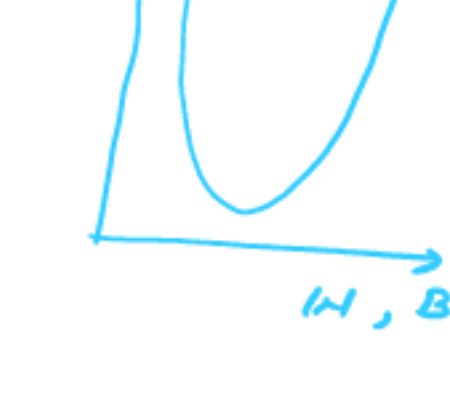
It enables our algorithm to take straight path in forward direction to global minima and dampen but the vertical oscillation.



B value ↑ smoothing ↑

B value ↓ smoothing ↓

Adagrad.



gD & constant
sgD & constant
MBSGD & constant
gDWM & constant

Core Logic

Dynamic Learning Rate.

Range

$$0 \text{ to } 1 \quad \beta = 0.01$$

sgD cont, Age cont, grade M/P

In Realworld Model will learn dense feature 1st and then sparse feature.

$$w_{new} = w_{old} - \alpha \frac{\partial L}{\partial w_{old}}$$

$$w_{new} = w_{old} - \alpha' \frac{\partial L}{\partial w_{old}}$$

Dynamic LR

$$\alpha' = \frac{\alpha}{\sqrt{nt + \epsilon}}$$

small positive no.

Derivative of loss w.r.t weights

A small positive number to avoid zero division error.

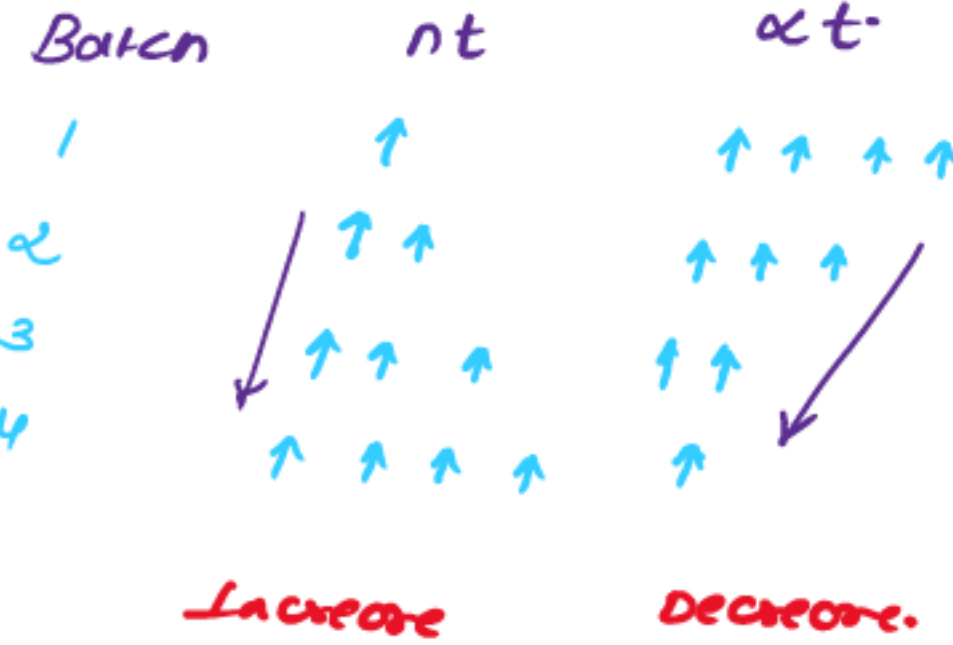
$$\epsilon = 10^{-7}, 10^{-8}$$

$$nt = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_i} \right)^2$$

summation.

$$\beta_1 + \dots + \beta_5 + \dots + \beta_{10}$$

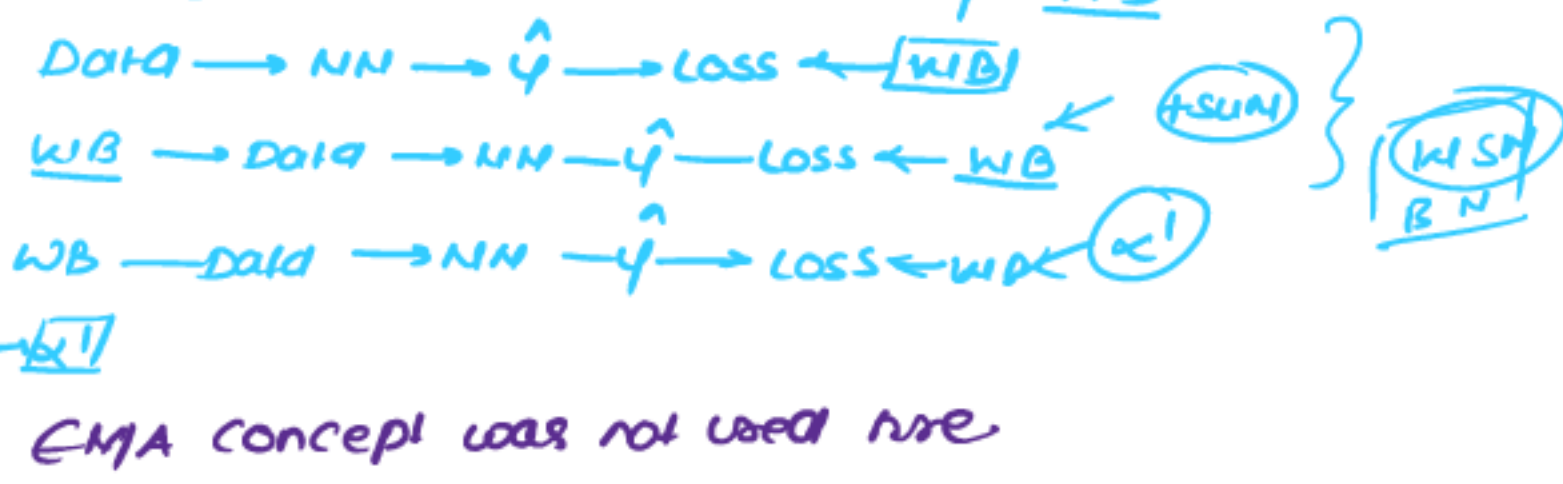
$$nt \propto \frac{1}{\alpha'}$$



$$\begin{matrix} 0.01 \\ 0.009 \\ 0.007 \\ 0.004 \\ 0.001 \end{matrix}$$



Problem faced here is called Learning Rate Decay.



EMA concept was not used here.

Learning Rate reduces Monotonically.

Adadelta or RMS prop.

$$w_{new} = w_{old} - \alpha \frac{\partial L}{\partial w_{old}}$$

$$w_{new} = w_{old} - \alpha' \frac{\partial L}{\partial w_{old}}$$

Dynamic Learning Rate.

$$\alpha' = \frac{\alpha}{\sqrt{s_{dw} + \epsilon}}$$

so now they have

introduced the concept of EMA again to get rid of Learning Rate Decay.

$$s_{dw} = \beta \times s_{dw,t-1} + (1-\beta) \left(\frac{\partial L}{\partial w} \right)^2 \quad \beta = 0.9$$

It will not allow decay of Learning Rate to zero

$$1.8 \rightarrow 0.8^{th} \rightarrow 50^{th}$$

Adadelta is nothing but extension of Adagrad and also trying to solve the problem Learning Rate Decay.

In Adadelta we will not try to accumulate all past sq. gradient values instead tries to restrict the window of accumulated past gradient