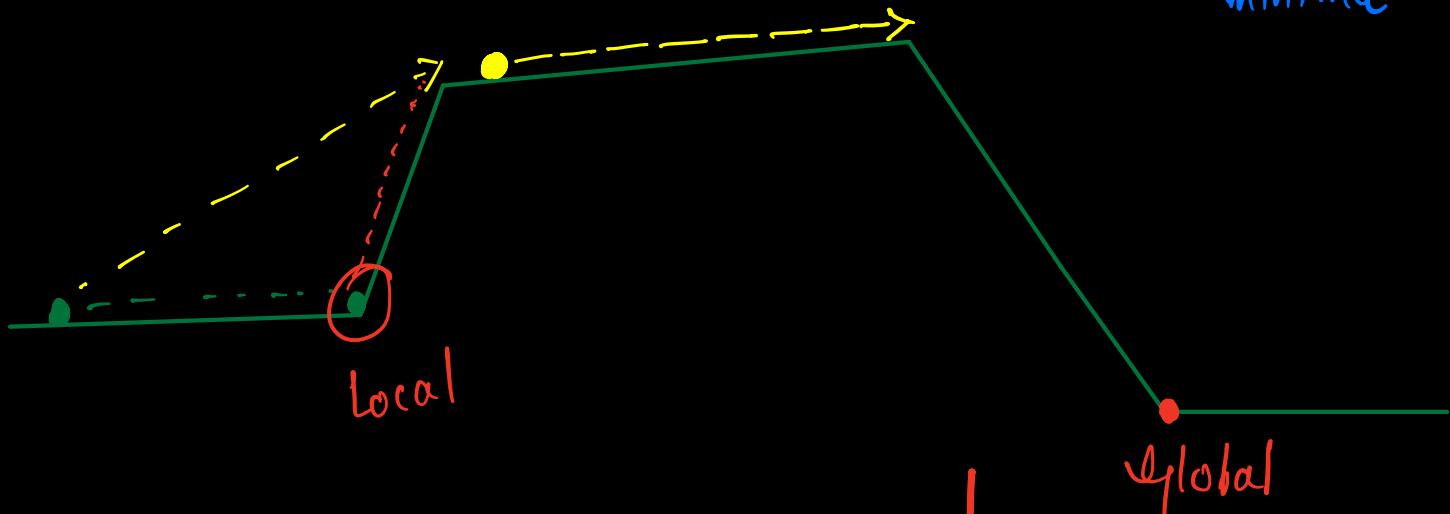
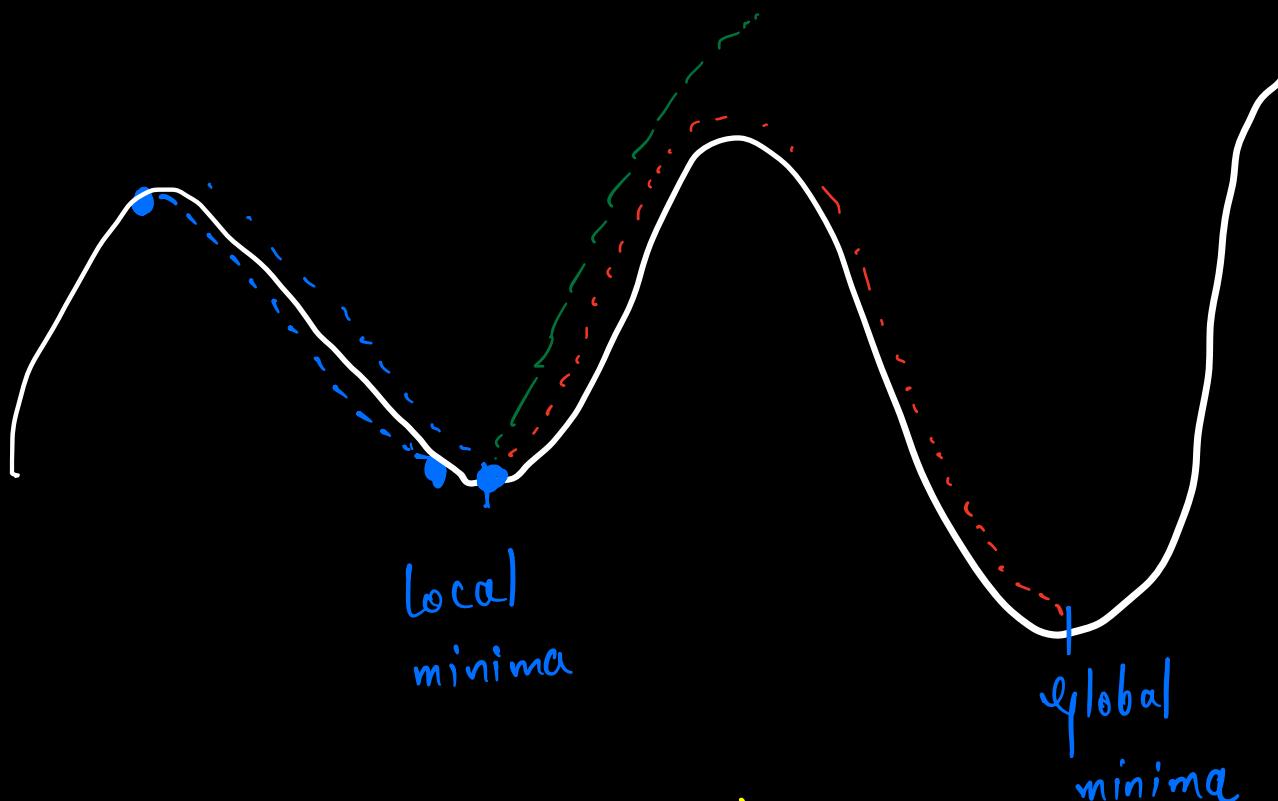


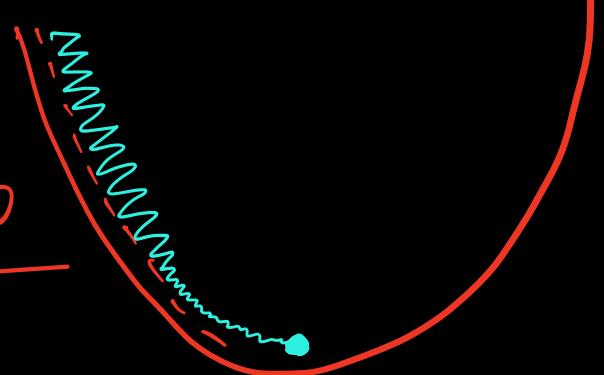
Today's Agenda

Optimizers



Steps:-

updating $\omega \& b$



Types

i) Batch Gradient Descent

i) Pass the entire dataset once.

ii) High computational power

Cons :- Slow,

Weight Update :-

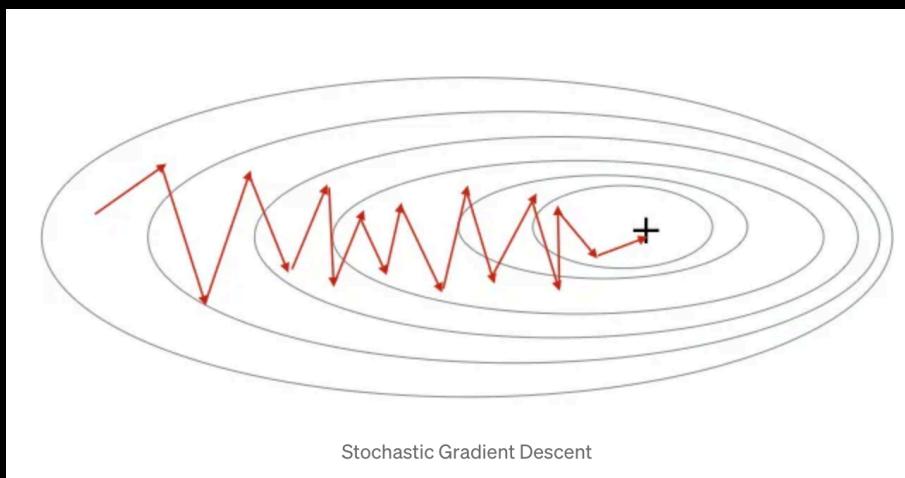
$(\alpha) \rightarrow L \cdot R$

$$w_N = w_0 - \underbrace{\alpha \cdot \frac{1}{N} \sum_{i=0}^N}_{\text{L.R}} \underbrace{\frac{\nabla_w L_i(w)}{\left\{ \frac{\partial \text{LOSS}}{\partial w_{\text{old}}} \right\}}}_{\text{Loss}}$$

Pro:- Very accurate updates

(ii) Stochastic Gradient Descent

(i) Pass one data point at one time.



Con:- Noisy updates $\rightarrow \{ \text{Jitter} \}$

Weight Update

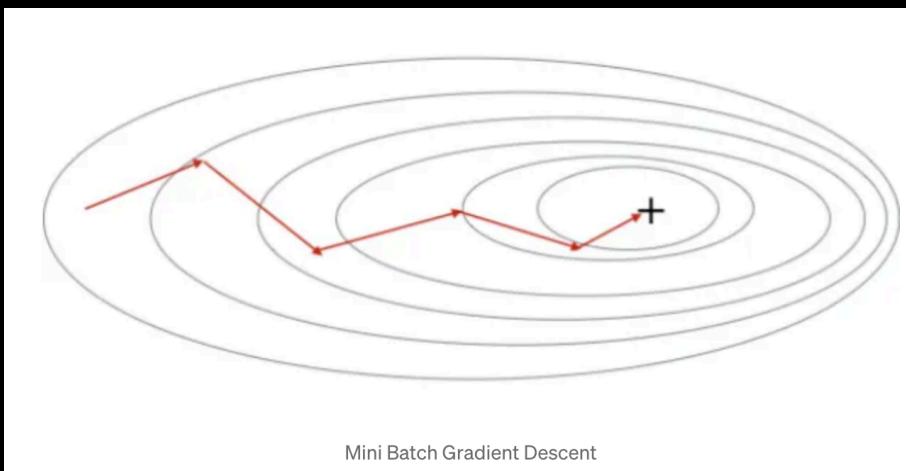
$$w_N = w_0 - \gamma \nabla_w L(w) \rightarrow \frac{\partial \text{Loss}}{\partial w_0}$$

L-R \rightarrow WRT Time , if should decrease

$\{ \text{Dynamic LR} \} \rightarrow \text{Epochs}$

- 1
- 0.1
- 0.01
- 0.001
- 0.0001
- 0.00001
- 0.000001

Mini Batch Gradient Descent



OOD

Batch \rightarrow group of datapoints

$$2^n \rightarrow 8, 16, 32, 64, 128, 256 \quad \frac{\text{VRAM}}{\text{GPU}}$$

Pro:- Smooth convergence, lesser noise
move stable than SGD.

Con:- Needs Tuning of batch size

Weight Update Rule:-

$$\omega_{\text{new}} = \omega_{\text{old}} - \eta \cdot \frac{1}{m} \sum_{n=1}^m \nabla_{\omega} L(\omega)$$

m = mini batch size, $m \leq N$

1000 data points

batch size = 100

Total batches = 10

Data Representation

i) Batch

2) Iteration
one update of the model using batch

SGD :- 1000 data points
 1000 iterations

Mini :- Batch size = 100
 Total iteration 10

Epochs

Full cycle of the entire dataset.

Model \rightarrow all training samples

1 Epoch \rightarrow entire dataset

10 Epochs \rightarrow entire dataset (10X)

1000 \rightarrow DP

100 \rightarrow BS

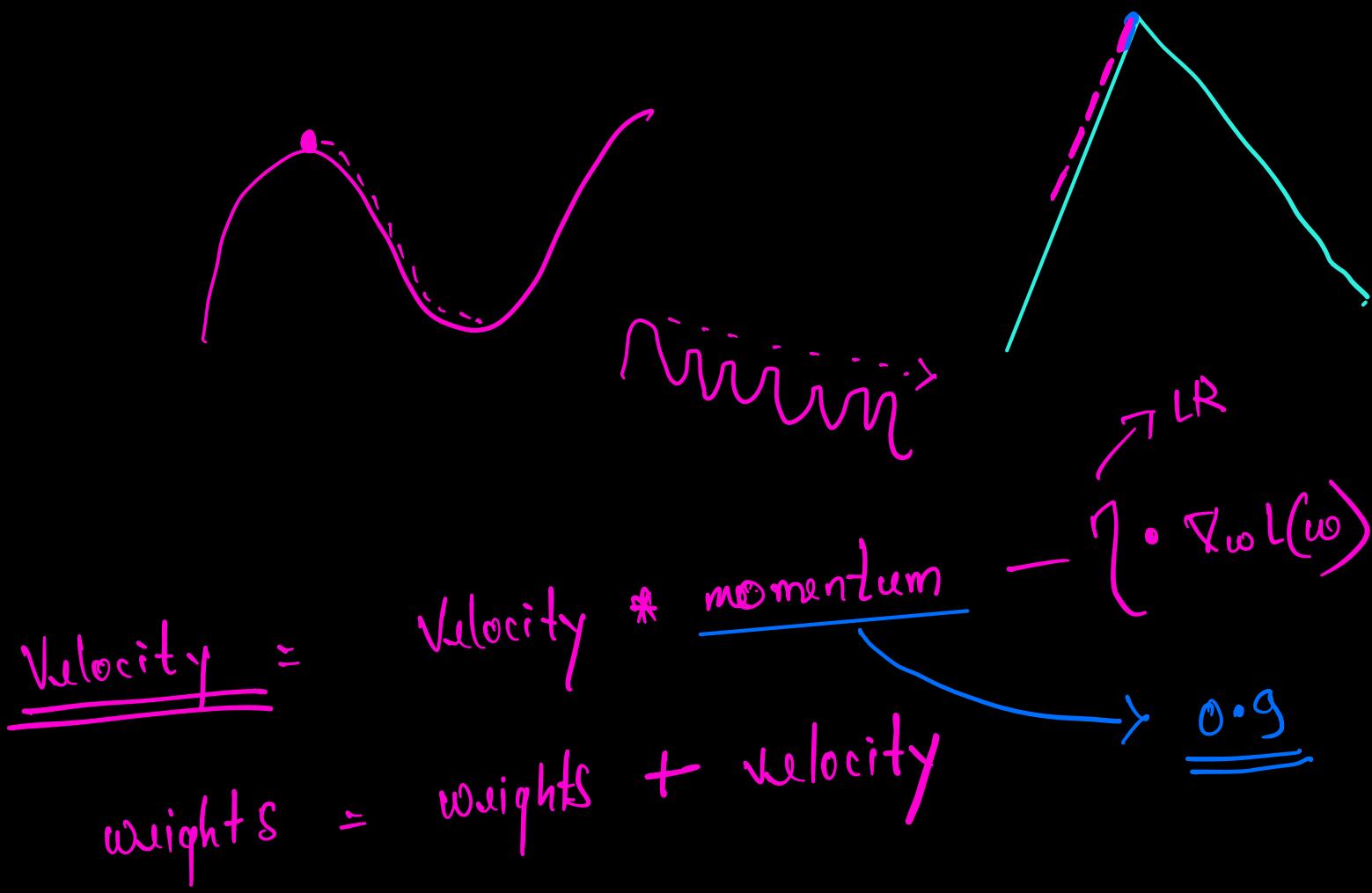
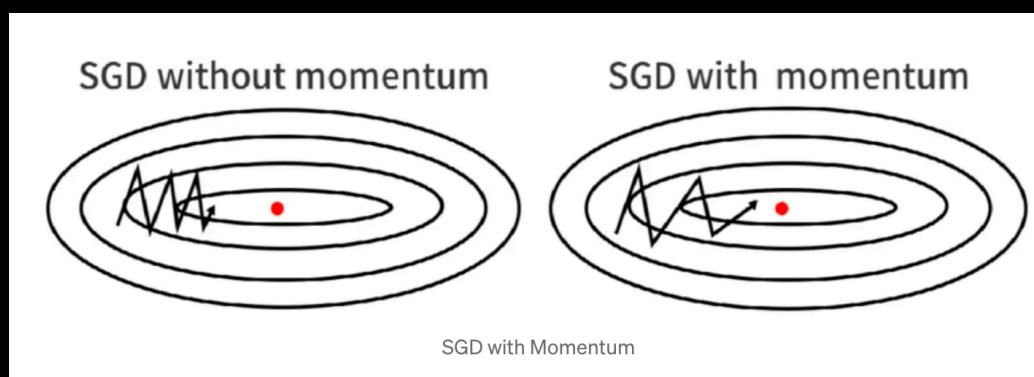
Total iterations for 1 epoch = 10
10 epochs = 10×10
= 100

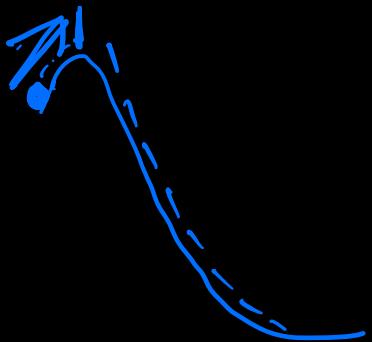
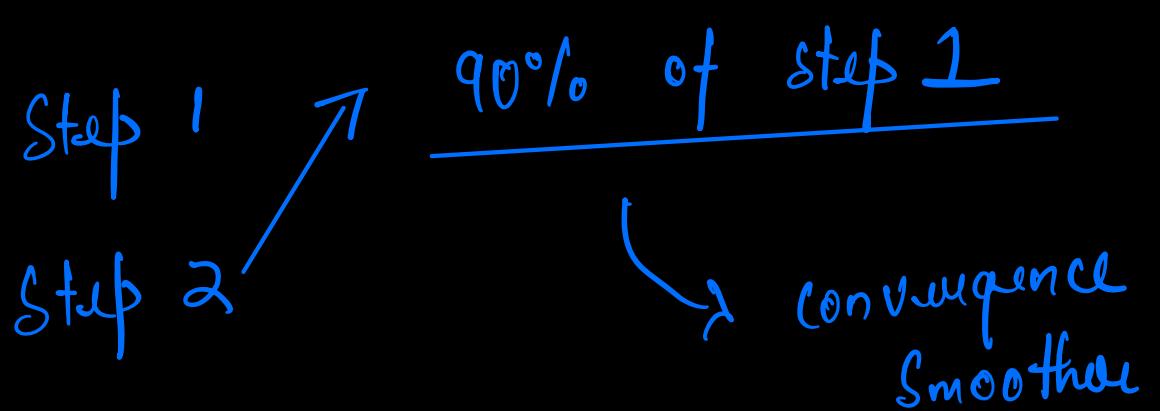
4) Step \rightarrow DL Frameworks (Keras, Pytorch)

Step = iteration

\rightarrow weight update

SGrD with Momentum





SGD with Momentum :-

$$v = \beta \times v - \underbrace{\gamma * \nabla_w L(w)}_{\text{L.R}} \rightarrow \text{gradient}$$

v = velocity

β = momentum

(i) Updating the weights

$$w = w + v$$

Momentum Co-Efficient

Initialising the, $v = 0$

0.9 over 0.99

AdaGrad

Adaptive Optimization Algorithm

(i) Adjust of your LR

Calculation

Variables :-

$$\eta = LR$$

g_t = gradient at time step t

G_t = Sum of squares of all past gradients

$$① \quad G_t = G_{t-1} + g_t^2 \quad \begin{matrix} \text{Accumulated} \\ \text{Squared} \\ \text{gradients} \end{matrix}$$

② Update Rule:-

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \cdot g_t$$

where

$$\epsilon = 1e^{-8}$$

Examples:-

$$\text{Step 1}, g_1 = 0.5$$

$$\text{Step 2} \quad g_2 = 0.5$$

$$g_1 = (0.5)^2 = 0.25$$

$$\begin{aligned} g_2 &= 0.25 + (0.5)^2 \\ &= 0.25 + 0.25 \\ &= 0.5 \end{aligned}$$

Adaptive LR :- $\frac{n}{\sqrt{G_{1t}}} = \frac{n}{\sqrt{0.5}} < \frac{n}{\sqrt{0.25}}$

Cons

Pros

i) No LR Tuning

L.R will keep shrinking towards zero value.

RMS PROP

Root Mean Square Propagation

Definition

g_t : gradient at time step

g_t
 g_{t+1}
 g_{t+2}

\vdots
 \vdots

$E[g^2]_t$:- Exponentially decaying average g_{t+100}
of past squared gradients

β :- decay rate (usually 0.9)

ϵ :- very small term

1. Update the Moving Average of Squared gradients.

$$E[g^2]_t = \beta \cdot E[g^2]_{t-1} + (1 - \beta) \cdot g_t^2$$

Moving Average :- Smooth out fluctuations

i) Simple Moving Average

ii) Exponential Moving Average

$$EMA_t = \alpha \cdot x + (1 - \alpha) EMA_{t-1}$$

α gives more weightage to recent values.

$$\begin{aligned} g &= \alpha \cdot x + (1 - \alpha) EMA_{t-1} \\ &= \alpha \cdot x + (1 - \alpha) \cdot 1 \cdot EMA_{t-1} \end{aligned}$$

Weight Update :-

$$w_{t+1} = w_t - \frac{\gamma}{\sqrt{E[g^2]_t + \epsilon}} \cdot g_t$$

Adagrad vs RMSProp

Gradient History

All past Squared

Recent Squared gradients (decay)

L.R

Keeps shrinking

Stays Adaptive

Cons :-

- 1) More complex
- 2) Tuning of β parameter (decay rate)

Adam

Adaptive Moment Estimation

Combination :- (i) SGD with momentum
(ii) Adaptive LR (RMSprop)

Maths of Adam

g_t :- gradient at time step t

m_t :- first moment ($\text{moment} = \text{mean of gradients}$)

v_t :- second moment ($\text{mean of squared gradients}$)

β_1 :- decay rate of m_t (0.9)

β_2 :- decay rate of v_t (0.999)

$$\eta = L \cdot R$$

ϵ = small num.

Step four Weight Updates

i) Compute the moment estimates

$$m_t = \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$$

$$v_t = \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$$

EWA = EMA

Exponential Weighted Average
Exponential Moving Average

2) Bias Correction (when t is small)

$$\hat{v}_t = \frac{\sqrt{t}}{1 - \beta_2^t}$$

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

This will help in mitigating the bias caused by $m_0 = 0$, $\hat{v}_0 = 0$

3) Weight Update Rule

$$w_{t+1} = w_t - \eta \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t}} + \epsilon$$

- (i) Introducing the momentum
- (ii) Scale down the LR if the
gradient have high variability