

Prerequisites
① GD
② MLP

WHAT YOU WILL STUDY IN TODAY VIDEO ?

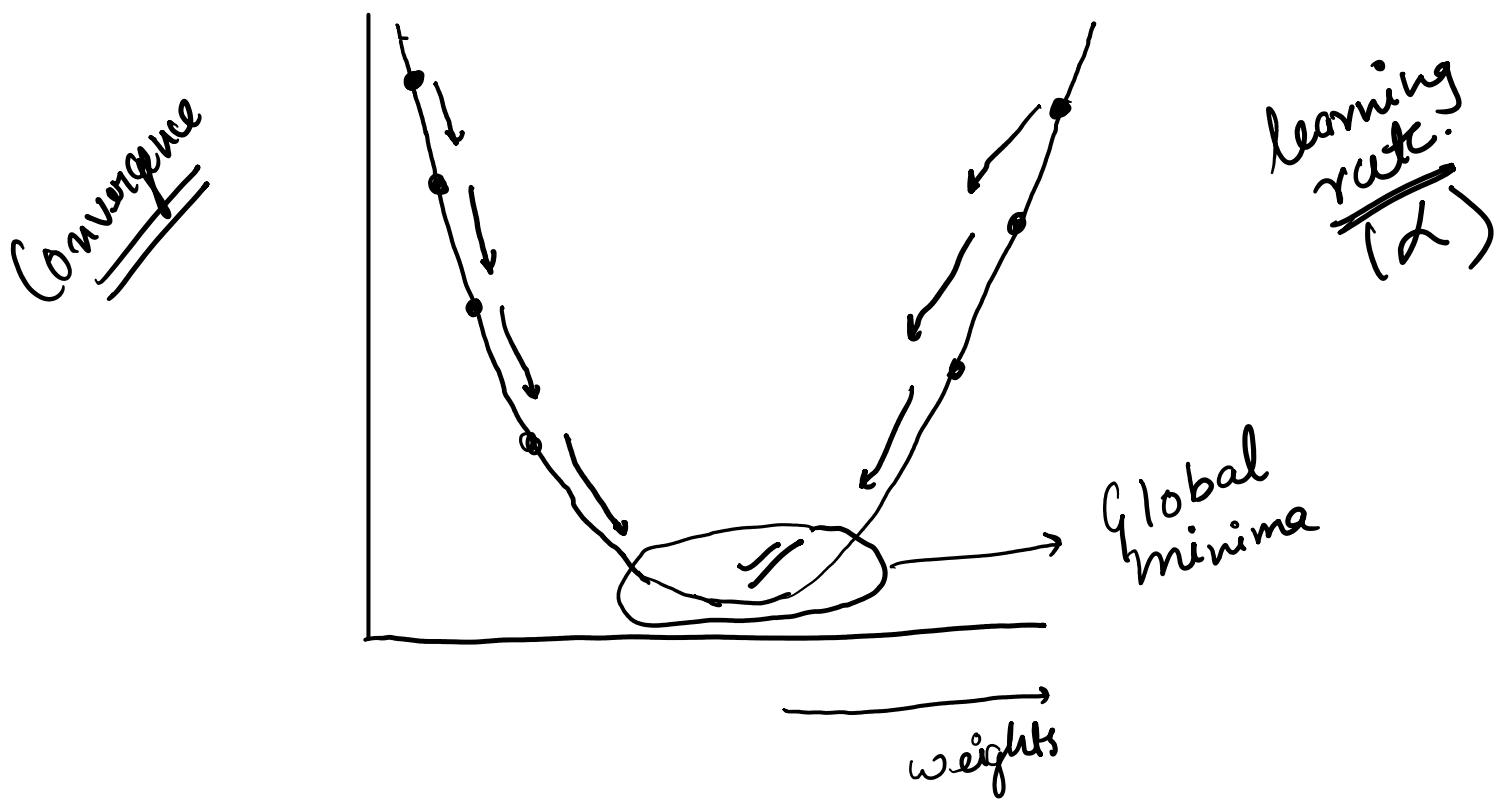


What is Optimizer?

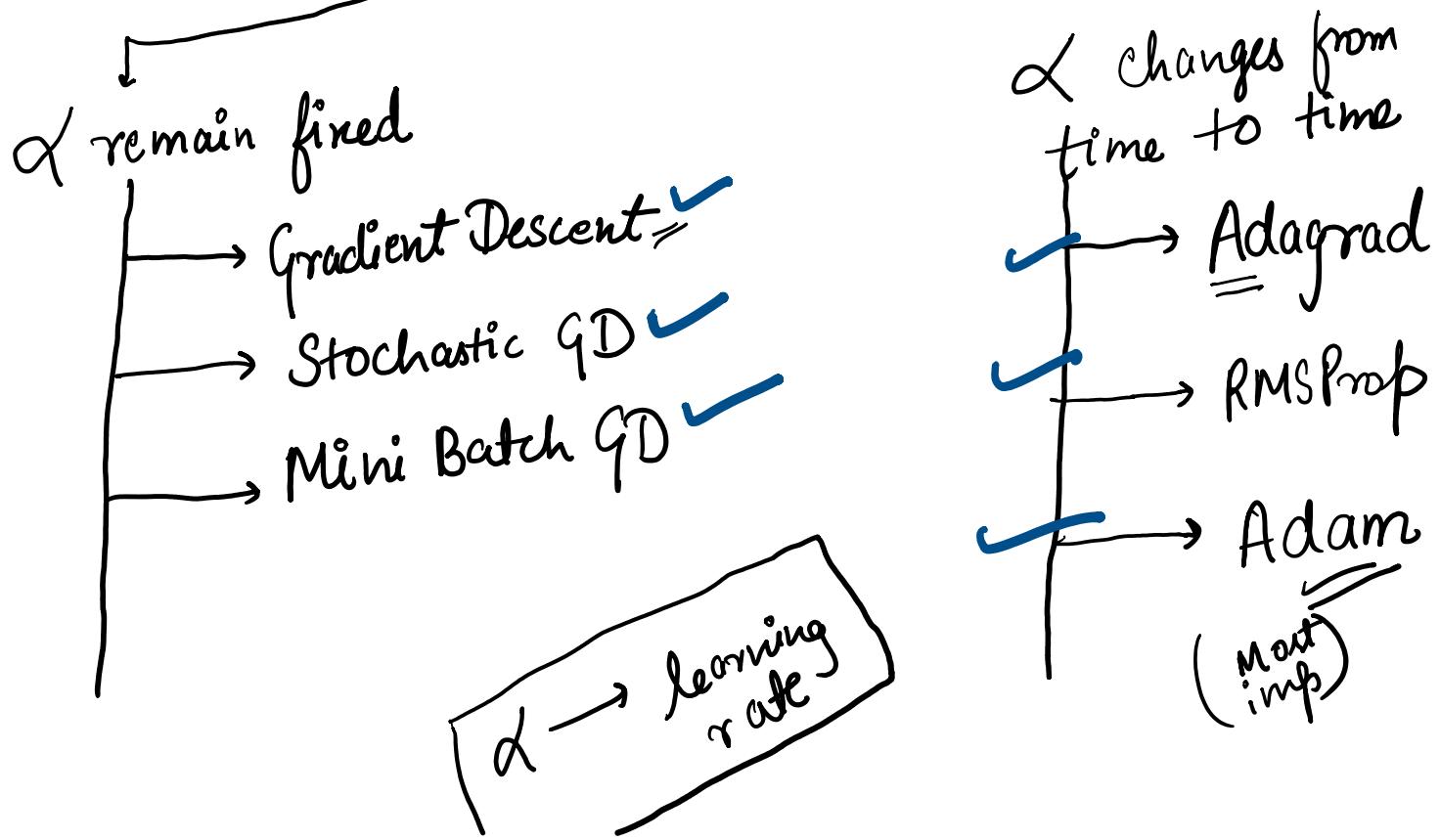
- ▶ Types of Optimizers
- ▶ Formula , Mathematical Working , Use Case of Each Optimizer
- ▶ Problems of Each Optimizer

Introduction

algorithms → weight adjust/update
during backpropagation



~~Types.~~



① GD → Gradient Descent

(Please watch GD video)

formula:

$$\omega_{\text{new}} = \omega_{\text{old}} - \alpha \cdot \frac{dL}{d\omega}$$

(loss gradient of whole dataset)

Working

Compute loss gradient for entire dataset → Weight updation

Use → Dataset is small
↳ So that they can fit in RAM

PC → 4GB RAM
Data → 500 MB

PC → 1GB RAM
Data → 8GB
X

Problem → Big dataset unable to work.

② Stochastic GD → SGD

formula:

$$\omega_{\text{new}} = \omega_{\text{old}} - \alpha \cdot \frac{dL}{d\omega}$$

loss gradient
for each row
of dataset

Working

Compute loss gd
for 1st row

Weight update

2nd row
loss gd

3rd
row

Use → Good for large data
 ↳ No load on RAM

Problem

① Time ↑

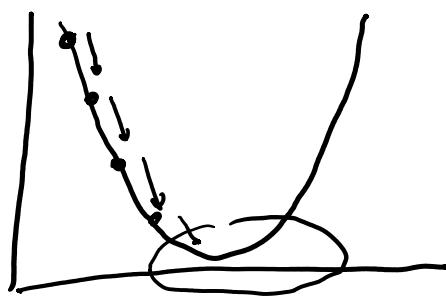
② Zig Zag behaviour
 ↳ Unstable convergence

③ Overshooting Problem

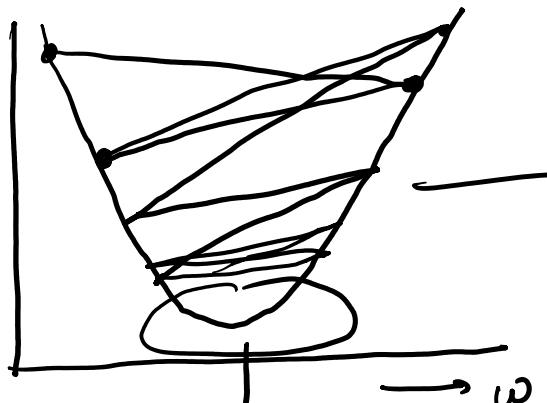
GD → Slowly/steadily towards global minima

GD \rightarrow slowly/converging

to minima



SQD \rightarrow



Zigzag behaviour
(slower convergence)

Optimal weight values

③ Mini Batch GD

$m \rightarrow$ batch size.

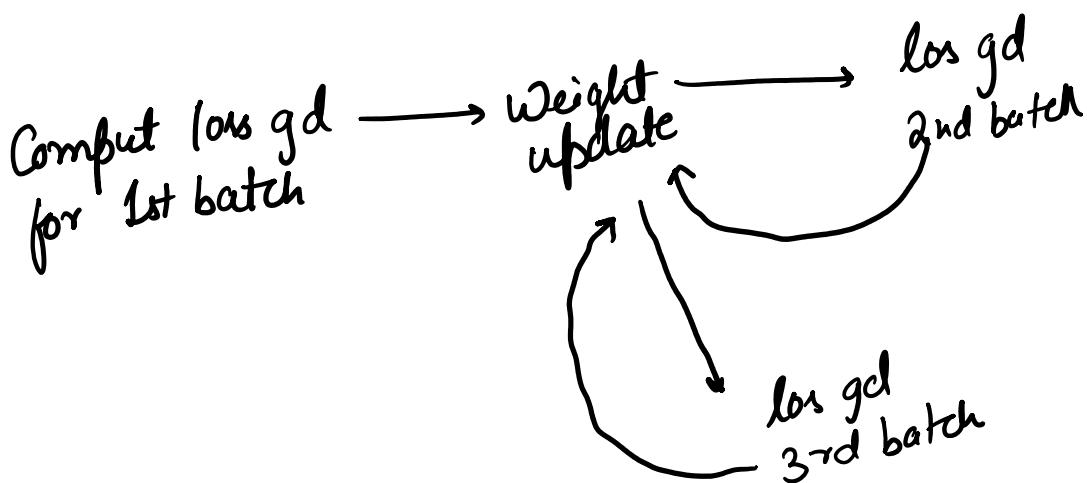
formula

$$\omega_{\text{new}} = \omega_{\text{old}} - \alpha \cdot \frac{1}{m} \sum_{i=1}^m \frac{dL}{d\omega}$$

$L \rightarrow$ loss gradient
for a batch

Working

Data \rightarrow 150 rows
Batch \rightarrow 3 \rightarrow 50 rows each



Use

Most commonly used GD. of all

Efficient
than SGD

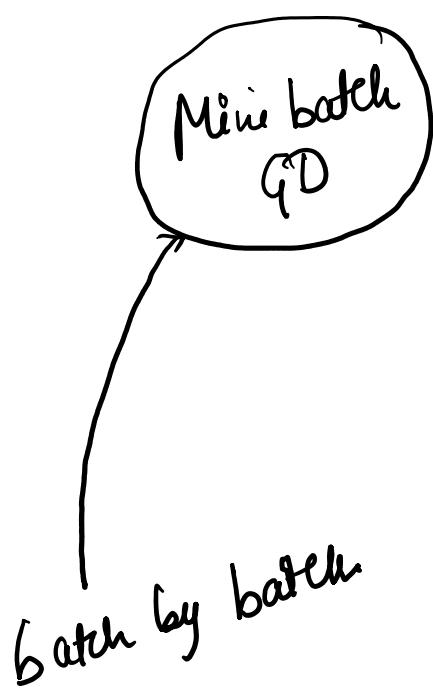
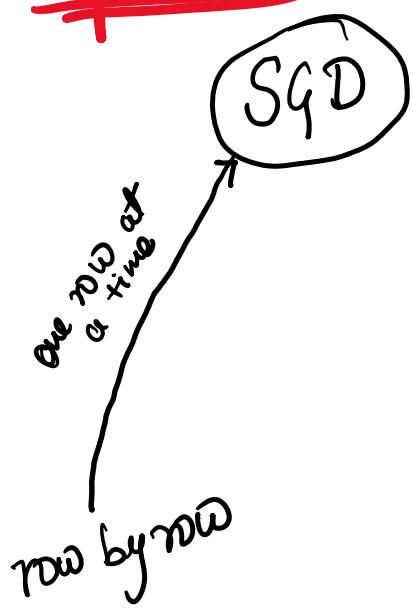
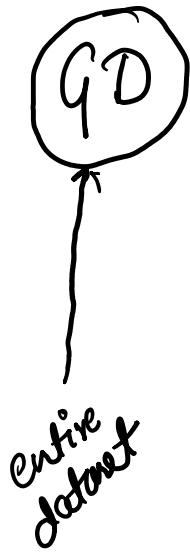
fast than
SGD

Stable
Convergence

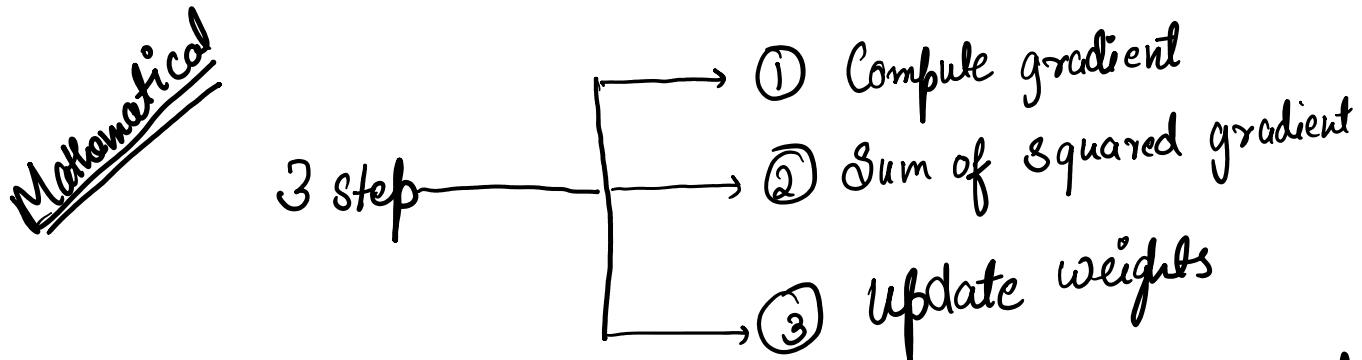
Problem → find best value of $\underline{\underline{m}}$.

Converge

Comparison



① ~~Adagrad~~ Adaptive gradient



eg Given,
 $\omega = 5$
 $\alpha = 0.1$

$\epsilon = 10^{-8}$ → very small
 (prevent division by 0)

loss fn = ω^2
 Derivative of loss fn $\Rightarrow g_t = 2\omega$

① Iteration 1 ($t=1$)

(a) $g_1 = 2\omega = 2 \times 5 = \underline{\underline{10}}$

(b) $G_1 = (10)^2 = \underline{\underline{100}}$ learning rate for this iteration

(c) $\omega_{new} = \omega_{old} - \frac{\alpha}{\sqrt{G_t} + \epsilon} \cdot g_t$

$$= 5 - \frac{0.1}{\sqrt{100} + 10^{-8}} \cdot 10$$

$$= 5 \quad \sqrt{100} + 10^{\circ}$$

$$\underline{\underline{\omega}}_{\text{new}} = \underline{\underline{4.9}}$$

② Iteration 2 ($t=2$)

$$\textcircled{1} \quad g_2 = 2 \times \omega = 2 \times 4.9 = \underline{\underline{9.8}}$$

$$\textcircled{2} \quad G_2 = 100 + (9.8)^2 = \underline{\underline{196.04}}$$

$$\textcircled{3} \quad \underline{\underline{\omega}}_{\text{new}} = 4.9 - \frac{0.1}{\sqrt{196.04}} \cdot 9.8$$

$$\underline{\underline{\omega}}_{\text{new}} = \underline{\underline{4.83}}$$

Do for iteration 3, 4 & so on

Working

① α decrease over time

Iteration 1.

$$\alpha = \frac{0.1}{10}$$

$$= \underline{\underline{0.01}}$$

Iteration 2

$$\alpha = \frac{0.1}{14}$$

$$= \underline{\underline{0.00714}}$$

decrease 

~~—~~ ~~=~~ decrease

- ② large gradient = smaller weight update
Small gradient = larger weight update

$$\Rightarrow \text{gradient}_{(g_t)} = \text{large}$$

$$g_t = \text{large}$$

g_t is in denominator

$$\frac{0.1}{100} \rightarrow \text{small}$$

Overall learning rate is also small

weight update is also small

use \longrightarrow Adjusting learning rate \equiv
 \hookrightarrow No need to define lr

use → "why" ↳ No need to define "

→ NLP projects

→ Sparse data

[1, 0, 0, 0, 0, 0]

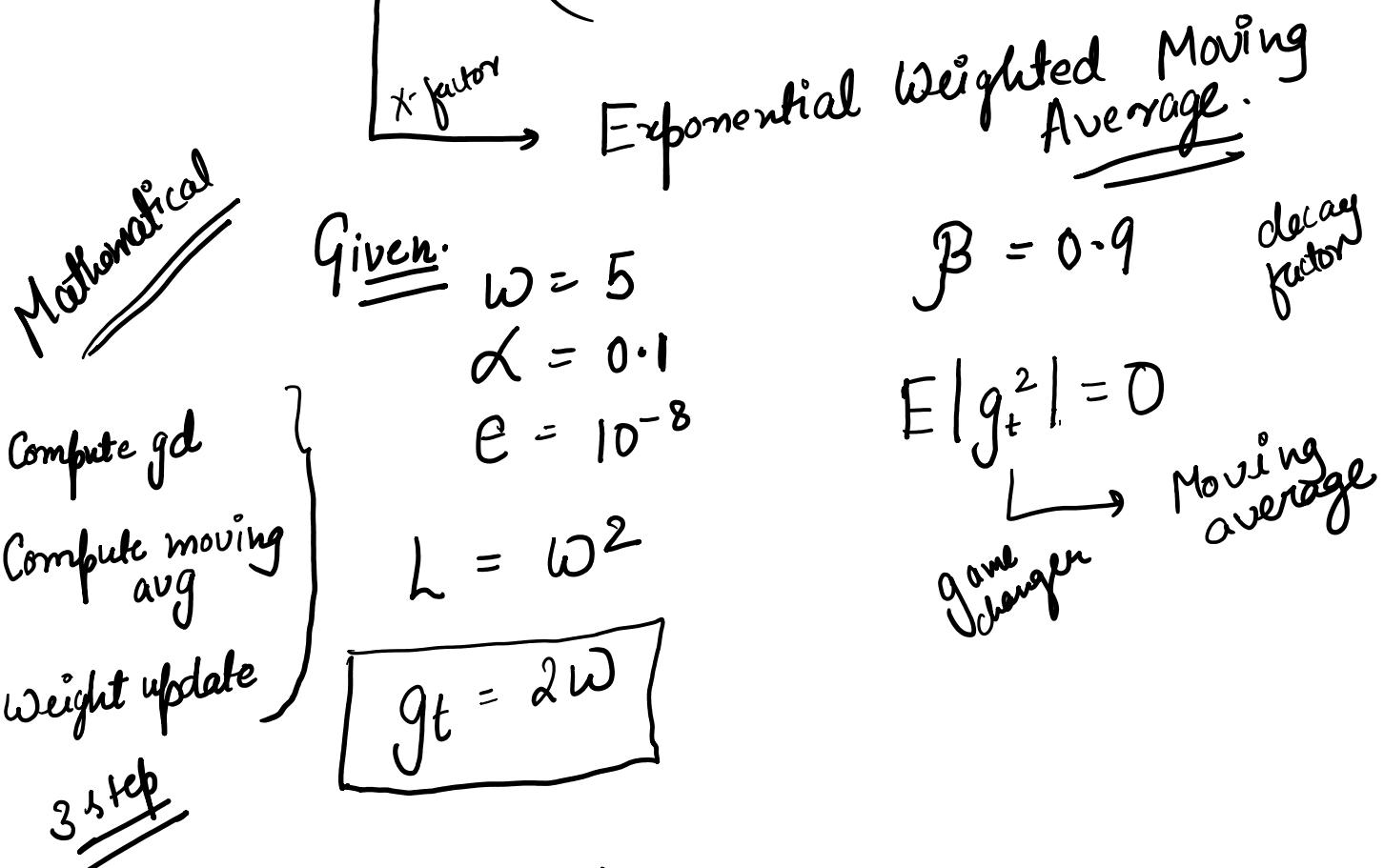
problem → Deep learning ~~X~~ ↳ Not used

↳ learning rate decay = main

↳ High Computation Cost. ==

fin → RMSProp ==
Adam

⑤ RMSProp



① Iteration 1 ($t=1$)

$$(a) g_1 = 2 \times \omega = 2 \times 5 = \underline{\underline{10}}$$

$$(b) E|g_1^2| = (\beta \times \text{Previous moving avg}) + (\alpha \times g_1^2)$$

$$= (0.9 \times 0) + (0.1 \times 100)$$

$$= \underline{\underline{10}}$$

$$(c) \underline{\underline{\omega_{\text{new}}}} = \underline{\underline{\omega_{\text{old}}}} - \frac{\alpha}{\sqrt{E|g_t^2| + \epsilon}} \times g_t$$

$$\tilde{\omega}_{\text{new}} = \omega_{\text{old}} - \frac{0.1}{\sqrt{10}} \times 10^{-8} \times 10$$

$$= 5 - \frac{0.1}{\sqrt{10}} + 10^{-8} \times 10$$

$$= 5 - \frac{0.1}{3.162} \times 10$$

$$\omega_{\text{new}} = 1.838$$

② Iteration 2 ($t=2$)

$$(a) g_2 = 2 \times \omega = 2 \times 1.838 = \underline{\underline{3.676}}$$

$$(b) Elg_2^2 = (0.9 \times 10) + (0.1 \times 3.676^2)$$

$$= 9 + 1.352$$

$$= \underline{\underline{10.352}}$$

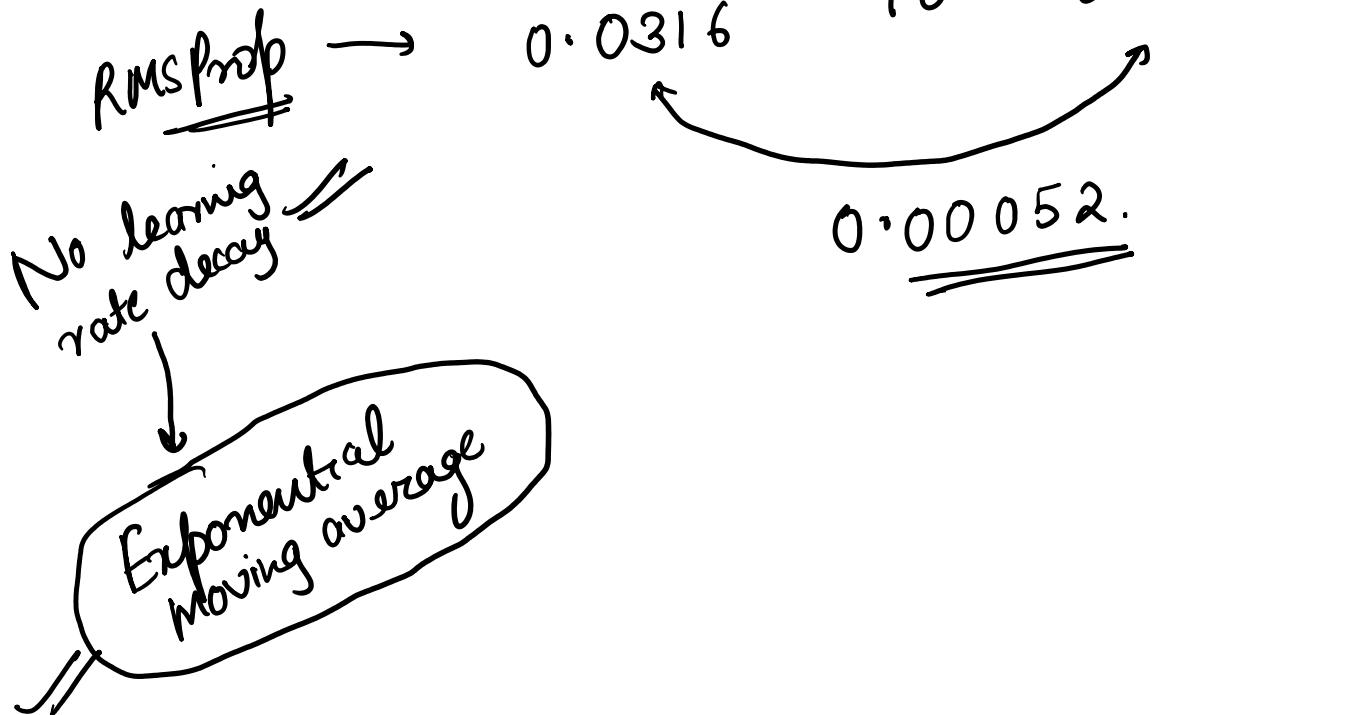
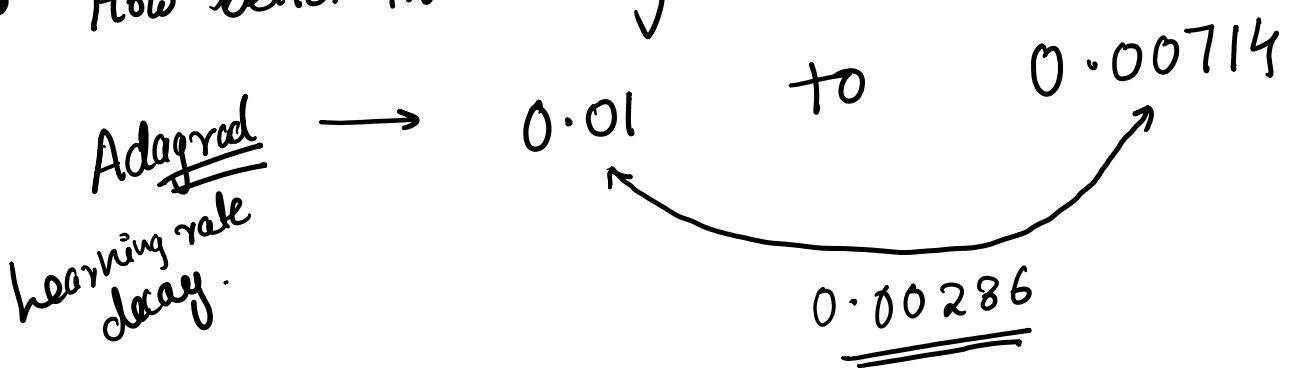
$$(c) \omega_{\text{new}} = 1.838 - \frac{0.1}{\sqrt{10.352} + 10^{-8}} \times 3.676$$

$$\omega_{\text{new}} = 1.724$$

on on

So on

- How better than Adagrad?



Use → Deep learning usable
Most popular till 2014
(Adam)

e.g. Inception
(Google)

e.g. Inception
(Google)

Problem:

There is no momentum in RMSProp
(technique in which optimizer remembers past gradients)

Adam has momentum

⑥

Adam

Combination

Momentum

+

RMSProp

Mathematical

Given: $\omega = 5$
 $\alpha = 0.1$

decay
factor

$$\beta_1 = 0.9$$

$$L = \omega^2$$
$$g_t = 2\omega$$
$$\epsilon = 10^{-8}$$

To keep track
of previous
gradients

$$\beta_2 = 0.999$$

① Compute gradient

$$g_1 = 2\omega$$
$$= 2 \times 5$$
$$= \underline{\underline{10.}}$$

② Compute \underline{m} (first moment estimate)

$$m_1 = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$m_1 = (0.9 \times 0) + (1 - 0.9) \times 10$$

$$\underline{m_1 = 1}$$

$$m_1 = 1$$

③ Compute V (second moment estimate)

$$V_1 = \beta_2 V_{t-1} + (1 - \beta_2) g_t^2$$

$$V_1 = (0.999 \times 0) + (1 - 0.999) \times 100$$

$$V_1 = 0.1$$

④ Bias Correction / Normalization
 (\hat{m}, \hat{v})

$$\hat{m} = \frac{m_1}{1 - \beta_1} = \frac{1}{1 - 0.9} = \underline{\underline{10}} \rightarrow \hat{m}$$

$$\hat{v} = \frac{V_1}{1 - \beta_2} = \frac{0.1}{1 - 0.999} = \underline{\underline{100}} \rightarrow \hat{v}$$

⑤ Weight updation

$$\omega_{\text{new}} = \omega_{\text{old}} - \frac{\alpha \hat{m}}{\sqrt{\hat{v}} + \epsilon} \cdot \times 10$$

$$= 5 - \frac{0.1 \times 10}{\sqrt{100} + 10^{-8}}$$

$\boxed{\omega_{\text{new}} = 4.9}$

//
Do for iteration 2 so on... //

Summary

① Compute gradient //

② Compute $m \equiv$ → weighted exponential moving avg of past gradient

③ Compute $v \equiv$ → moving avg of squared gradients //

RMS Prop
Learning rate adjustment

{ larger gradient → smaller update
smaller gradient → larger update }

④ Bias \hat{m}, \hat{v} → Normalization

① Bias Correction $\rightarrow \hat{m}, \hat{v} \text{ No. } //$

Helps in avoiding slow convergence //

⑤ Weight updation //

~~Note:~~ Variants of Adam

① AdamW \rightarrow Adam with weight decay
(most used) //

② AdaBound

③ RAdam

Use \rightarrow Perfect optimizer //