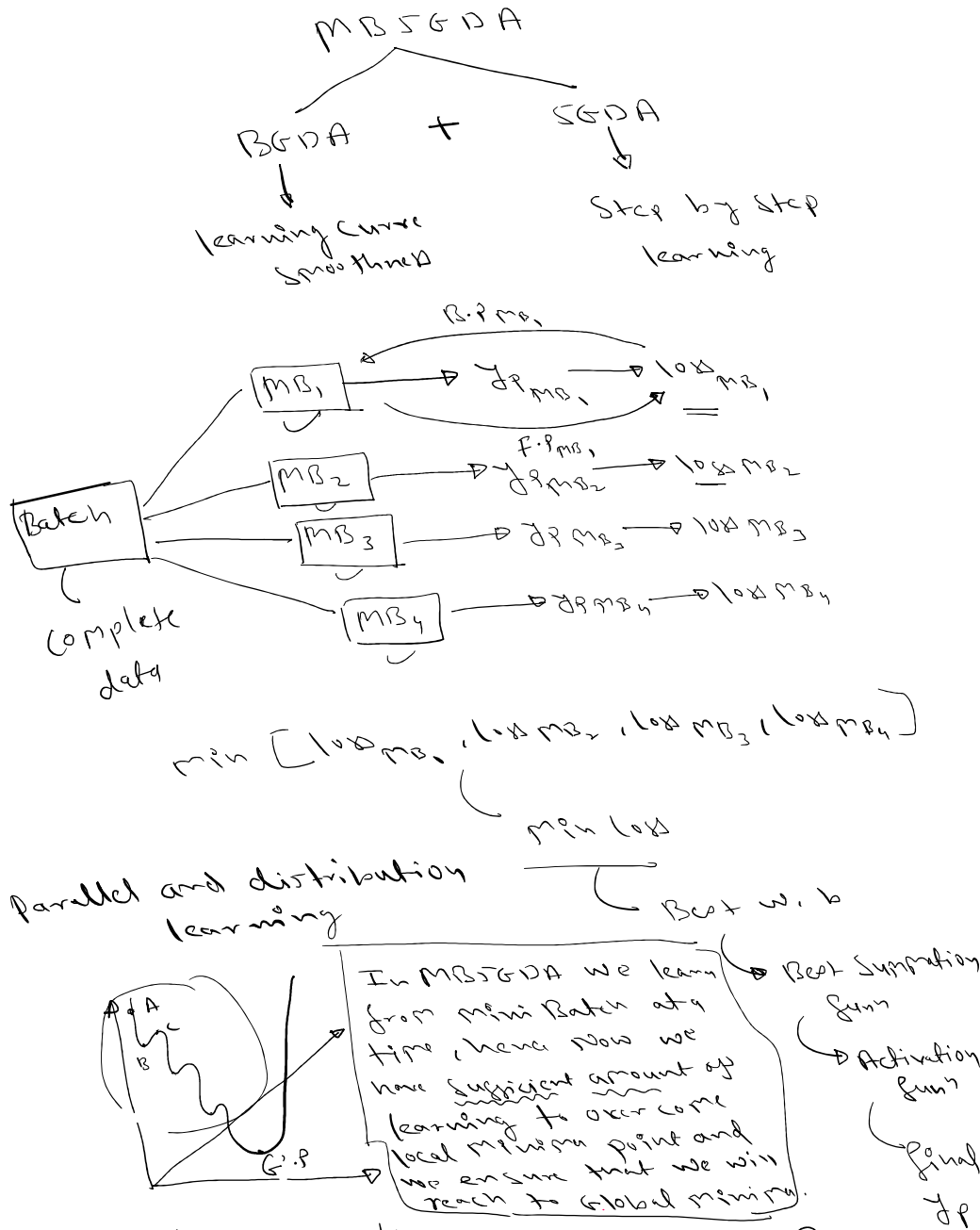Deep learning Day - 4

(3) Mini Batch Stochastic gradient Decent algorithm

MBSGDA

BGDA    +    SGDA

learning Curve          Step by step
Smoothness               learning

B.P MB₁

[MB₁] → $\frac{\partial P}{\partial MB_1}$ → loss MB₁ =

F.P MB₁
[MB₂] → $\frac{\partial P}{\partial MB_2}$ → loss MB₂

Batch

[MB₃] → $\frac{\partial P}{\partial MB_3}$ → loss MB₃

complete
data

[MB₄] → $\frac{\partial P}{\partial MB_4}$ → loss MB₄

min [ loss MB₁, loss MB₂, loss MB₃, loss MB₄ ]

min loss

Parallel and distribution
learning

Best w, b

In MBSGDA we learn          Best Summation
from mini Batch at a         funn
time, hence now we
have sufficient amount of    Activation
learning to overcome         funn
local minima point and
we ensure that we will       final
reach to Global minima.      $\frac{\partial P}{\partial}$

* Random sampling

Batch₁ (250)

Batch₂ (250)

Batch₃ (250)

Batch₄ (250)

Complete
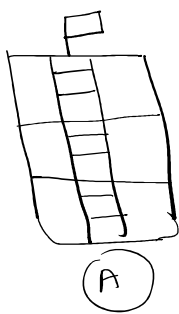data

(1000)

1000 → 10 Mini Batches — 1 Mini Batch

100

* Advantage of MBSGDA :-

(1) Compare to BGDA , MBSGDA required less resources and time to achieve convergence.

(2) Since we are learning from Mini Batch at a time, we have sufficient amount of learning to over come local minima hence we ensure that we will reach to Global minima.
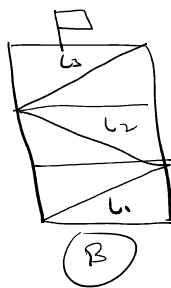
* Drawback of MBSGDA :-
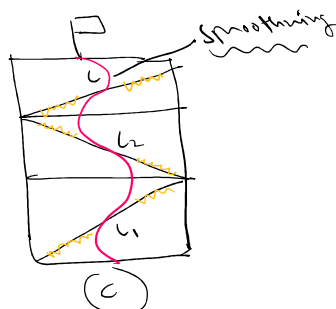
(1) Noise is still present in learning curve.
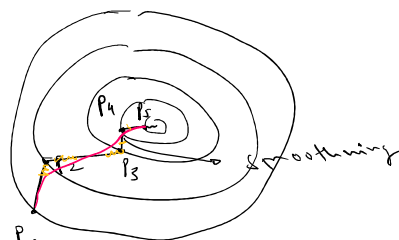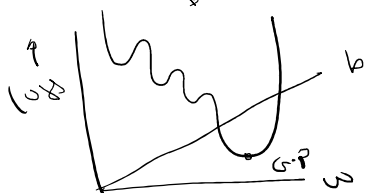
* Why we don't want noise in our learning curve ?

A

time ↓
effort ↑
↓

B

time ↑
effort ↓

C

Smoothing

time ↓
effort ↓

Noise → time ↑ → time complexity

Noise → time ↑ ⟶ P. ⟶ time complexity
        ⟶ effort ↑ ⟶ ⟶ resources.

④ Mini Batch stochastic gradient Decent algorithm
   with momentum.

                        smoothing

$$w_{new} = w_{old} - d \boxed{\dfrac{\partial L}{\partial w}}$$

$$b_{new} = b_{old} - d \boxed{\dfrac{\partial L}{\partial b}} \longrightarrow \text{momentum}$$

(Gradient) ⟶ Potential energy

$\dfrac{\partial L}{\partial w} \longrightarrow$ Changes in loss with respective weight.

$\dfrac{\partial L}{\partial w} =$

$\dfrac{\partial L}{\partial b}$

loss

b

G.P

w

learning rate

Noise

A    B

C

G.P

$$w_{new} = w_{old} - d \times \underline{vdw}$$

$$b_{new} = b_{old} - d \times \underline{vdb}$$

$vdw$ or $vdb \longrightarrow$ Velocity Component

$$vdw_t = \beta \times vdw_{t-1} + (1-\beta) \dfrac{\partial L}{\partial w}$$

$$vdb_t = \beta \times vdb_{t-1} + (1-\beta) \dfrac{\partial L}{\partial b}$$

where

$t \longrightarrow$ no. of hidden layer

Range  $\longrightarrow \beta \longrightarrow$ Smoothing parameter
$(\longrightarrow 0 \text{ to } 1)$          (0 to 1)

$vdw$ or $vdb \longrightarrow$ velocity component

Case I $\Longrightarrow \beta = 0$  ✗

$$vdw_t = \beta \times vdw_{t-1} + (1-\beta) \dfrac{\partial L}{\partial w}$$

$$= 0 \times vdw_{t-1} + 1 \times \dfrac{\partial L}{\partial w}$$

$$vdw = \dfrac{\partial L}{\partial w} \longrightarrow \text{No smoothing}$$

Case II $\rightarrow$ $\beta = 1$ ✗

$$Vdw_t = \beta \times Vdw_{t-1} + \underbrace{(1-\beta)\frac{dL}{dw}}_{0}$$

$$\boxed{Vdw_t = Vdw_{t-1}} \rightarrow \text{Full smoothing}$$

① ② ③
1 2 3

Case III $\rightarrow$ $\beta = 0.98$

$$Vdw_t = \beta \times Vdw_{t-1} + (1-\beta)\frac{dL}{dw}$$

$$\boxed{Vdw_t} = \underbrace{\boxed{0.98\, Vdw_{t-1}}}_{\underset{\text{Magnitude}}{I}} + \underbrace{\boxed{0.02\frac{dL}{dw}}}_{\underset{\text{Direction}}{II}} \qquad \boxed{G \cdot P}$$
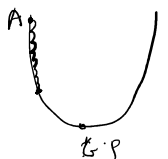
* Best optimization function for weight
  and bias $\rightarrow$ <u>MBSGDA with momentum</u>

* Why we need to update learning rate
  dynamically in Neural Networks ?



$\boxed{\text{Learning Rate}}$ —— <u>Constant</u>

$\boxed{\text{Learning Step}}$

$\boxed{\underline{\text{Rate}}}$

Case I $\rightarrow$ If learning rate is too small ?
                                          (0.001)

$\rightarrow$ training time increase
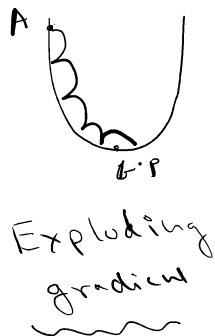              exponentially.

$\rightarrow$ we face vanishing gradient
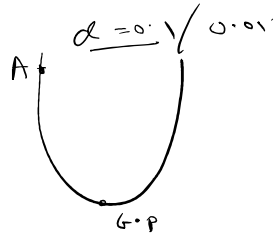              issue.

<u>Vanishing</u>
<u>Gradient</u>  $\rightarrow$ we never achieve global mining

Case II $\rightarrow$ If learning rate is too high ?
                                          (0.1)

Case II $\longrightarrow$ If learning rate is too high!.

$\longrightarrow$ we face over shooting issue

$\longrightarrow$ due to over shooting we face exploding gradient issue.

A

Exploding gradient

A $\quad \alpha = 0.1 \; / \; 0.01$

G.P

$\alpha = 0.1 \longrightarrow$ too high $\downarrow$

$\alpha = 0.01 \longrightarrow$ small $\uparrow$

learning rate (dynamically update.

(1) Ada grad. (Adaptive gradient)

(2) Ada Delta
$\quad \hookrightarrow$ it is also called as RMS prop (Root mean squared propagation

(1) Ada grad

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

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

$$\boxed{\alpha_{new} = \frac{\alpha_{old}}{\sqrt{\eta + \epsilon}}}$$

$\longrightarrow$ where $\epsilon \simeq e^{-6}$ (small positive number)

we use $\epsilon$ in equation to avoid zero divisional error.

$\frac{\partial L}{\partial w}, \quad \frac{\partial L}{\partial w_2} \quad \frac{\partial L}{\partial w_3}$

$$\eta = \sum_{i=1}^{5} \left( \frac{\partial L}{\partial w_i} \right)^2$$

A $\eta = \left( \frac{\partial L}{\partial w_1} \right)^2 + \left( \frac{\partial L}{\partial w_2} \right)^2 + \left( \frac{\partial L}{\partial w_3} \right)^2$

$\eta_{L_1} \quad h_{L_2} \quad h_{L_3}$

$\eta \uparrow \longrightarrow$ $\alpha_{new}$ ↓

↘ vanishing gradient

\* Drawback of Ada grad

→ We cannot use Ada grad in case of deep Neural Networks.

② Ada Delta (RMS Prop)

(Root mean squared propagation)

↑ $\eta \longrightarrow \alpha_{new}$ ↓

$\eta \longrightarrow$ Replace $\longrightarrow S_{dwt}$

$$w_{new} = w_{old} - \alpha_{new} \frac{dc}{dw}$$

$$b_{new} = b_{old} - \alpha_{new} \frac{dc}{db}$$

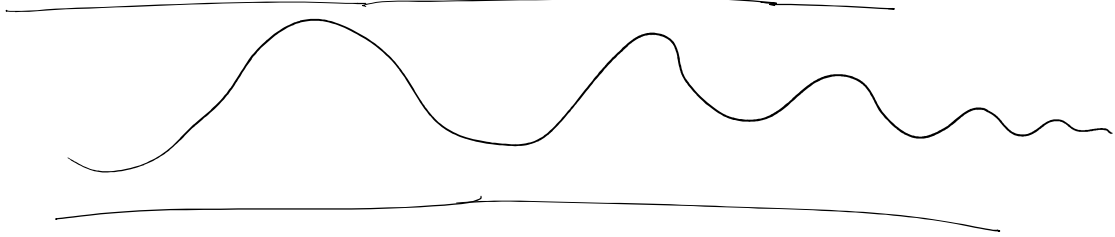$$\alpha_{new} = \frac{\alpha_{old}}{\sqrt{S_{dwt} + \epsilon}}$$

$S_{dwt}$ = velocity component

$$V_{dwt} = \beta \times V_{dwt-1} + (1-\beta) \frac{dc}{dw}$$

$$\boxed{S_{dwt} = \beta \times S_{dwt-1} + (1-\beta) \left(\frac{dc}{dw}\right)^2}$$

$$Sdw_t = \beta \times Sdw_{t-1} + (1-\beta)\left(\frac{\partial c}{\partial w}\right)$$

$\beta \longrightarrow$ Smoothing parameter

$\xi$

$\longrightarrow$ Best optimization function for learning rate $\longrightarrow$ Ada Delta

( RMS prop

* Best optimization function for weight and bias ?
  $\longrightarrow$ MBSGDA with momentum

+

* Best optimization function for learning rate ?
  $\longrightarrow$ Ada Delta (RMS prop)

$=$ ADAM