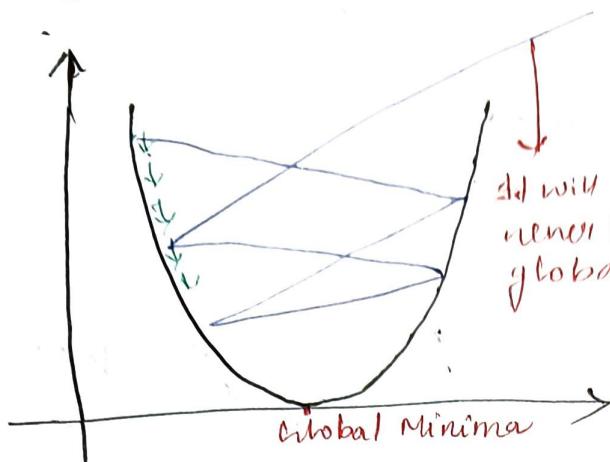


15 - Jan - 2023

Exploding Gradient Problem



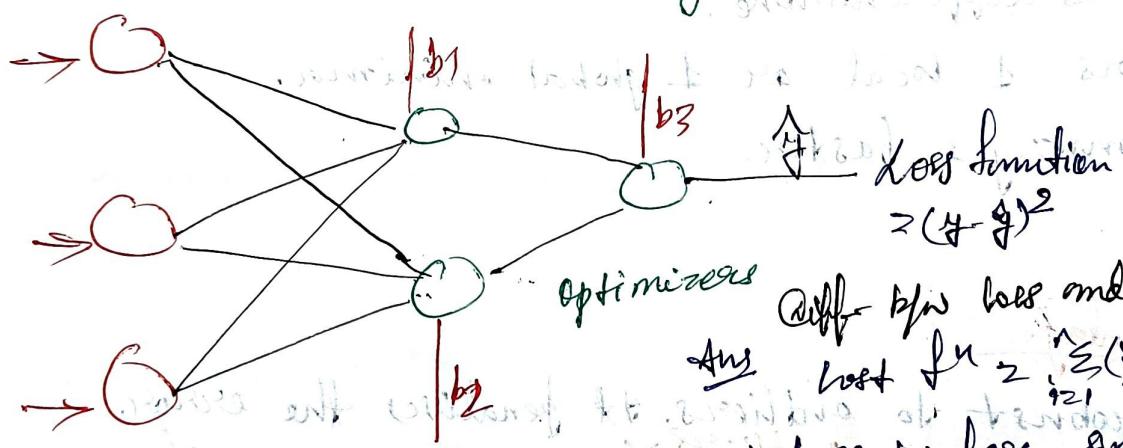
→ weight updation in such a way that it is jumping here and there.

it will never come to global minima.

Q Why this jumping problem happens?
Ans Because of ~~higher~~ weight initialization.

Forward Propagation

Binary classification



Diff b/w loss and cost fn.

$$\text{Ans} \quad \text{cost fn} = \sum_{i=1}^n (H_i - g_i)^2$$

- loss is for single data point
→ loss is for batch data points.

Loss Functions

ANN

- ① Classification
- ② Regression

Regression

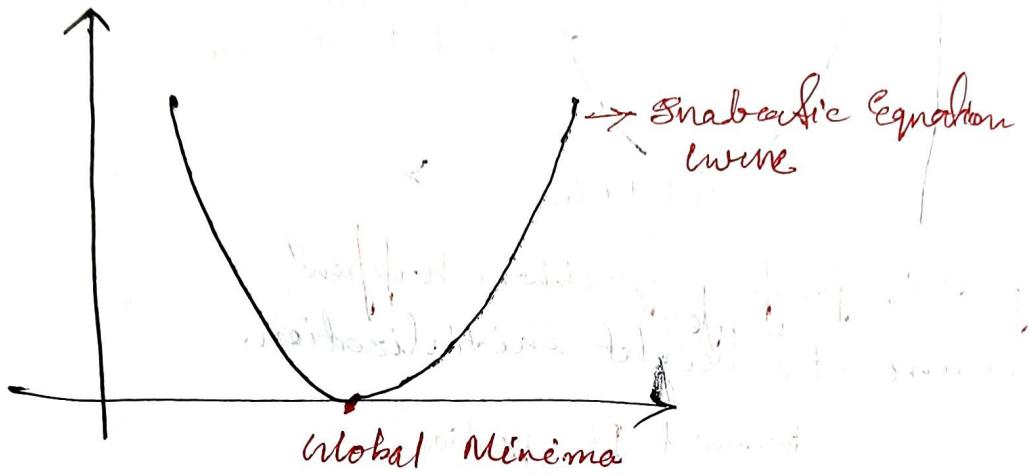
- ① MSE (Mean Squared Error)
- ② MAE (Mean absolute error)
- ③ Huber loss
- ④ RMSE (Root mean Squared Error)

① MSB (Mean Squared Error)

$$\text{Loss} J^L = (y - \hat{y})^2$$

$$\text{Cost} J^M = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

\downarrow
Quadratic Equation

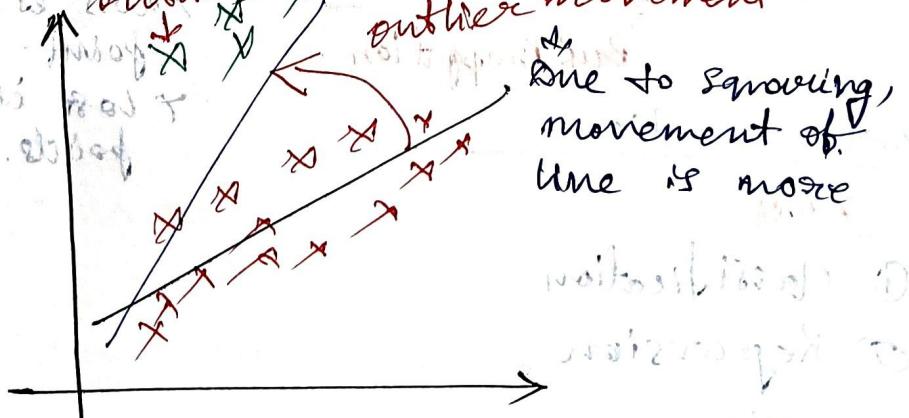


Advantage

- ① MSB is differentiable.
- ② It has 1 local & 1 global minima.
- ③ It converges faster.

Disadvantage

- ① Not robust to outliers. It penalizes the error.



Outlier is a single point that is far away from other points.
Outlier is an outlier because it is far away from other points.
Outlier is an outlier because it is far away from other points.

② MAE (Mean Absolute Error)

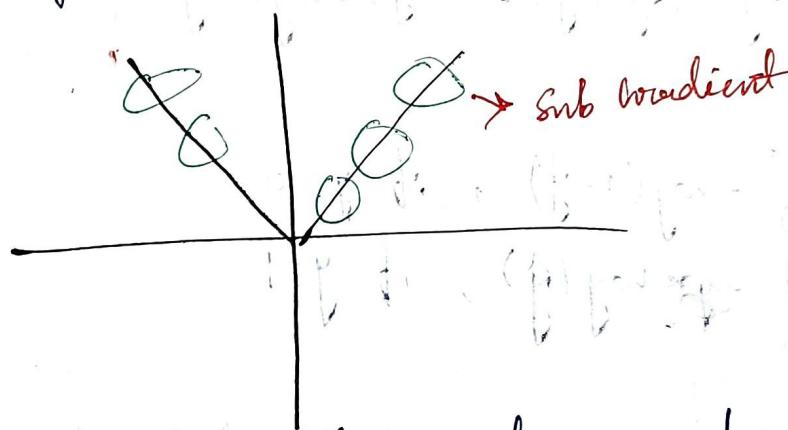
$$\text{loss } f^M = |y - \hat{y}| \quad \text{cost } f^M = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Advantage

- ① Robust to outliers, because movement of one is slow. We are not squaring here.

Disadvantage

- ② Convergence is slow.



- ③ It is not differentiable at every point.

③ Huber loss

> combination of MSE and MAE

$$\text{cost } f^H = \begin{cases} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

④ RMS E

$$\text{RMS E} = \sqrt{\text{MSE}}$$

$$(\text{RMS E})^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = (\mathbf{y}, \mathbf{y}) - (\mathbf{y}, \mathbf{\hat{y}}) + (\mathbf{\hat{y}}, \mathbf{\hat{y}})$$

Loss or Cost fn for classification Problem

CROSS ENTROPY

- ① BINARY CROSS ENTROPY (Binary class)
- ② CATEGORICAL CROSS ENTROPY (Multi class)
- ③ SPARSE CATEGORICAL ENTROPY (Multi class)

① Binary class Entropy (Binary classification)

$$\text{loss fn} = -y \log(\hat{y}) - (1-y) \log(1-\hat{y}) \quad \text{← log loss}$$

$$\text{loss fn} = \begin{cases} -\log(1-\hat{y}), & \text{if } y=0 \\ -\log(\hat{y}), & \text{if } y=1 \end{cases}$$

$$\hat{y} = \frac{1}{1+e^{-z}}$$

$z = \text{activation fn}(\text{input} \star \text{weight})$

↓
Sigmoid Activation Function

② Category Categorical Cross Entropy (Multi-class classification)

	f_1	f_2	f_3	O/P	Label	$j=1$	$j=2$	$j=3$	\dots	$c=3$ no. of category
$i=1$	2	3	4	hood	1	1	0	0		
$i=2$	5	6	7	bad	0	1	0	0		
$i=3$	8	9	10	Neutral	0	0	0	1		

> Categorical cross entropy will apply one-hot encoding.

$$\text{loss fn}(x_i, y_i) = - \sum_{j=1}^c y_{ij} \ln(\hat{y}_{ij})$$

Actual values $\Rightarrow \hat{y}_{ij} [y_1 \ y_{12} \ y_{13} \dots \ y_{1c}]$

$[y_2 \ y_{22} \ y_{23} \dots \ y_{2c}]$

$y_{ij} = \begin{cases} 1, & \text{if the element is in the class} \\ 0, & \text{otherwise} \end{cases}$

Prediction $\Rightarrow \hat{y}_{ij} \Rightarrow \text{Softmax Activation} \Rightarrow \text{soft}(z)$

$$= \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

$$= \frac{e^{z_i}}{e^1 + e^2 + e^3 + \dots + e^k}$$

⑧ SPARSE Categorical Entropy

$$\text{opp} = [0.2, 0.3, 0.5]$$

\downarrow highest
2nd index

\Rightarrow It will give the highest value index as output.

ex

$$[0.2, 0.3, 0.1, 0.2, 0.2]$$

\downarrow highest 1st index.

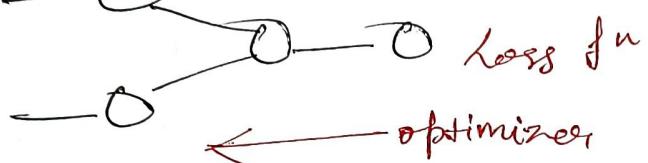
\Rightarrow losing other information.

Right combination of HL, O/P layer, Problem Statement & loss fn

<u>Hidden layers</u>	<u>O/P layer</u>	<u>Problem Statement</u>	<u>loss function</u>
① ReLU and its variants	Sigmoid	Binary classification	Binary cross entropy
② ReLU and its variants	Softmax	Multi class	Categorical cross entropy
③ ReLU or its variants	Linear	Regression	MSLE, MAE, Huber loss, RMSE

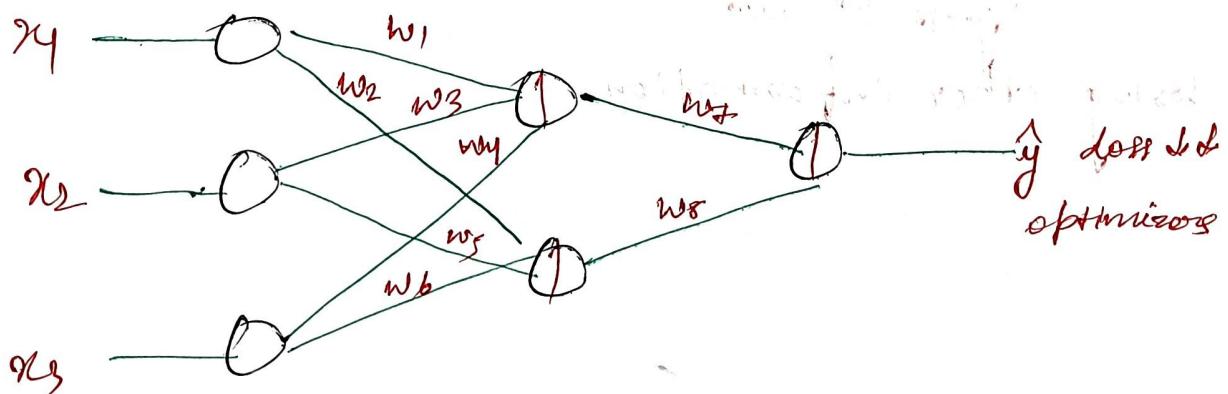
Optimizers

- To update weights.



- ① Gradient Descent
- ② SGD (Stochastic Gradient Descent)
- ③ SGD with momentum
- ④ Mini batch SGD
- ⑤ Adagrad and RMSprop
- ⑥ Adam optimizers
- ⑦ Adadelta

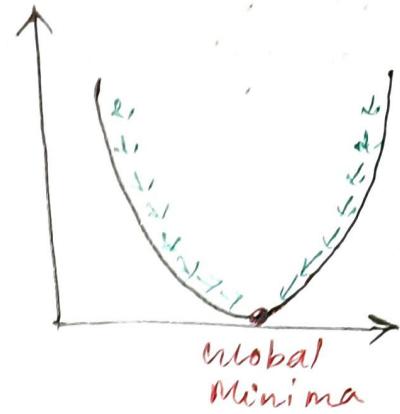
① Gradient Descent optimizers (all the records are taken together in BT FIP C BFO)



weight updation formula

$$W_{\text{new}} = W_{\text{old}} - \gamma \left(\frac{\partial \text{loss}}{\partial W_{\text{old}}} \right)$$

slope
chain rule of
Backwarding



MSE

$$\text{loss}_{f_n} = (y - \hat{y})^2$$

$$\text{cost } f_n = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

Epochs vs Iterations

Data points = 1000

Epochs 1 {
 ↓
 1000 data points
 ↓
 P.P
 ↓
 iterations
 ↓
 weight will
 get updated

Epochs 2 {
 ↓
 weight will get
 updated

→ This will continue till cost f_n is minimized.

Example

cost f_n decreased till 98th ~~iterations~~ &
 after that it is not decreasing, then stop
 at that point.

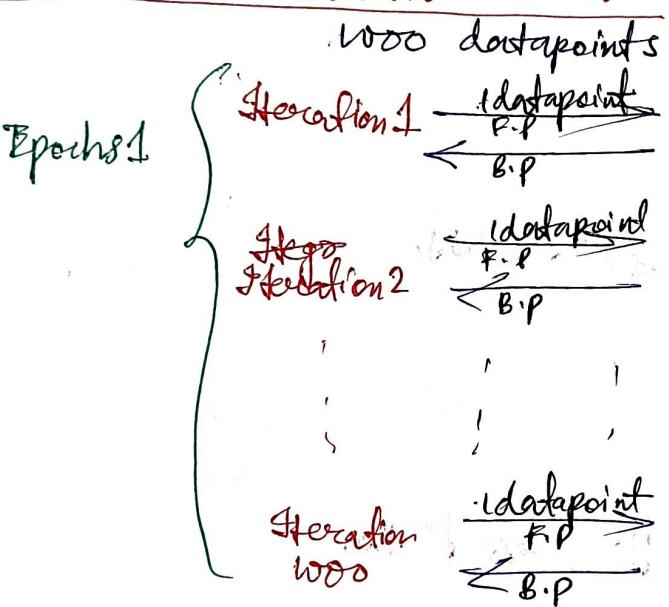
Advantage

- ① Convergence will happen.

Disadvantage

- ② It will require high RAM and CPU. It is resource intensive.

Stochastic Gradient Descent



Epochs 1000

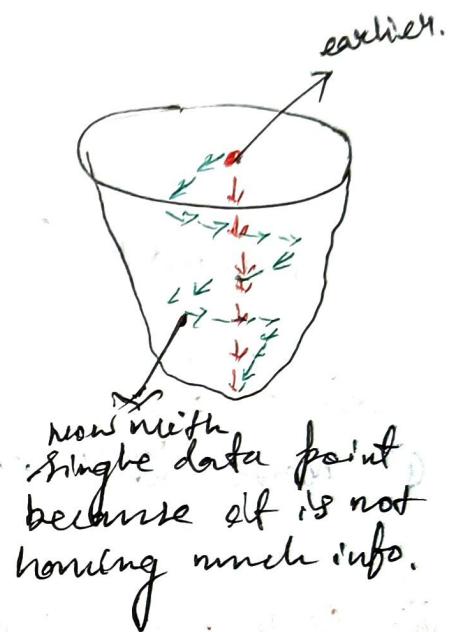
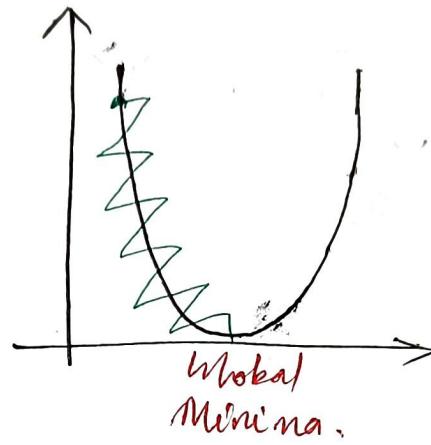
The gradient of the cost function is calculated from a single data point. This makes it very slow and prone to oscillations. It's only useful for small datasets.

Advantage

- ① Solves resource issue

Disadvantage

- ① Time complexity will be very high.
- ② Noise will be introduced because of single data point.
Movement inf in the graph is called noise.



② Mini Batch SGD

> Along with epochs and iterations, we give batch size.

$$\text{Data points} = 1000$$

$$\text{Batch size} = 100$$

$$\text{Epoch} = \frac{1000}{100} = 10$$

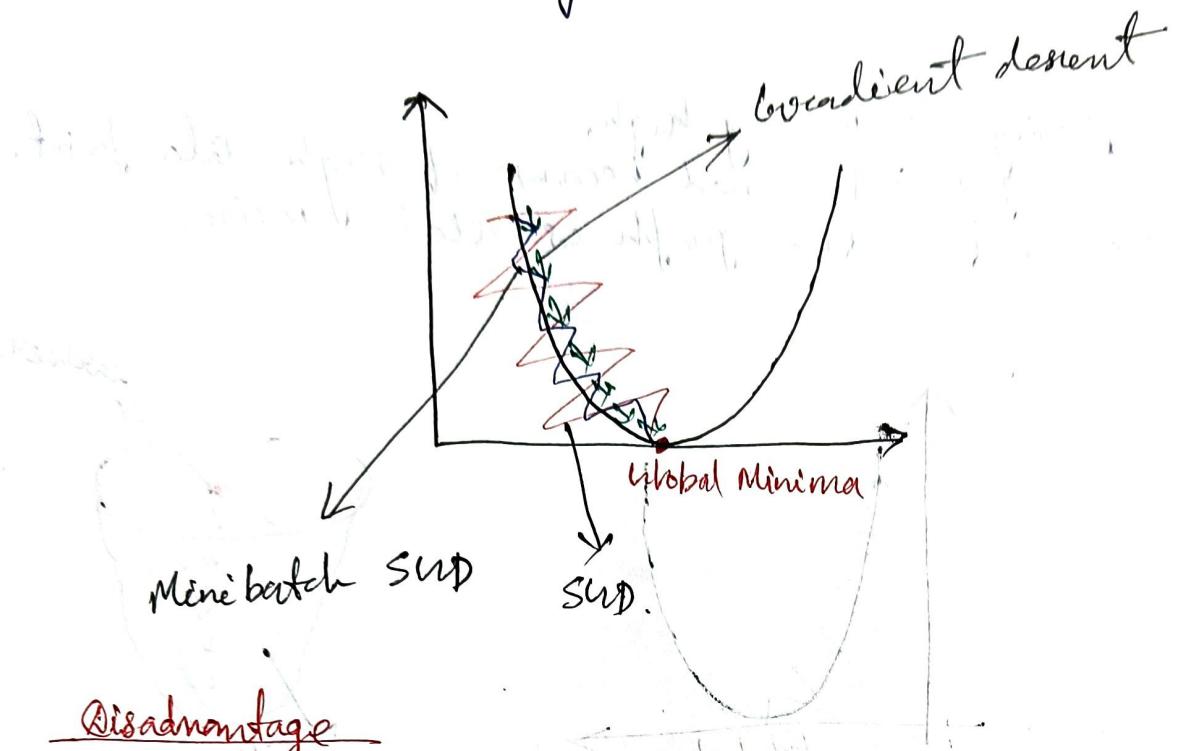
Weights and biases using 100 data points as batch

Epoch 1 { Start P.P. → Step 1 of J, cost function $J = 1.00 \times 10^6$ → End P.P. }

Epoch 1000

Advantages

- ① Shared resource dist., with convergence speed increases.
- ② Noise will get reduced, when compared to SGD.
- ③ Efficient resource usage.



Disadvantage

- ① Noise still exist.
- ② Step with momentum

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial \text{loss}}{\partial w_{\text{old}}}$$

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

$$w_t = w_{t-1} - \eta \frac{\partial \text{loss}}{\partial w_{t-1}}$$

Using this to smoothen

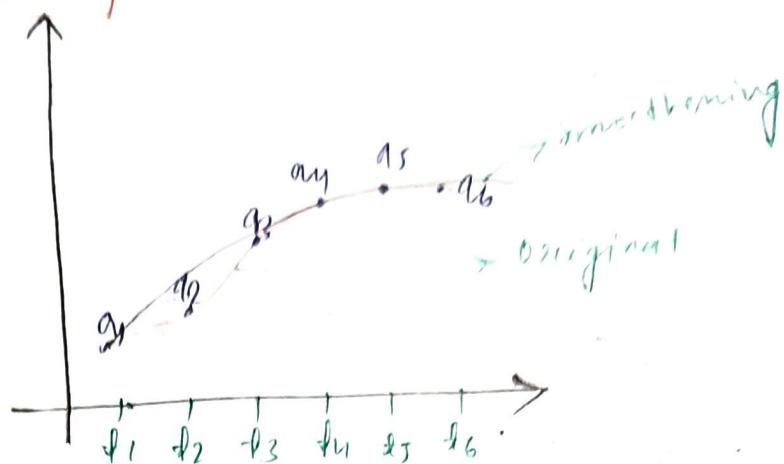
Exponential Weight Average (Smoothing)

Time $t_1 t_2 t_3 t_4 \dots t_n$

value $a_1 a_2 a_3 a_4 \dots a_n$

$$V_{t1} = a_1$$

$$V_{t2} = \beta \cdot V_{t1} + (1-\beta) \cdot a_2 \rightarrow \text{Exponential Smoothing}$$



Let's $\beta = 0.95$

$$V_{t2} = 0.95a_1 + 0.05a_2$$

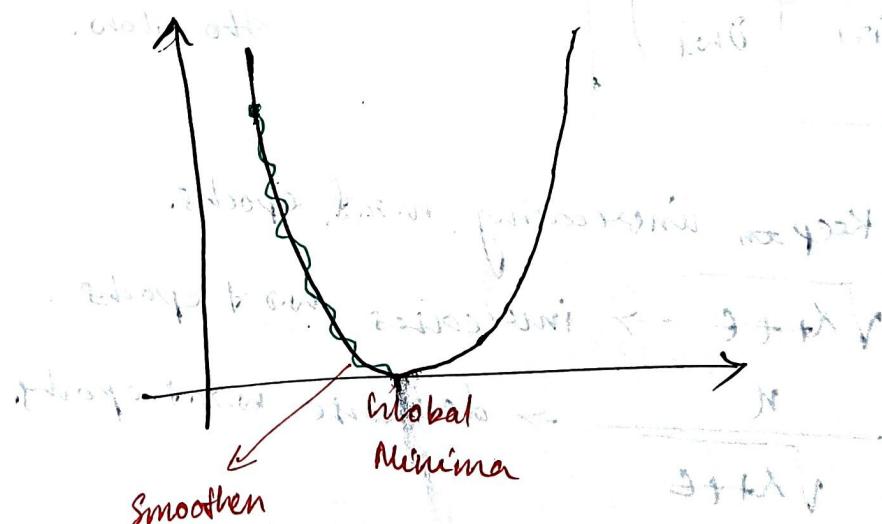
↓
has more control

→ In exponential smoothening, we have at $t-1$ more control
the value of t .

$$V_{t3} = \beta \cdot V_{t2} + (1-\beta) \cdot a_3$$

$$= 0.95 \cdot 92 + 0.05 \cdot 98$$

After Smoothing



→ Noise will get reduced.

Advantages

- ① It reduces the noise.
- ② It smoothen the noise.
- ③ Quick convergence.

(*)

⑤ Adagrad (Adaptive Gradient Descent)

$$w_t = w_{t-1} + \gamma'$$

$$w_t = w_{t-1} - \eta \frac{\partial \text{loss}}{\partial w_{t-1}}$$

$$w_t = w_{t-1} - \eta' \left[\frac{\partial \text{loss}}{\partial w_{t-1}} \right] \rightarrow \text{Updated}$$

loss = 0.5 * ||w||^2

$$\eta' = \eta$$

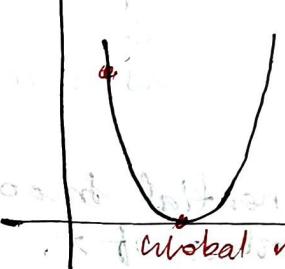
Dynamic learning rate

$$\sqrt{d_t + \epsilon} \rightarrow \text{Epsilon = small value to prevent the zero condition for } d_t$$

$$d_t = \sum_{i=1}^n \left(\frac{\partial \text{loss}}{\partial w_i} \right)^2$$

$\eta = \text{fixed}$

loss = 0.5 * ||w||^2



- > Convergence may happen at the constant rate.
- > But we want convergence should be faster at the beginning and after that it should be slow.

> d_t will keep increasing w.r.t epochs.

$$\sqrt{d_t + \epsilon} \rightarrow \text{increases w.r.t epochs}$$

$$\frac{\eta}{\sqrt{d_t + \epsilon}} \rightarrow \text{decrease w.r.t epochs.}$$

- Increase fig. size affect

Disadvantages

- ① If α becomes large in very large neural network
 $w_f \approx w_{f-1}$
- ② γ' possibility to become a small value ≈ 0 .

② Adadelta and RMS Prop (*Bringing exponential weighted average in learning rate*)

$$\gamma' = \frac{\alpha}{\sqrt{S_{dwt} + \epsilon}} \rightarrow \text{Dynamic Learning Rate}$$

Exponential Weight Average

Exponential Weighted Average

$$S_{dwt} = \dots$$

$$S_{dwt} = \beta \times S_{dwt-1} + (1-\beta) \left(\frac{\partial \text{loss}}{\partial w_f} \right)^2$$

③ Adam Optimizer

combination of SGD with momentum + RMS Prop (*Dynamic Learning Rate*)

$$w_f = w_{f-1} - \gamma' V_{dw}$$

$$V_{dwf} = \beta \times V_{dwf-1} - (1-\beta) \frac{\partial \text{loss}}{\partial w_{f-1}}$$

$$\gamma' = \frac{\alpha}{\sqrt{S_{dwt} + \epsilon}}$$

Biased weight updation

$$b_t = b_{t-1} - \eta' V_{db}$$

$$V_{db} = \beta * V_{db,t-1}$$

$$V_{db} = \beta * V_{db,t-1} + (1-\beta) \frac{\text{dloss}}{\partial b_{t-1}}$$

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

$$\left(\frac{\text{dloss}}{\partial b_{t-1}} \right) (q-1) \rightarrow \text{dloss} \times q = \text{peak}$$

comes up with a reinforcement view of the update rule

$$\left(\frac{\text{dloss}}{\partial b_{t-1}} \right)^2 \rightarrow \text{peak}^2$$

$$\left(\frac{\text{dloss}}{\partial b_{t-1}} \right) \rightarrow \text{peak} \times q = \text{peak}^2$$

$$\left(\frac{\eta}{\sqrt{S_{db} + \epsilon}} \right)^2 \rightarrow \eta^2$$