

9D Disadvantages.

- ① weights updation takes time.
- ② computationally expensive
- ③ take time to convergence.

9D advantages.

- ① optimal solution is guaranteed.
- ② smooth curvature.

Epoch

Whenever entire data set is being pass to the neural network is called epoch.

$$1000 \rightarrow NN \rightarrow \hat{y} \rightarrow \text{Loss} \quad \text{1 output 1 Epoch}$$

$$\begin{matrix} 10 \text{ Epochs} \\ 1000 \rightarrow NN \rightarrow \hat{y} \rightarrow \text{Loss} \quad \text{High} \\ \vdots \\ 1000 \quad \text{10th} \rightarrow NN \rightarrow \hat{y} \rightarrow \text{Loss} \quad \text{Low} \end{matrix}$$

Iteration.

$$\begin{matrix} \text{batch size} & \text{no. of batches} & \text{bias} \\ 100 & \nearrow \frac{10}{1000} & \downarrow \\ 1000 & & w_{\text{new}} = w_{\text{old}} - \alpha \frac{\partial L}{\partial w_{\text{old}}} \end{matrix}$$

$$\text{1st batch } 100 \rightarrow NN \rightarrow \hat{y} \rightarrow \text{Loss} \leftarrow |WB| \quad \text{Iteration 1}$$

$$\text{Updated } WB \rightarrow NN \rightarrow \hat{y} \rightarrow \text{Loss} \leftarrow |WB| \quad \text{Iteration 2}$$

$$\begin{matrix} \text{and then} & & & & \text{Iteration 10} \\ 1000 & 1000 & 1000 & \nearrow \frac{100}{1000} & \\ \text{1 Epoch} & 10 \text{ Epoch} & \text{Epoch 1} & \text{Iteration 10} & \\ 1000 & 1000 \times \frac{10}{1000} & & & \end{matrix}$$

Stochastic gradient Descent

Random.

	x_1	x_2	x_3	x_4	Yact
②	^	^	^	^	^
	^	^	^	^	^
①	^	^	^	^	^
	^	^	^	^	^

9D Entire data set

SGD 1 Record at a time which is randomly picked.

$$1 \text{ Record} \rightarrow NN \rightarrow \hat{y} \rightarrow \text{Loss} \rightarrow \text{update } |WB|$$

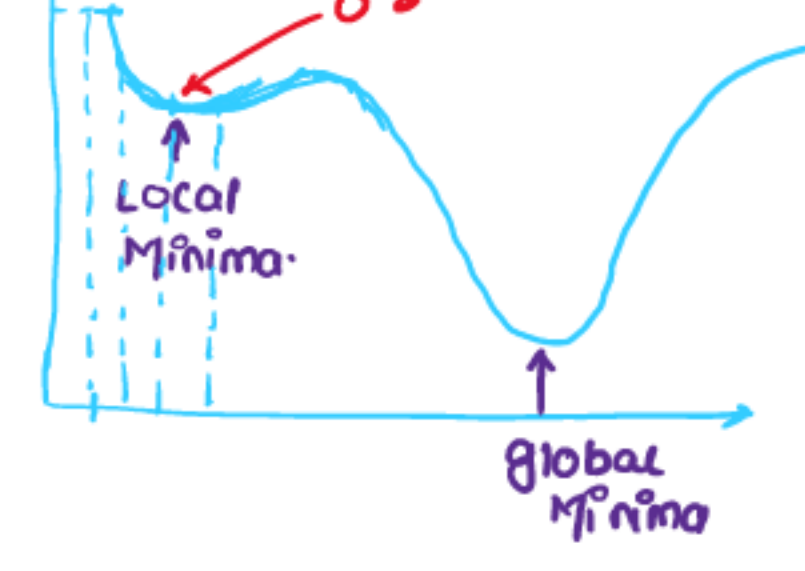
$$\text{Updated } WB \text{ and Record} \rightarrow NN \rightarrow \hat{y} \rightarrow \text{Loss} \leftarrow |WB|$$

Entire Dataset

Epoch 1, Iteration 1000, WB = 1000 times.

$$\begin{matrix} \times 10 & 1000 \times 10 & 1000 \times 10 \end{matrix}$$

SGD used to have multiple Local minima.



Slope

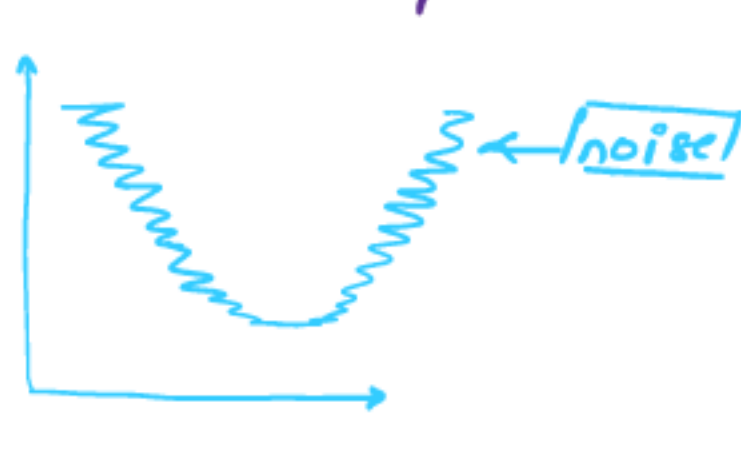
$$w_{\text{new}} = w_{\text{old}} - \alpha \frac{\partial L}{\partial w_{\text{old}}}$$

- ① Record 1 good Loss ↓
- ② Record 2 good Loss ↓
- ③ Record 3 good Loss ↓
- ④ Record 4 bad Loss ↑

- ① ES
- ② $w_{\text{new}} = \text{old } w$

Slope Reduce 0

2nd Disadvantage.



$$\begin{matrix} \text{Loss} \uparrow \\ \text{Loss} \downarrow \end{matrix} \leftarrow \begin{matrix} w \uparrow \\ w \downarrow \end{matrix}$$



as one sample from dataset is chosen randomly path taken to reach global minima is usually noiser than typical gradient Descent

Mini batch stochastic gradient Descent

GD ⇒ Entire Dataset
SGD ⇒ 1 by 1
MBSGD ⇒ Batch by batch

Data set 1000

$$\begin{matrix} 100 & 100 & 100 & 100 & 100 \\ \checkmark & \checkmark & \checkmark & \checkmark & \checkmark \end{matrix}$$

1 Epoch

5 Iteration

$$100 \text{ Batch} \rightarrow NN \rightarrow \hat{y} \rightarrow \text{Loss} \leftarrow |WB|$$

$$\text{Updated } WB \rightarrow NN \rightarrow \hat{y} \rightarrow \text{Loss} \leftarrow |WB|$$

200

$$\begin{matrix} \text{stochastic} & \text{Record } RP \\ MB & \rightarrow \text{Batch} \rightarrow \text{100 read time} \\ 100 & 100 \end{matrix}$$

$$\begin{matrix} SGD \\ \text{Batch size 1} \\ \text{Batch no 1000} \end{matrix}$$

$$\begin{matrix} \text{Wood } 1000 & \text{Wood } 1 & \text{Wood } 200 \\ \uparrow & \uparrow & \uparrow \\ \text{strong} & \text{noise} & \text{noise} \\ \downarrow & \downarrow & \downarrow \\ \text{noise} & \text{noise} & \text{noise} \end{matrix}$$

- ① Here we have solved the problem of noise to some extent from SGD.
- ② Batch sizes are often tuned to computational architecture.

such as power of 2

$$4 \quad 8 \quad 16 \quad 32 \quad 64 \quad 128 \rightarrow$$

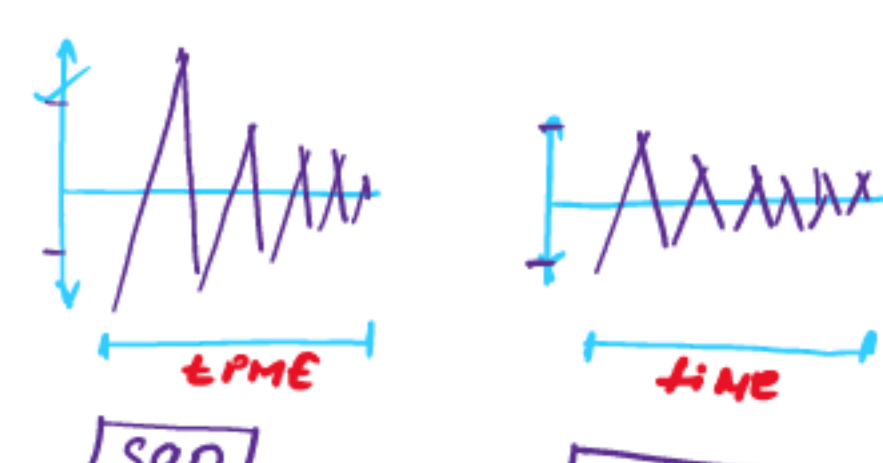
128 ← batch size

Depending upon batch size the updates can be made less noisy.

$$\text{batch } \uparrow \quad \text{noise } \downarrow$$

$$\text{no. of batch } \downarrow \rightarrow \text{9D}$$

$$\begin{matrix} 1000 \\ \text{no. of batch } \downarrow & \text{batch } \downarrow \\ 1 & 1000 \end{matrix} \leftarrow \text{MBSGD} \rightarrow \text{9D}$$



Exponential Moving Avg

$$\begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \\ 20 & 20 & 20 & 20 & 20 & 20 & 20 & 20 & 40 & 40 & P \end{matrix} \rightarrow \begin{matrix} \text{MA } 33.3 \\ \text{Exp } 27.2 \end{matrix}$$

Moving avg

$$\begin{matrix} 8 & 9 & 10 & 11 \\ 20 & 40 & 40 & 33.33 \end{matrix} \checkmark$$

EMA

$$P = B \times \text{previous avg} + (1-B) \times \text{Today}$$

$$= 0.8 \times 24 + (1-0.8) \times 40$$

$$= 0.8 \times 24 + 0.2 \times 40$$

$$= 19.2 + 8$$

$$P = 27.2 \quad \text{11th}$$

smoothing windows

$$\frac{20+20+20}{3} = 60/3 = 20$$

$$= \frac{60}{3} = 20$$

$$= \frac{100}{3} = 33.33$$

$$\frac{120}{5} = 24$$

$$= \frac{1}{1-B}$$

$$= \frac{1}{1-0.8} = \frac{1}{0.2}$$

$$= 5 \text{ days}$$

$$\frac{1}{1-B} = 10 \text{ days}$$

$$\uparrow \quad 0.7 = 3 \text{ days}$$

$$1 \text{ yes } 268 \leftarrow$$

$$\frac{1}{1-0.99} \checkmark$$

Exponential smoothing is basically used in time related problems. stock price prediction, climate forecasting.