

Tutorial-3

Optimizers and ANN

Day-3 - Deep Learning:

Agenda

① Optimizers

- (i) Gradient Descent
 - (u) SGD (Stochastic Gradient Descent)
 - (uu) Mini Batch SGD
 - (iv) SGD with momentum
 - (v) Adagrad (Adaptive Gradient Descent)
 - (vi) RMSPROP
 - (vii) Adam Optimizer

what is Batch, Epochs, Iterations

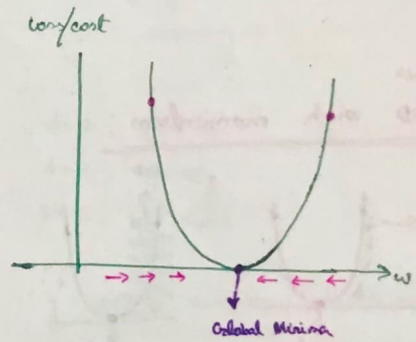
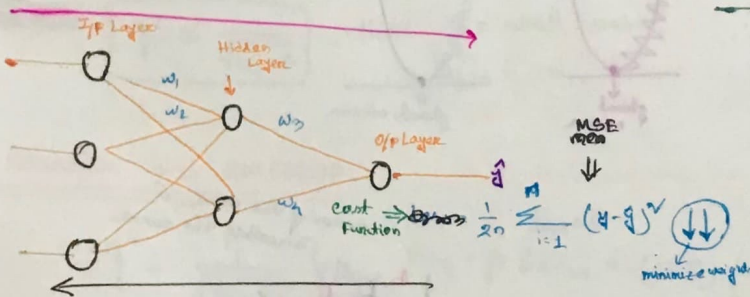
↓
ANN

① Gradient Descent:

Weight updation formula-

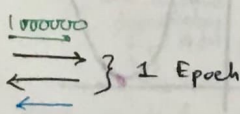
$$w_{new} = w_{old} - \eta \frac{\partial L}{\partial w_{old}}$$

learning rate



optimizers

Epoch



it means that {1000000}
1 F/w and B/w propagation.
Combination of forward and backward propagation.

Disadvantage

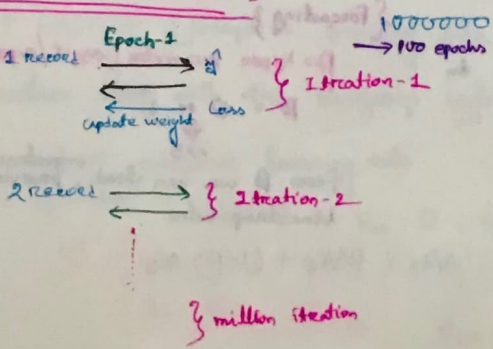
- ① Resource Extenese (Huge RAM)

② SGD (Stochastic Gradient Descent):

- ① Ram ↓↓

Disadvantages:

- ① Convergence will be very very slow
- ② Time complexity will be very high

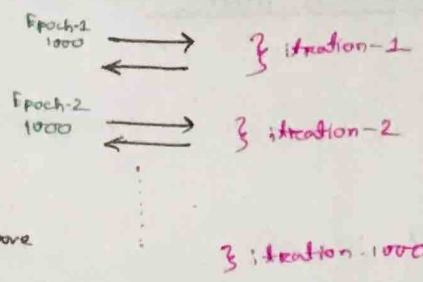


Mini Batch SGD:

Advantages:

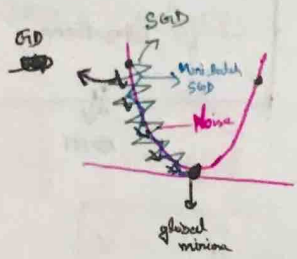
1. Doesn't resource intensive
2. Convergence will be better
3. Time complexity will reduce

10000000, batch-size = 1000



$$\text{iteration} = \frac{\text{records}}{\text{batch size}} = \frac{10000000}{1000} = 10000$$

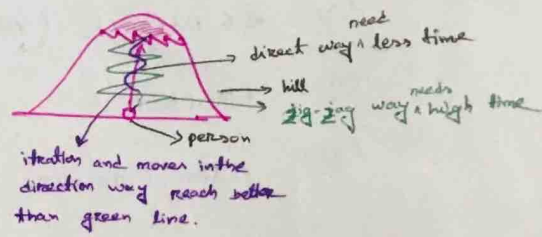
1st iteration's = 1000 records
2nd " = 1000 "
3rd " = 1000 "
...
10000 " = 10000 records



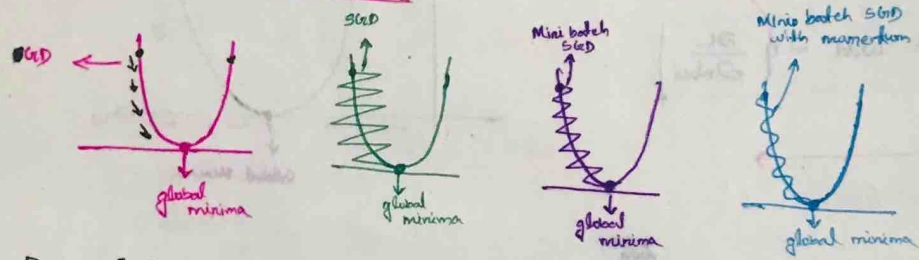
SGD → zig-zag is very high. Noise also high

Mini Batch SGD → zig-zag is very low. Noise also minimal
zig-zag is called noise

How do we remove Noise?
we use momentum.



Mini Batch SGD with momentum

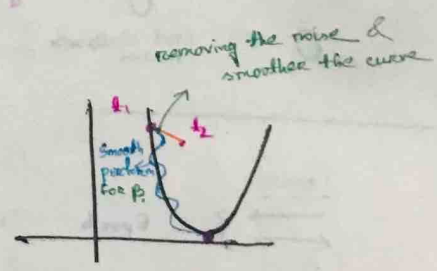


Exponential Weighted Average (moving)

used in
Time Series
model

ARMA, ARIMA

Weight Updation — $w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$
Bias — $b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}}$



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

current time previous time

Example:

- t1 t2 t3 t4 ... tn
a1 a2 a3 a4 ... an

$v_{t1} = a_1$ (previous)
 $v_{t2} = \beta v_{t1} + (1-\beta) a_2$ (current)
 $= (0.95) v_{t1} + (0.05) a_2$
very important value is high less important value is low

{ Forecasting }

β hyper parameter (weight/exponential moving avg)

$\beta = 0 \text{ to } 1$
↓
0.95

From β we can find importance, previous or current timestamp value.

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

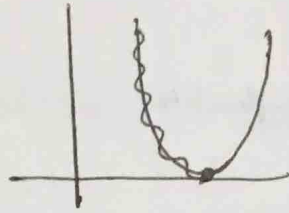
Exponential weighted / Moving Avg,

$$w_t = w_{t-1} - \eta \nabla w$$

$$\nabla w_t = \beta \nabla w_{t-1} + (1-\beta) \frac{\partial L}{\partial w_{t-1}}$$

Advantages:

- (i) Reducing the noise
- (ii) Mini Batch
- (iii) Quicker convergence



Adagrad (Adaptive Gradient descent)

$$\eta = \text{fixed} \Rightarrow \text{Adaptive (change } \eta) \Rightarrow \text{Learning rate decreasing} \Rightarrow \text{global minima}$$

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

$$w_t = w_{t-1} - \eta' \frac{\partial L}{\partial w_{t-1}}$$

$\eta' = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$

η' is a small number (0.01) is a small number. if $\alpha_t = 0$, $\eta' = \infty$. to remove this problem ϵ is used.

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

α_t is increasing can't control it

When α_t is increasing, η' is decreasing.

because $\eta' \propto \frac{1}{\sqrt{\alpha_t + \epsilon}}$. Here $\eta = \text{small number}$

When model was standing to train we should keep on decreasing it as we reach the global minima.

$$t=1, \eta = 0.01$$

$$t=2, \eta = 0.05$$

$$t=3, \eta = 0.003$$

$\therefore \eta'$ is small number, that's why $w_t \approx w_{t-1}$

Adadelta and RMSPROP:

Exponential weighted average

$$S_{dw} = 0$$

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

there will be a slow decreasing

$$S_{dw_t} = \beta S_{dw_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

$$\beta = 0.95$$

$$S_{dw_t} = (0.95) S_{dw_{t-1}} + (0.05) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

it can control small values

We can control S_{dw} so that there will be a

slow decreasing of this learning rate, η' and we will be able to reach the global minima.

Adam Optimizer: (best optimizer forever)

momentum + RMSPROP (Adaptive learning rate)

Combine with momentum along with the rmsprop intuition.

Intuition means, we have adaptive learning rate

diff types of Adam optimizer

$$V_{dw} = 0, V_{db} = 0, S_{dw} = 0, S_{db} = 0$$

(i) adam max

Adam optimizer

$$w_t = w_{t-1} - \eta' \nabla w$$

$$b_t = b_{t-1} - \eta' \nabla b$$

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

$$V_{dw_t} = \beta V_{dw_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

$$V_{dw_t} = \beta V_{dw_{t-1}} + (1-\beta) \frac{\partial L}{\partial w_{t-1}}$$

(ii) Smoothing

(iii) Learning Rate Adaptive