# Optimisation in Deep Learning

GD

SGD

Minibatch SGD

SGD with momentum

Adagrad (Adaptive stochastic Gradient)

RMS prop

Adadelta

Adam.

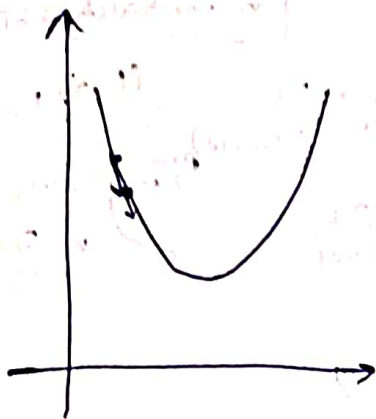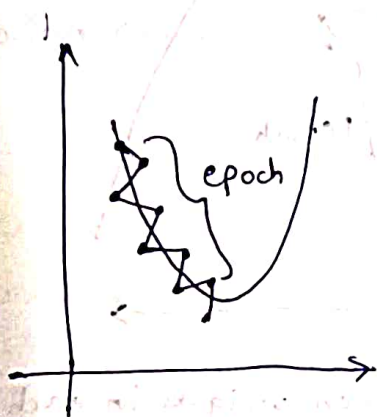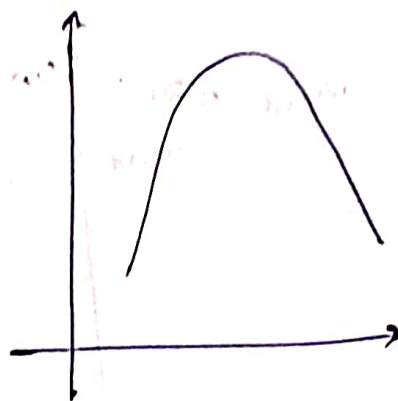| GD | SGD | Minibatch SGradient descent |
|---|---|---|
| -Takes whole data in 1 iteration | -Say we have 100 samples, then in (each sample individually) | total = 100 samples |
| | epoch-1 | |
| epoch = epoch-1 _____ iter-1 | iter1   1st sample ⇄ forward prop backward prop | batch size = 10 |
| iter = iter-1 sample Considers ~~each~~ all samples as single unit. | iter2   2nd sample ⇄ prop | no. of ~~samples~~ ~~for each~~ batches = 100/10 = 10 |
| | : : | iter = 10 |
| | iter 100   100th sample ⇄ | epoch1  iter's |
| | | 1st iter   1 to 10 ᵉ 2nd iter   11 to 20 ᵉ 3rd iter, 21 to 30ᵉ : 10th iter  91 to 100 |

---
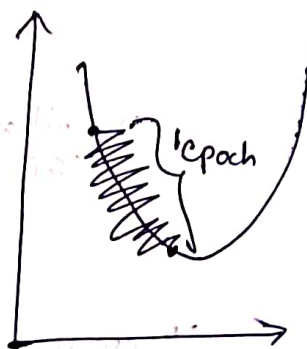
**disadv:-**

- more RAM
- computationally expensive

-It will take less RAM
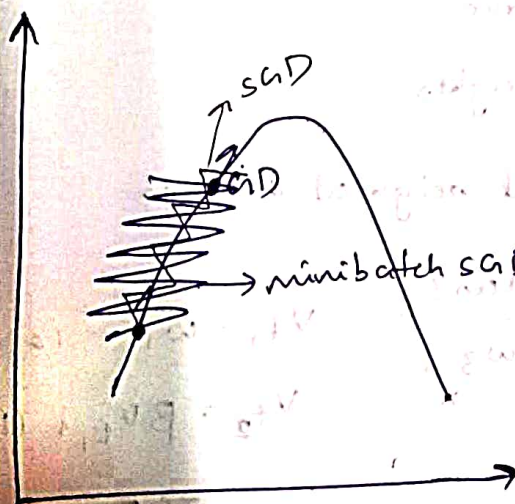-But still computationally expensive.

(For GD we get only one value)



For SGD, for 1 epoch, it will go in zig-zag fashion

more zigzag in minibatch SGD.



minibatch SGD (takes extreme left/right) does not give a smooth curve like (in all 3)

→ Disadv for all three :- It does not give optimum path to reach the goal.

| GD | SGD | minibatch, SGD |
|---|---|---|

**GD**

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

**SGD**

$$(\omega_{new})_{iter} = (\omega_{old})_{prev \ iter} - \eta \frac{\partial L}{\partial (\omega_{old})_{prev \ iter}}$$

$$\omega_t = \omega_{t-1} - \eta \frac{\partial L}{\partial \omega_{t-1}}$$

**minibatch, SGD**

$$((\omega) \ batch)_t = (\omega \ batch)_{t-1} - \eta \frac{\partial L}{(\partial \omega batch)_t}$$



minibatch SGD

There is no relation between weights in each batch
So, it is giving steep values

## SGD with momentum

→ This gives a smooth curve

→ It considers the new weights

→ It calculates exponential weighted average

| item1 | item2 | item3 |
|---|---|---|
| $\omega_1$ | $\omega_2$ | $\omega_3$ |

equivalent to

physics acceleration.

$$V_{t_1} = \omega_1 \qquad (a_1)$$

$$V_{t_2} = \beta V_{t_1} + (1-\beta)\overset{(a_2)}{\omega_2}$$

If $\beta = 0.9$,
90% importance to
previous iteration
and 10% to current
iteration.

$$V_{t_3} = \beta V_{t_2} + (1-\beta)\overset{(a_3)}{\omega_3}$$

$$= \beta(\beta V_{t_1} + (1-\beta)\omega_2) + (1-\beta)\omega_3$$

For SGD, $\omega_t = \omega_{t-1} - \eta \dfrac{\partial L}{\partial \omega_{t-1}}$

For SGD with momentum, $\omega_t = \omega_{t-1} \eta \boxed{V_{d\omega_t}}$ → This term is derivative of $V_{t_i}$

$$V_{d\omega_1} = \frac{\partial V}{\partial \omega_t}(\omega_1)$$

$$V_{d\omega_2} = \beta V_{\partial \omega_1} + (1-\beta)\frac{\partial L}{(\partial \omega_2)}$$

For iteration 1, $\omega_1 = \omega_0 - \eta V_{d\omega_1}$

For iteration 2, $\omega_2 = \omega_1 - \eta V_{d\omega_2}$

$V_{d\omega_3} = \beta V_{d\omega_2} + (1-\beta)\dfrac{\partial L}{\partial \omega_3}$

→ SGD is independent of previous weights

→ SGD with momentum gives importance to previous weights (90%) and remaining 10% to current weight.

SGD with momentum gives smooth curve



11/10/22

- Adagrad
- RMS prop
- Adadelta
- Adam

## Adagrad

- Adaptive gradient descent

→ $\eta$ is dynamic

At each iteration, if we change learning rate, it will give a better result.

→ So, works better than SGD with momentum

$$w_t = w_{t-1} - \eta_t\left(\frac{\partial L}{\partial w_{t-1}}\right)$$

$$\eta_t = \frac{\eta}{\sqrt{\alpha_t + \varepsilon}} \qquad \varepsilon?$$

$$\alpha_t = \sum_{i=1}^{t}\left(\frac{\partial L}{\partial w_i}\right)^2 \longrightarrow \begin{array}{l}\text{Cumulative} \\ \text{sum of current} \\ \text{plus previous} \\ \text{values}\end{array}$$

e.g:- if t=3,

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

---

## RMS prop

→ Root mean square propagation optimization.

$$w_t = w_{t-1} - \frac{\eta}{\sqrt{v_t + \varepsilon}}\frac{\partial L}{\partial w_{t-1}}$$

$$v_t = \beta v_{t-1} + (1-\beta)\left(\frac{\partial L}{\partial w_t}\right)^2$$

→ improvement to adagrad

→ Instead of taking cumulative sum, it takes exponential weighted average.

---

## Adadelta

→ Adaptive delta optimization

→ $\eta$ term is completely removed and delta term is added.

$$w_t = w_{t-1} - \frac{\sqrt{D_{t-1} + \varepsilon}}{\sqrt{v_t + \varepsilon}} \times (\Delta w_t)$$

$$\Delta w_t = w_t - w_{t-1}$$

$$D_t = \beta D_{t-1} + (1-\beta_1)(\Delta w_t)^2$$

$$v_t = \beta_2 v_{t-1} + (1-\beta_2)\left(\frac{\partial L}{\partial w_t}\right)^2$$

## Adam

→ adaptive moment estimation

→ momentum + RMS prop
  (combines both)

$$w_t = w_{t-1} - \frac{\eta * V_{dw_t}}{\sqrt{S_{dw_t} + \varepsilon}}$$

momentum ↗

RMS prop ↘

$$V_{dw_t} = \beta_1 V_{dw_{t-1}} + (1-\beta_1)\frac{\partial L}{\partial w_t}$$

(momentum) ↑

$$S_{dw_t} = \beta_2 S_{dw_{t-1}} + (1-\beta_2)\left(\frac{\partial L}{\partial w_t}\right)^2$$

↳ RMS prop

This technique is further
improved by bias correction

$$w_t = w_{t-1} - \frac{\eta * \cancel{V}_{dw_t}}{\cancel{\sqrt{S}}}$$

$$w_t = w_{t-1} - \frac{\dfrac{\eta * V_{dw_t}}{(1-\beta_1)t}}{\sqrt{\dfrac{S_{dw_t}}{(1-\beta_2)^t} + \varepsilon}}$$

...
...

$$w_t = w_{t-1} - \frac{\eta * V_{dw_t}}{\sqrt{S_{dw_t} + \varepsilon}}$$

?

$$V_{dw_t} = \frac{\beta_1}{(1-\beta_1)} V_{dw_t} + \frac{\partial L}{\partial w_t}$$   ?

$$S_{dw_t} = \frac{\beta_2}{(1-\beta_2)} S_{dw_t} + \left(\frac{\partial L}{\partial w_t}\right)^2$$   ?

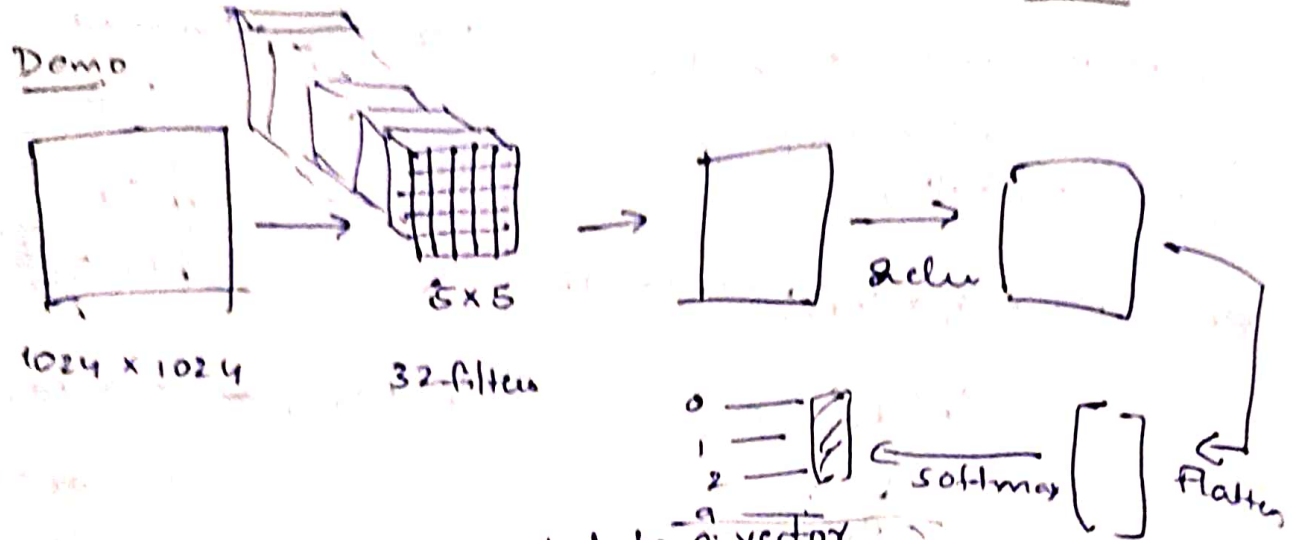Disadvantage

① It may suffer from
vanishing gradient problem

$$\eta_{t+1} = \frac{\eta_b}{\sqrt{\alpha_{t+1}} + \varepsilon}$$   ?

→ vanishing
gradient → $\alpha_t = \alpha_{t+1}$

⇒ $\eta_t = \eta_{t-1}$

advantage

1) It may provides better
convergence than SGD with
momentum.

Demo



1024 × 1024                5 × 5

32 filters
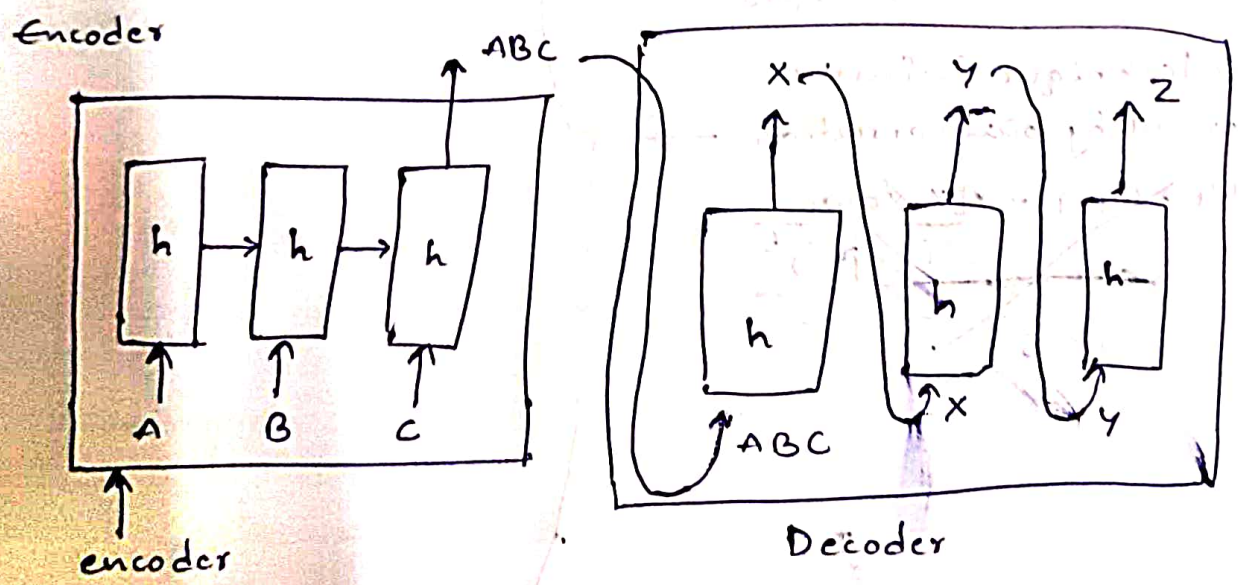
2 clu

softmax

Flatten

- flatten :- matrix converted to a vector
- Densenet :- each and every neuron is connected with every other
  (fully connected)

conv 2D — grayscale image

## Encoder - Decoder RNN

Encoder



encoder

Decoder

ABC → XYZ (translation)

ABC → WXY (It might not translate accurately)

Incorrect translations → noise

$\left( \begin{matrix} RNN \\ LSTM \end{matrix} 3f \right)$
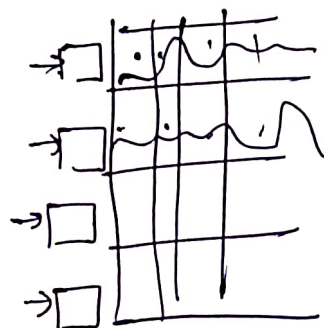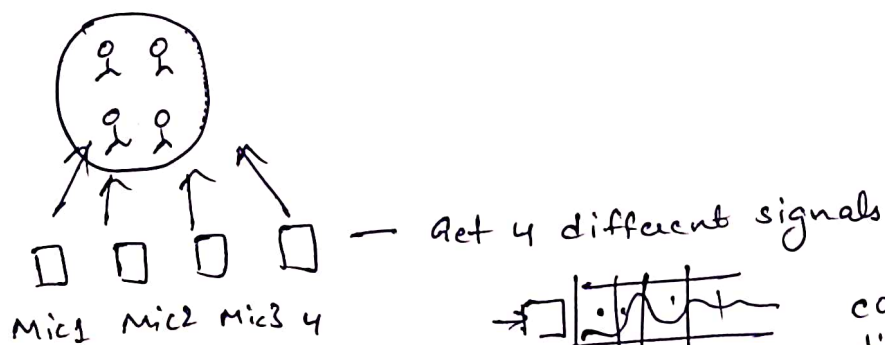
Encoder Decoder LSTM (is also there).

## ||UNIT-3||

Probabilistic models of the input $P_{model}(x)$

For example, to rate mess food,
→ ambience
→ arrangement  } service
→ cleanliness
→ quantity of food
→ quality  } food
→ variety in menu

( Reduction in meaningful less-features $\hat{no}$.of )

In ML, we do this using dimensionality reduction



— get 4 different signals

Mic1  Mic2  Mic3  4

convert into discrete values, extract individual signals.

From the observed variable, getting non-observed variable — <u>latent variables</u>.

$$P_{model}(x) = E_h \, P_{model}\left(\frac{x}{h}\right)$$

Latent variables: applns in economics, ...

## <u>Linear factor models</u>

→ Data generation in LFMs:

1. Sample the explanatory