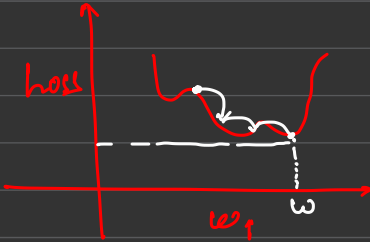
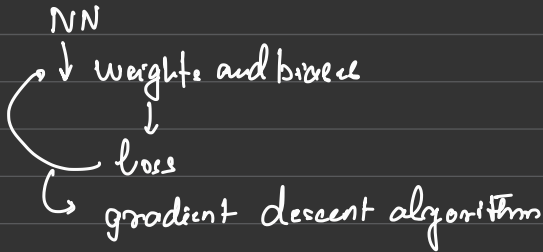
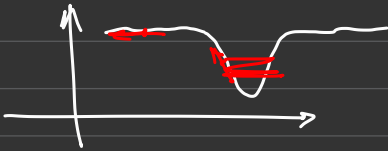


Optimizers



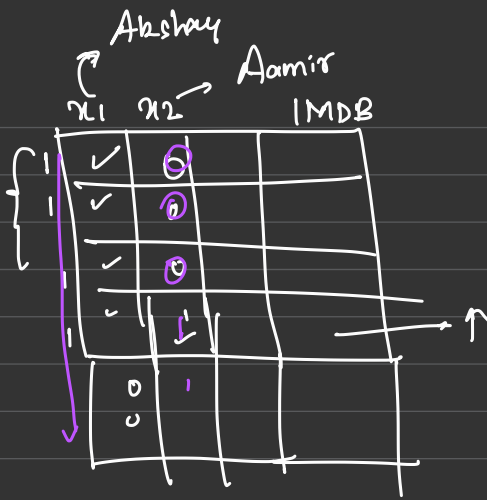
loss curves are extremely complicated functions \rightarrow multiple minimas

- \rightarrow be flat at certain regions
- be very sharp at certain regions

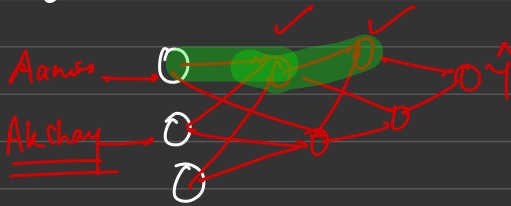


Gradient descent by design has various issues.

- \rightarrow Multiple local minimas.
- \rightarrow The gradients are not consistent
- \rightarrow Sparse features create a problem in GP



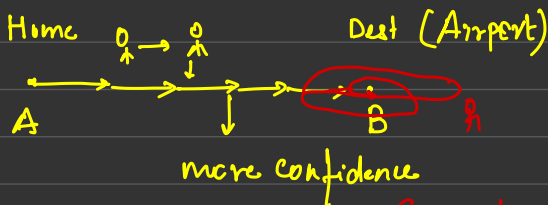
most of the rows of aamir Khan's features
 $\rightarrow 0$



weights related to Aamir Khan
 weights related to Akshay Kumar

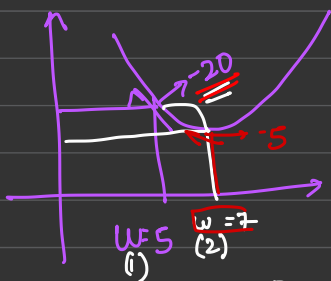
$$w_{i+1} = w_i - \alpha \frac{dJ}{dw_i}$$

Momentum



GD

GD with momentum



$$w(2) = w(1) - \alpha \left(\frac{dw}{dw} \right) \rightarrow 0.1$$

$$w(2) = w(1) - \alpha \times -20$$

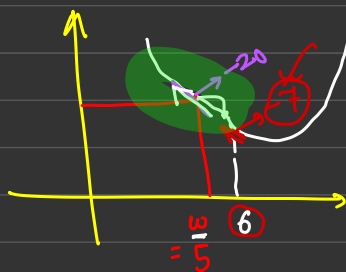
$$w(2) = 5 + 2$$

$$w(2) = 7$$

$$w(3) = w(2) - 0.1 \times -5$$

$$w(3) = 7 + 0.5$$

$$= 7.5$$



$$m(0) = 0$$

$$m(t) = \beta m(t-1) + (1-\beta) \times \frac{dw}{dw}$$

$$\text{update} \rightarrow w(2) = w(1) - \alpha(m(t))$$

$$m(1) = 0.5 \times 0 + 0.5 \times -20$$

$$m(1) = -10$$

$$w(2) = 5 - 0.1 \times -10$$

$$w(2) = 6$$

$$m(2) = 0.5 \times -10 - 0.5 \times -7$$

$$m(2) = [5 + 3.5] = -8.5$$

$$w(3) = w(2) - 0.1 \times -8.5$$

$$w(3) = 6 + 0.85$$

$$= 6.85$$

Momentum based gradient update actually
takes into account the historical gradient \times weightage
+ the current gradient \times weightage.

↓
accelerates the process of the update



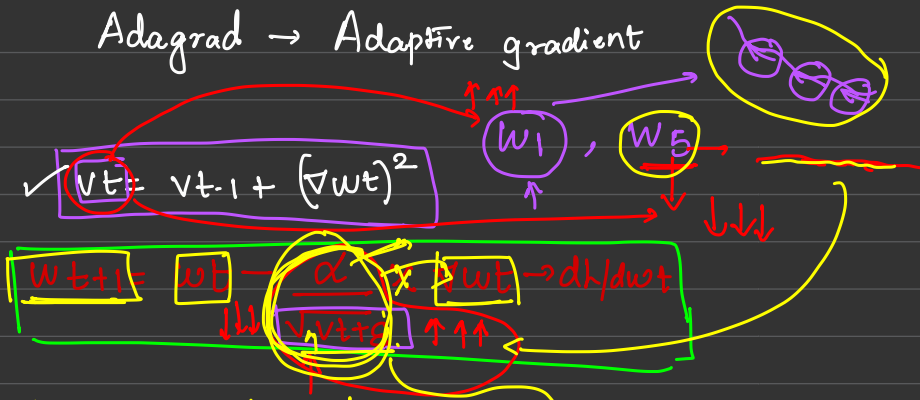
$$wb_{t+1} = wb_t - \alpha mb_t$$

$$mb_t = \beta \times mb_{t-1} + (1-\beta) \times \frac{dL}{dw_t}$$

history
point/current

near the global min
it takes longer to cross
it is so fast it
crosses / ignores local
min

Adagrad \rightarrow Adaptive gradient



Weights which will see big updates will become smaller

Adaptive gradients start punishing weight updates for weights which have already seen big gradient updates
 \rightarrow don't punish/reward weight updates for weights which have not seen big gradient updates

Rewarding (Amir Khan) $\uparrow \uparrow$
 Punish (Ashay Kumar) $\downarrow \downarrow$ \rightarrow balanced hearing

$$v_t = v_{t-1} + (dw_t)^2 \rightarrow \text{too high, too fast}$$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{v_{t+1} + \epsilon}} \times dh/dw_t$$

Rms Prop.

$$v_t = \boxed{B \times v_{t-1} + (1-B) \times (w_t)^2}$$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{v_t + \epsilon}} \times dh/dwb$$

ϵ
Epsilon

Adam optimizer.

Adaptive moment

momentum

$$\begin{aligned} \rightarrow m_t &= B_1 \times m_{t-1} + (1-B_1) \times \nabla w_t \\ \rightarrow v_t &= B_2 \times v_{t-1} + (1-B_2) \times (w_t)^2 \rightarrow \text{RMS} \end{aligned}$$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{v_t + \epsilon}} \times \hat{m}_t$$

$$\hat{m}_t = \frac{m_t}{1-B_1} \quad \hat{v}_t = \frac{v_t}{1-B_2}$$