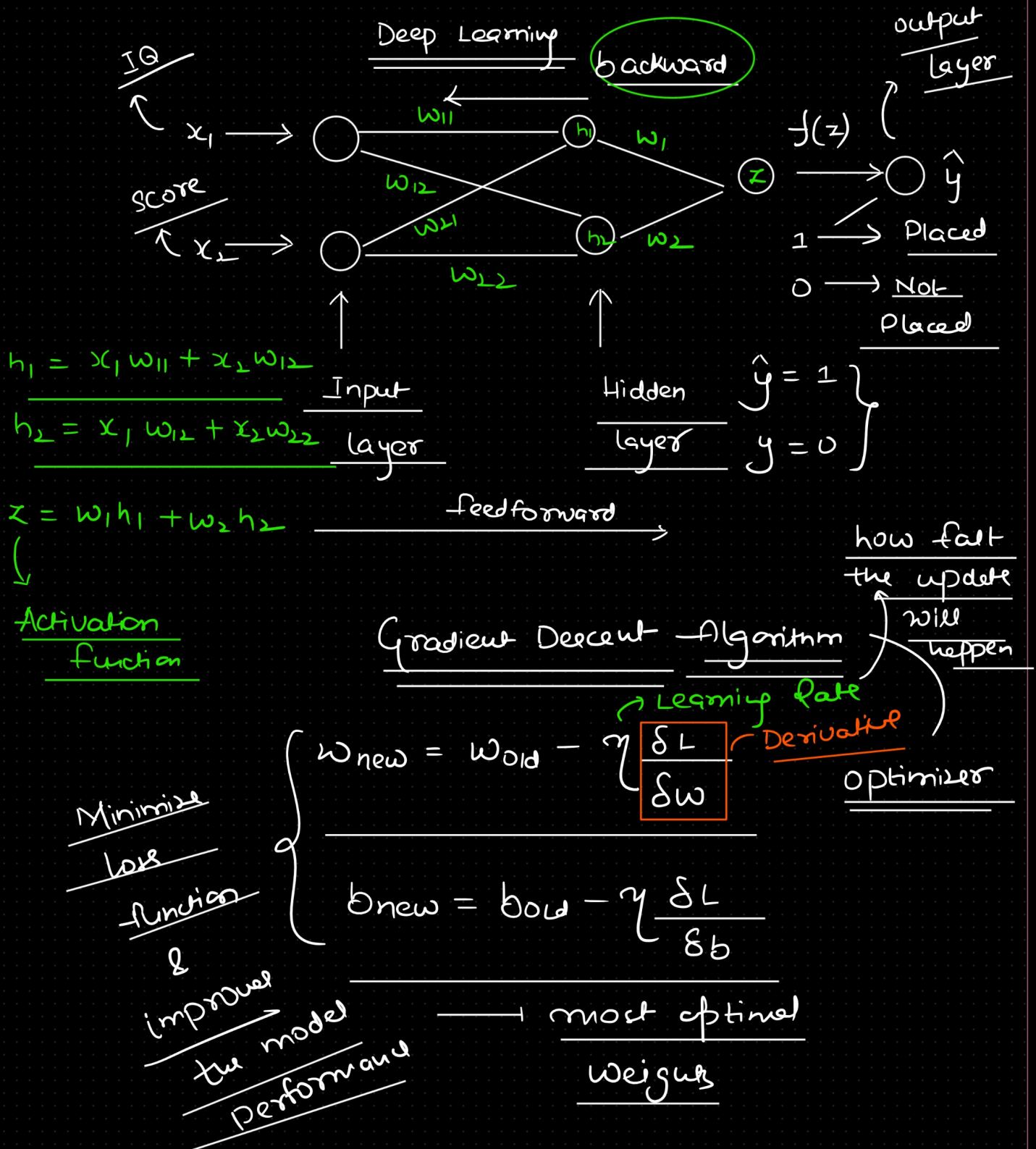
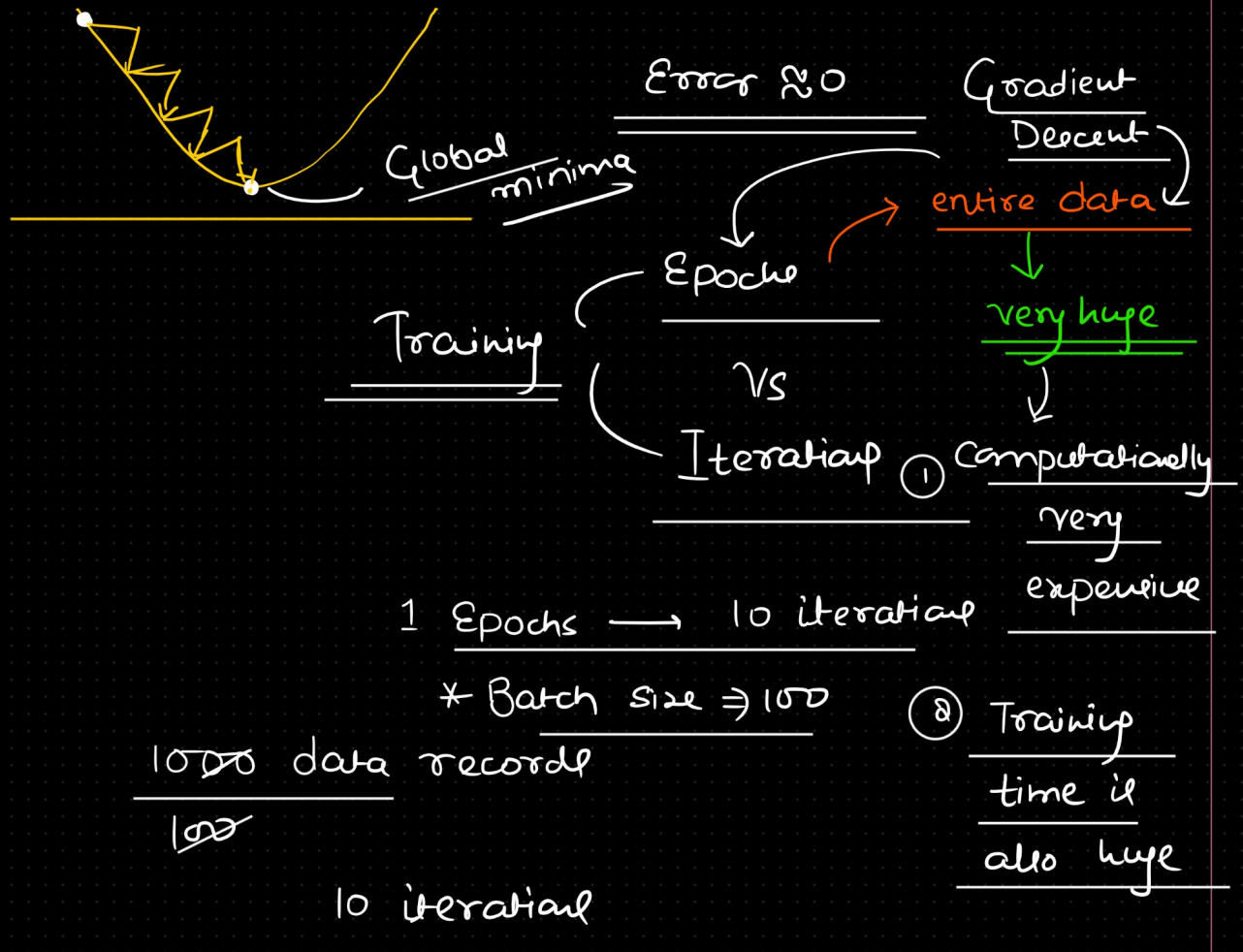
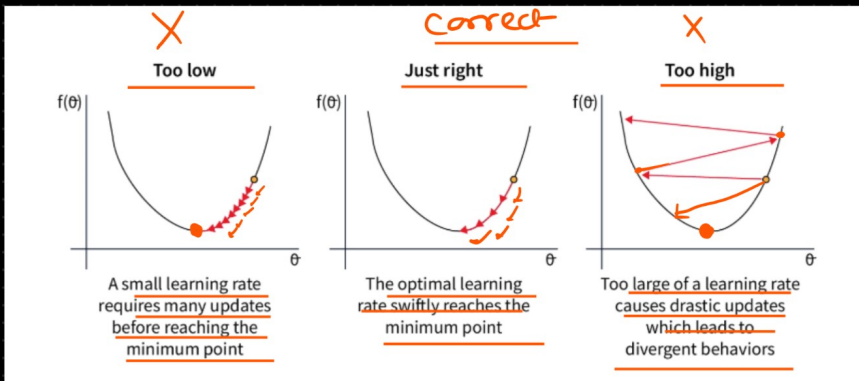


Optimizers

Optimizers are algorithms that adjust a model's parameters during training to minimize a loss function and improve performance.



Learning Rate \rightarrow hyperparameter
hyperparameter tuning ≈ 0.001 (before training)



SGD (Stochastic Gradient Descent)

1 training sample at a time

1000 records

100 iterations \rightarrow 1 epoch

Reduce \leftarrow
the Computational
task

Mini batch Gradient Descent

batch \Rightarrow 100

num of records \Rightarrow 1000

$\#$ iterations $\Rightarrow \frac{1000}{100} = 10$

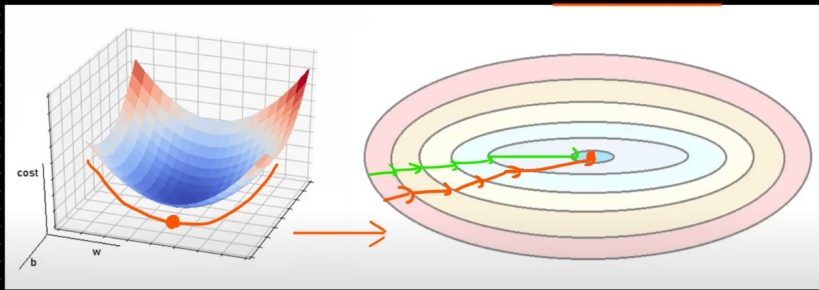
20 epochs

Version of Gradient Descent

for every single epoch,
10 iterations

- Gradient Descent \rightarrow entire data
- Stochastic Gradient Descent \rightarrow single training sample at every iteration
- ** Mini batch Gradient Descent

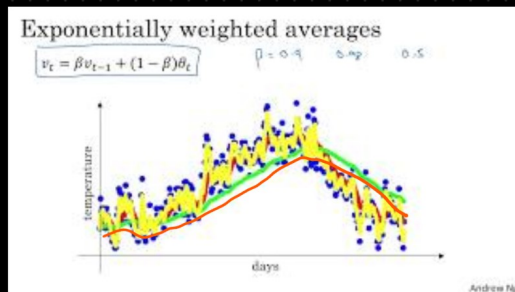
Contour



Exponential average

↓
Momentum
optimizer

Smoothing



$0 < \beta < 1$

$\beta = 0.9$

$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$

$v_0 = 0$

$v_1 = \beta v_0 + (1 - \beta) \theta_1$

$v_1 = (1 - \beta) \theta_1$

$v_2 = \beta [(1 - \beta) \theta_1] + (1 - \beta) \theta_2$

$v_4 = (1 - \beta) (\theta_4 + \beta \theta_3 + \beta^2 \theta_2 + \beta^3 \theta_1)$

$w_{\text{new}} = w_{\text{old}} - \eta v_{dw}$

What is the impact of β ?

Adagrad (Adaptive Gradient Boosting)

↳ $\eta = 0.001$ (fixed before training)

dynamic $\alpha \leftarrow \eta$ (Adaptive)

Adam \rightarrow Maximum Realtime industry work

\rightarrow Momentum + RMS Prop

\downarrow

better generalizability

Note:- In every optimizer, the only difference is in the update of two parameters.