

Data Science & AI



Machine Learning



Artificial Neural Network

One Shot



By- Krish Naik Sir

Recap of Previous Lecture



Topic

Deep Learning

Topic

Topic

Topic

Topic

Topics to be Covered



Topic

Optimizers

Topic

Drop outliers

Topic

Topic

Topic

Deep learning

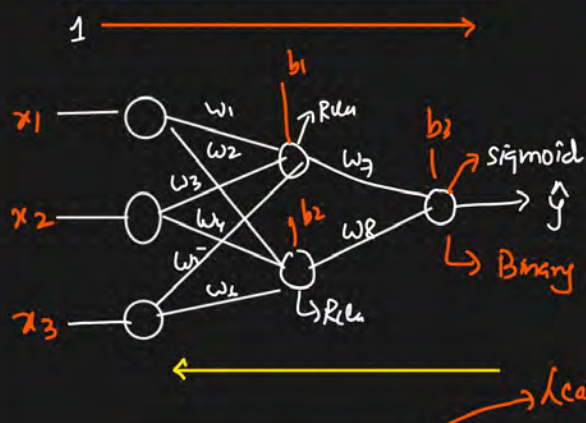
- ① Optimizers ✓
- ② Weight Initialization Technique ✓
- ③ DROPOUT LAYERS ✓

{ BATCH NORMALIZATION }

① Optimizers

- ① Gradient Descent ✓
- ② Stochastic Gradient Descent ✓
- ③ Mini batch SGD ✓
- ④ SGD With Momentum ✓
- ⑤ Adagrad and Rmsprop
- ⑥ Adam Optimizers. ← Best

① GRADIENT DESCENT Optimizer



MSE, RMSE, MAE

$$\text{Loss} = [\text{Error}] \downarrow \downarrow$$

Gradient Descent



$$w_{\text{new}} = w_{\text{old}} - \eta \left[\frac{\partial L}{\partial w_{\text{old}}} \right]$$

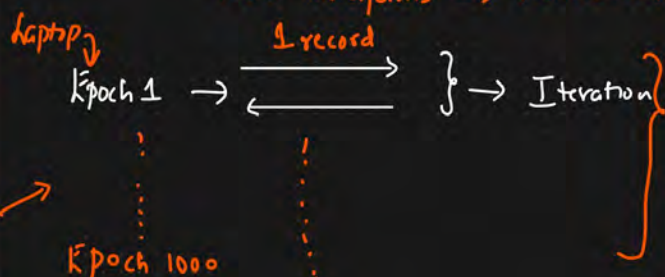
DATASET = 1000 datapoints

$$\text{Batch size} = 1000/10 = 100$$

1000 iteration

10 Iteration

1000 datapoints → TRAINING



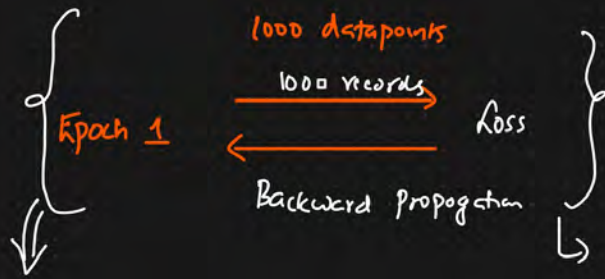
200

Epoch 1

Epoch 2

$$\sum_{i=1}^{100} (y_i - \hat{y}_i)^2$$

Time Training ↑↑



TRAINING SET = 1 Million Records.

System Hungry \Rightarrow Nvidia GPU data center.

Cost $\uparrow\uparrow$

\Downarrow
Highly Expensive

Gradient Descent \Rightarrow Convergence

2 Epoch = 1 iteration

100 Epochs.

Advantages

① Convergence will happen

Disadvantage

① Huge Resource, RAM, GPU

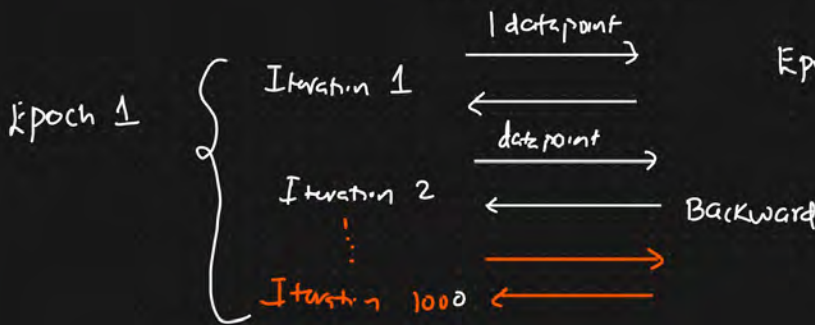
② Stochastic Gradient Descent

batch_size = 1



Epoch & Iteration

Training = 1000 records

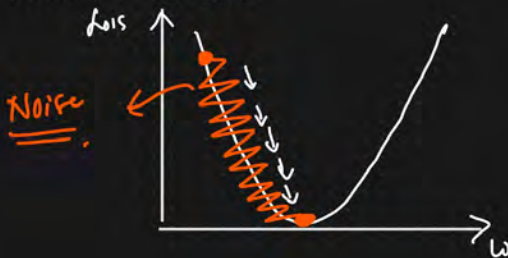


Epoch = 50, 100, 150, 200

Optimizers \rightarrow Changing weight And bias.

Advantage

① Solves Resource Issue



Disadvantage



① Time Complexity $\uparrow\uparrow$

② Convergence will take more time

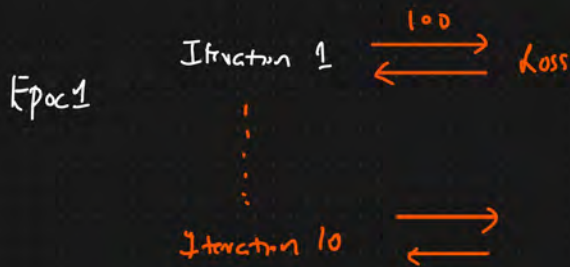
③ Noise get introduced

③ Mini batch SGD

Batch size = 100

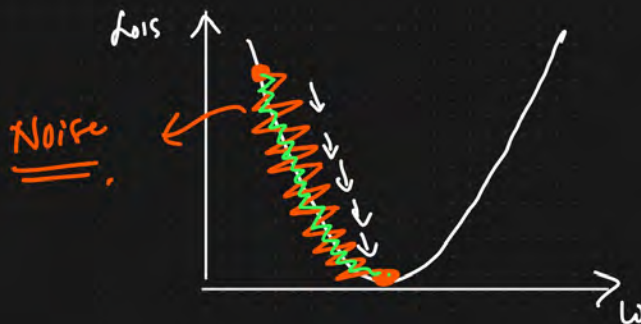
Training = 1000

$$\text{iteration} = \frac{1000}{100} = 10 \text{ iteration}$$



50, 100,

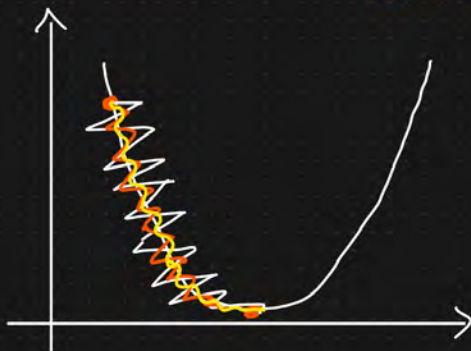
Early Stopping Loss ↓ X



④ SGD With Momentum

Noise
[Smoothing]

Exponential Weighted Average = ARIMA



Time $t_1 \quad t_2 \quad t_3 \quad t_4 \quad \dots \quad t_n$

Values $a_1 \quad a_2 \quad a_3 \quad a_4 \quad \dots \quad a_n$

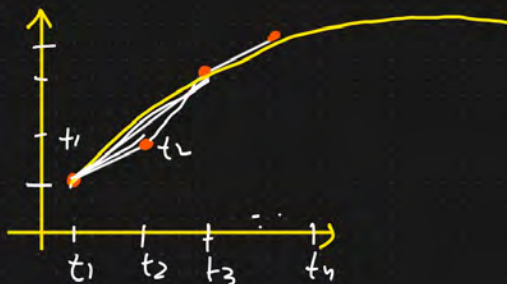
β = Influence

$$V_{t_1} = a_1$$

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

$$= [0.95 * a_1 + (0.05) * a_2]$$

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



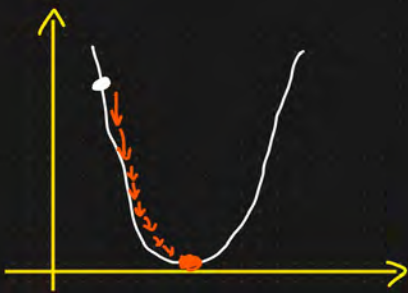
$$\boxed{w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}} \Rightarrow \boxed{w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}}$$

④ Adagrad → Adaptive Gradient Descent

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

Learning Rate fixed.

$\eta = \text{fixed} \Rightarrow$ Dynamic Learning Rate.



\Rightarrow As the convergence happens the learning rate should change

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

$$\eta' = \frac{\eta}{\sqrt{d_t + \epsilon}} \quad \eta' \downarrow \downarrow \downarrow \quad d_t \uparrow \uparrow \uparrow$$

$\sqrt{d_t + \epsilon} \Rightarrow$ Small value

$$\begin{matrix} \rightarrow & \rightarrow & \rightarrow \\ t=1 & t=2 & t=3 \end{matrix}$$

$$\eta = 0.01 \quad \eta = 0.005 \quad \eta = 0.003 \quad \dots$$

$$\boxed{0.000001}$$

Disadvantage

① $\eta' \rightarrow$ Possibility to become a very small value ≈ 0 .

⑥ Adadelta OR RMSPROP

$$\eta' = \frac{\eta}{\sqrt{S_{dw_t} + \epsilon}} \quad \underline{\underline{\beta = 0.95}}$$

η' slowly reduce $\eta' \downarrow \downarrow$

Exponential Weight Average

$$S_{dw_t} = 0$$

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

$$b_{\text{new}} = b_{\text{old}} - \eta' \frac{\partial L}{\partial b_{\text{old}}}$$

⑦ Adam Optimizer

SGD with Momentum + **RMSPROP** [Dynamic LR + Smoothing]

$$\begin{aligned} w_t &= w_{t-1} - \eta' v_{dw} \\ b_t &= b_{t-1} - \eta' v_{db} \end{aligned}$$

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

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

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

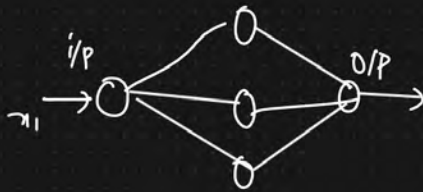
$$v_{db_t} = \beta * v_{db_{t-1}} - (1-\beta) \frac{\partial L}{\partial b_{t-1}}$$

Weight Initializing Technique

Key Points

- ① Uniform Distribution
- ② Xavier/Glorot Initialization
- ③ Kalming He Initialization

- ① Weights should be small
- ② Weights should not be same
- ③ Weight should have good Variance



$$\begin{aligned} i/p &= 1 \\ o/p &= 1 \end{aligned}$$

① Uniform Distribution

$$W_{ij} \sim \text{Uniform Distribution} \left[\frac{-1}{\sqrt{\text{input}}}, \frac{1}{\sqrt{\text{input}}} \right]$$

② Xavier/Glorot Initialization

Researcher \rightarrow Xavier Glorot

① Xavier Normal Init

$$W_{ij} \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{(\text{input} + \text{output})}}$$

② Xavier Uniform

$$W_{ij} \sim \text{Uniform Distr} \left[\frac{-\sqrt{6}}{\sqrt{i/p + o/p}}, \frac{\sqrt{6}}{\sqrt{i/p + o/p}} \right]$$

③ Kaiming He Initialization

① He Normal

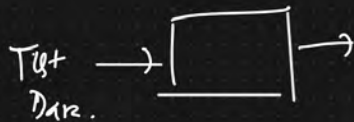
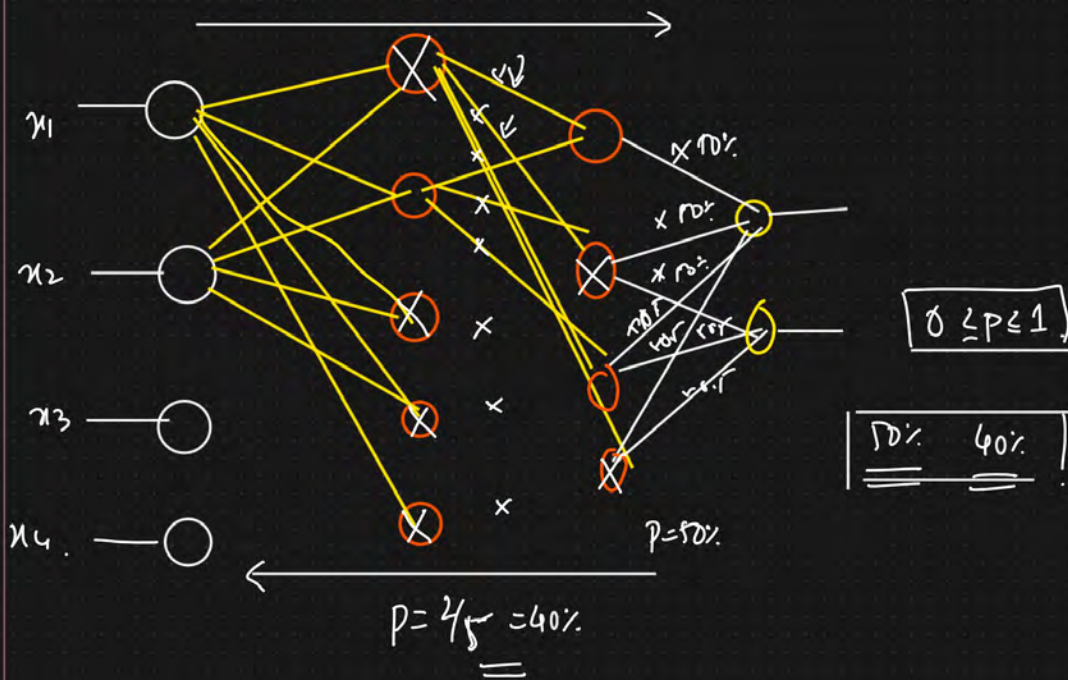
$$w_{ij} \sim \mathcal{N}(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{i/p}}$$

② He Uniform

$$w_{ij} \sim \text{Uniform}_{\text{Dist}} \left[-\sqrt{\frac{6}{i/p}}, \sqrt{\frac{6}{i/p}} \right]$$

Drop out layer





2 mins Summary



Topic

One

Topic

Two

Topic

Three

Topic

Four

Topic

Five

THANK - YOU