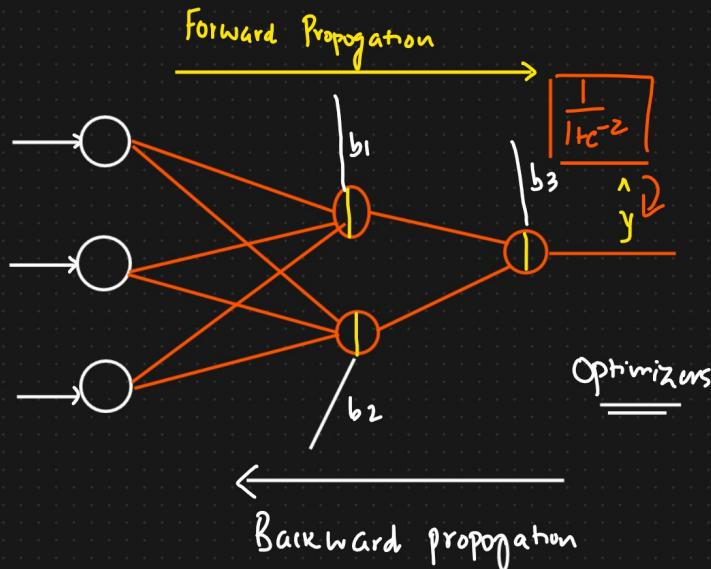


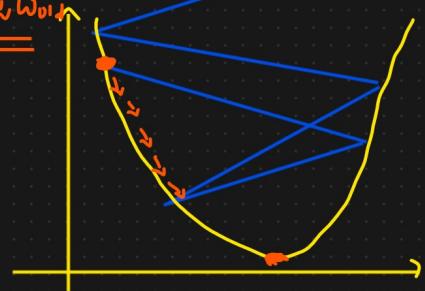
① Loss function And Cost function

② Optimizers



Weight Initialization $w_{new} \approx w_{old}$

10000



Loss function {Single datapoint} MSE
= $(y - \hat{y})^2$

Cost function = $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ {Batch data points}.



① MSE {Mean Squared Error}

② MAE {Mean Absolute Error}

③ Huber Loss

④ RMSE

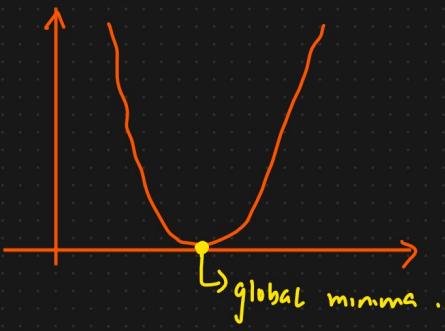
① Mean Squared Error

$$\text{Loss fn} = (y - \hat{y})^2$$



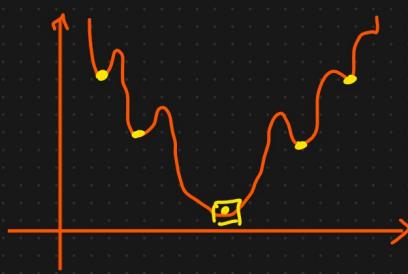
Quadratic Equation

$$\text{Cost fn} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



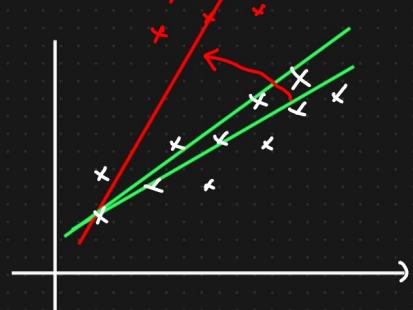
Advantages

- ① MSE is Differentiable
- ② It has 1 local or 1 global Minima
- ③ It converges faster



Disadvantages

- ① Not Robust to outliers



② Mean Absolute Error (MAE)

$$\text{Loss fn} = |y - \hat{y}|$$

$$\text{Cost fn} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

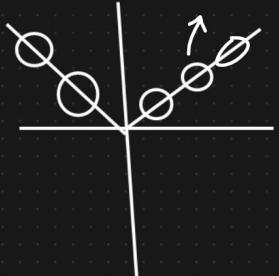
Subgradient

Advantages

- ① Robust to outliers

Disadvantage

- ① Convergence is slow



③ Huber Loss

① MSE

② MAE

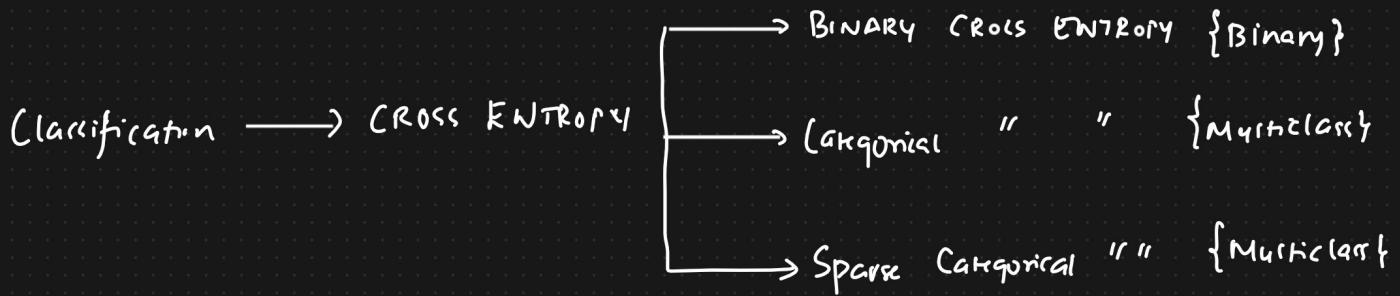
$$\text{Cost fn} =$$

$$\begin{cases} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise.} \end{cases}$$

\Downarrow Hyperparameter

\Downarrow
MAE

* Loss or Cost function for Classification Problem



① Binary Cross Entropy

$$\text{loss} = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y})$$

\Downarrow *log loss*

$$\text{loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

$$\hat{y} = \frac{1}{1+e^{-z}}$$

\Rightarrow Sigmoid
Activation
fn

② Categorical Cross Entropy (Multiclass Classification)

(-3)

	f_1	f_2	f_3	O/P	$j=1$ Good	$j=2$ Bad	$j=3$ Neutral	$C = \text{No. of Categories}$
$i=1$	2	3	4	Good	1	0	0	
$i=2$	5	6	7	Bad	0	1	0	
$i=3$	8	9	10	Neutral	0	0	1	

C

$$\text{loss}(x_i, y_i) = - \sum_{j=1}^C y_{ij} * \ln(\hat{y}_{ij})$$

Actual values $\leftarrow y_{ij} = [y_{11} \quad y_{12} \quad y_{13} \dots y_{1C}]$

$$[y_{21} \quad y_{22} \quad y_{23} \quad \dots \quad y_{2c}]$$

$$y_{ij} = \begin{cases} 1 & \text{if the element is in the class} \\ 0 & \text{Otherwise} \end{cases}$$

Prediction $\Rightarrow \hat{y}_{ij} \Rightarrow$ Softmax Activation = $Sof(z) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$

$$= \frac{e^{z_i}}{e^1 + e^2 + e^3 + \dots + e^K}$$

③ Sparse Categorical Entropy

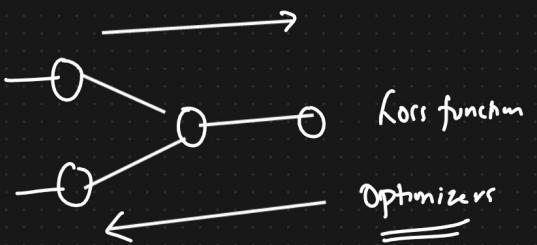
$$\left[\begin{smallmatrix} 0 & 1^{st} & 2^{nd} \\ 0.2, & 0.3, & 0.5 \end{smallmatrix} \right] \downarrow \quad \left[\begin{smallmatrix} 0.2 \downarrow 0.3, & 0.1, & 0.2, & 0.2 \end{smallmatrix} \right] \downarrow$$

2nd Index \Rightarrow Op. $\rightarrow 0.3 \rightarrow 1^{st}$ Index

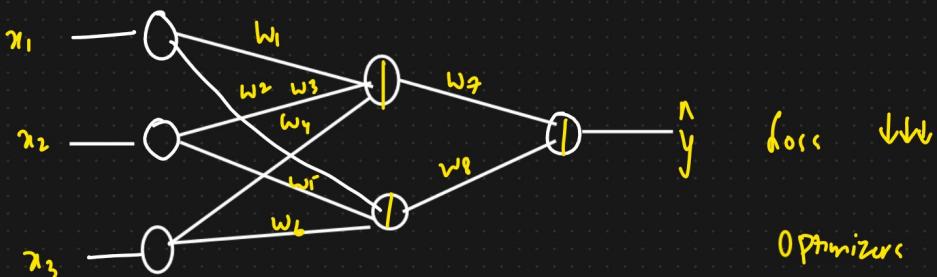
	<u>Right Combination</u>	Activation applied		
	<u>Hidden Layers</u>	<u>O/P Layer</u>	<u>Problem Statement</u>	<u>Loss function</u>
①	ReLU or its Variants	Sigmoid	Binary Classification	Binary Cross Entropy
②	ReLU or its Variants	Softmax	Multi Class	Categorical or Sparse CE
③	ReLU or its Variants	Linear	Regression	MSE, MAE, RMSE, R-squared, RMSE

Optimizers

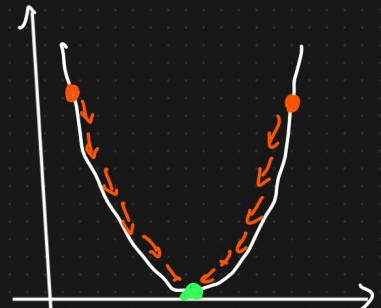
- (1) Gradient Descent
- (2) SGD (Stochastic Gradient Descent)
- (3) Mini batch SGD
- (4) SGD with Momentum
- (5) Adagrad and RMSProp
- (6) Adam Optimizers



① Gradient Descent Optimizer



loss fn
Optimizer



Weight Update Formula

$$w_{old} = w_{old} - \eta \left[\frac{\partial h}{\partial w_{old}} \right] \Rightarrow \text{Slope} \Rightarrow \text{Chain Rule of Derivation}$$

MSF

$$\text{Loss} = (y - \hat{y})^2$$

$$\text{Cost fn} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

KEpochs vs Iteration

$$[\text{DATAPOINTS} = 1000] \Rightarrow \underline{\underline{1000000}}$$

$$\begin{aligned} \text{Epoch 1} & \left\{ \begin{array}{l} \xrightarrow{1000 \text{ datapoints}} \\ \xleftarrow{1 \text{ Iteration}} \end{array} \right. & \hat{y}_i &= \text{Cost} \downarrow \downarrow \\ & \text{Weights will get update} & & \end{aligned}$$

100 Epochs

Advantage

895 R9m

Disadvantage

- ① Convergence will happen
 - ① Huge Ram And GPU \Rightarrow Resource Intensive.

② Stochastic Gradient Descent (SGD)

4th Ram

1000 datapoints

$$100 \times 100^\circ = 1$$

Iteration 1

Los ss

1

1

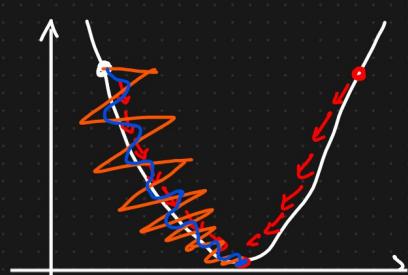
Jterghm 1000

A hand-drawn green outline of a Christmas tree with a zigzag pattern on its trunk and branches, enclosed within a circle.

Epoche 2

11

Epoch 100



Advantage

Disadvantage

- ④ Solve Resource Issues
 - ① Convergence will take time
 - ② Noise will get introduced

③ Mini Batch SGD

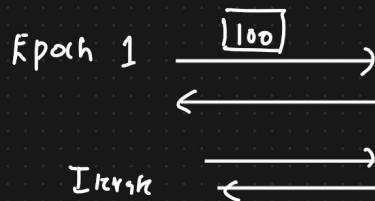
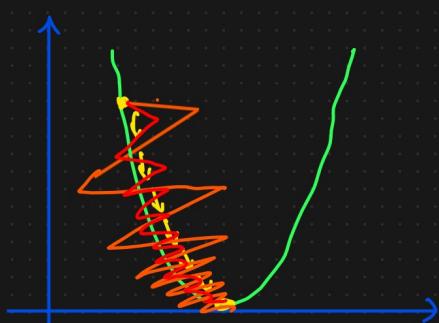
8gb RAM

10000
1000

Batch size

100

Epoch, Iteration, Batch size



Advantage

- ① Convergence speed increases
- ② Noise will be less when compared to SGD
- ③ Efficient Resource Usage

Disadvantage

- ① Noise still exists

④ SGD With Momentum

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

$$b_{new} = b_{old} - \eta \left[\frac{\partial h}{\partial b_{old}} \right]$$

$$w_t = w_{t-1} - \eta \left[\frac{\partial h}{\partial w_{t-1}} \right]$$

Exponential Weight Average {Smoothing}

Time t₁ t₂ t₃ t₄ ... t_n

Value a₁ a₂ a₃ a₄ ... a_n

Exponential Weighted Average

$$\underline{\beta=0.05}$$

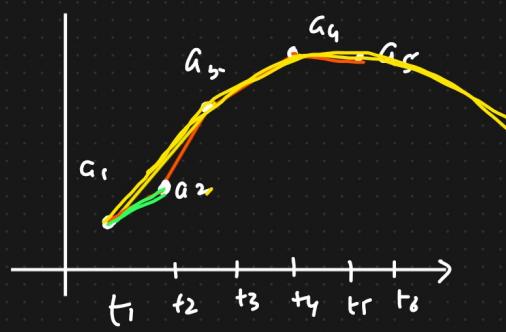
$$V_{t_1} = a_1$$

$$\boxed{0.95} \leftarrow \boxed{\beta=0.95}$$

$$\underline{\beta>0}$$

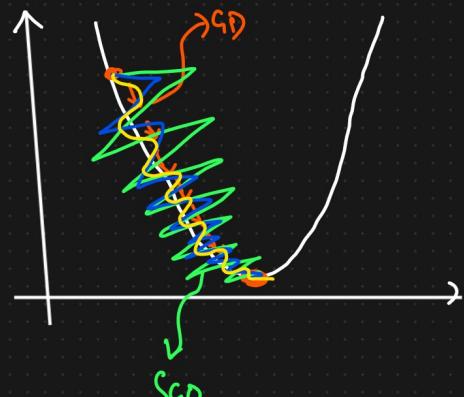
$$V_{t_2} = \beta * V_{t_1} + (1-\beta) * a_2 \quad \beta \downarrow$$

$$= \boxed{0.95 * a_1} + \boxed{0.05 * a_2}$$



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

$$= 0.95 [$$



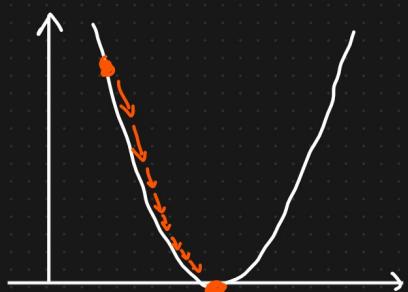
Advantage

- ① Reduces the noise
- ② Smoothen the noise
- ③ Quick Convergence

(f) Adagrad $\div \{ \text{Adaptive Gradient Descent} \}$

$\eta = \text{fixed} \Rightarrow \text{Dynamic LR}$

$$W_t = W_{t-1} - \eta \frac{\partial h}{\partial w_{t-1}}$$



$$W_t = W_{t-1} - \eta \frac{\partial h}{\partial w_{t-1}}$$

$$\eta^1 = \eta \quad \eta^1 \downarrow \downarrow \downarrow$$

$$W_t \approx W_{t-1}$$

$$d_f = \sum_{i=1}^t \left(\frac{\partial h}{\partial w_i} \right)^2 \Rightarrow \begin{cases} \sqrt{d_f + \epsilon} \xrightarrow{\epsilon \text{ small value}} \text{Large value} \\ \epsilon \text{ large value} \end{cases}$$

Disadvantages

① η^1 = Possibility to become a small value ≈ 0

② Adadelta And RMS Prop [Exponent weighted Exponential weighted Avg]

Learning Rate $\boxed{S_{dw_t} = 0}$ \downarrow

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

$\eta^1 = \frac{\eta}{\sqrt{S_{dw_t} + \epsilon}}$

S_{dw_2}

Adam Optimizer

SGD with Momentum + RMSprop [Dynamic LR]

$$w_t = w_{t-1} - \eta^1 v_{dw}$$

$$\eta^1 = \frac{\eta}{\sqrt{S_{dw_t} + \epsilon}}$$

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

$$\sqrt{S_{dw_t} + \epsilon}$$

$$b_t = b_{t-1} - \eta^1 v_{db}$$

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

$$\eta^1 = \frac{\eta}{\sqrt{S_{db_t} + \epsilon}}$$

$$\eta^1 = \frac{\eta}{\sqrt{S_{db_t} + \epsilon}}$$