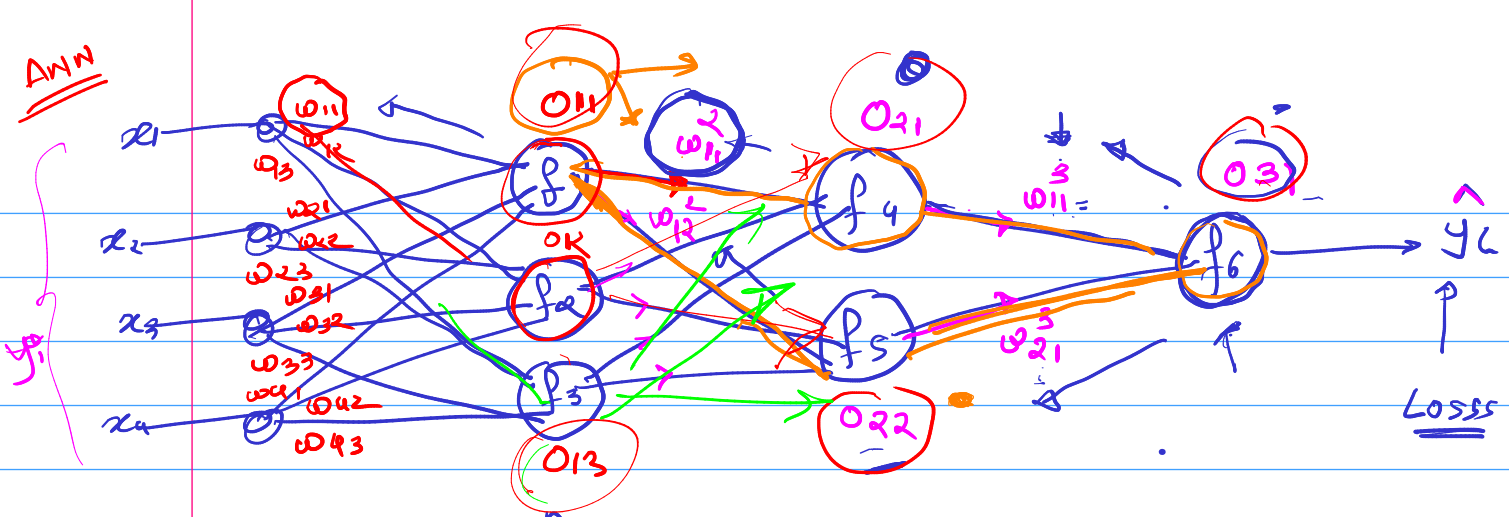


ANN



$$\text{Loss} = \sum_{i=1} (y - \hat{y})^2$$

to minimize the loss, we use optimizer.

Back propagation → change / update, weights & bias.

$$w_{11}^3_{\text{new}} = w_{11}^3_{\text{old}} - \eta \frac{dL}{dw_{11}^3_{\text{old}}}$$

chain Rule

$$\frac{dL}{dw_{11}^3_{\text{old}}} = \frac{dL}{d031} \cdot \frac{d031}{dw_{11}^3} = \text{chain rule}$$

$$\frac{dL}{dw_{21}^3} = \frac{dL}{d031} \cdot \frac{d031}{dw_{21}^3}$$

$$\frac{dL}{dw_{11}^2} = \left[\frac{dL}{d031} \cdot \frac{d031}{d021} \cdot \frac{d021}{dw_{11}^2} \right] + \left[\frac{dL}{d031} \cdot \frac{d031}{d022} \cdot \frac{d022}{dw_{11}^2} \right]$$

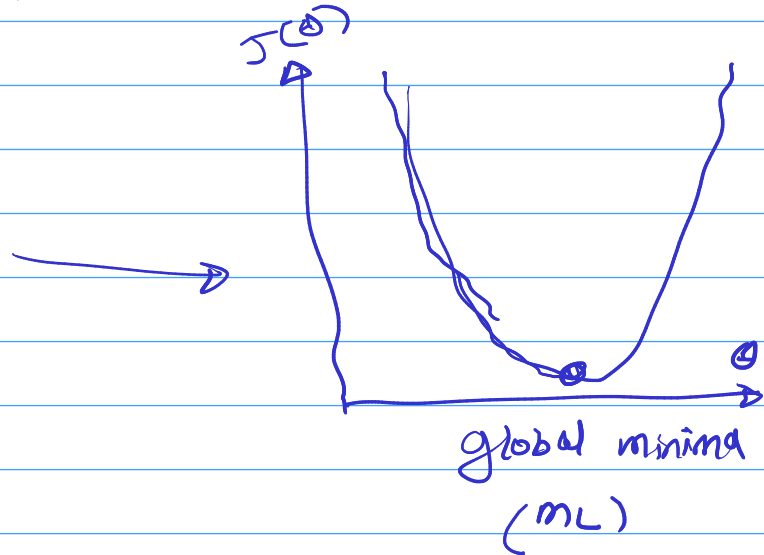
chain Rule = we find $\frac{dL}{dw_{11}}$ in chain rule

epoch \rightarrow 1 forward & 1 backward.

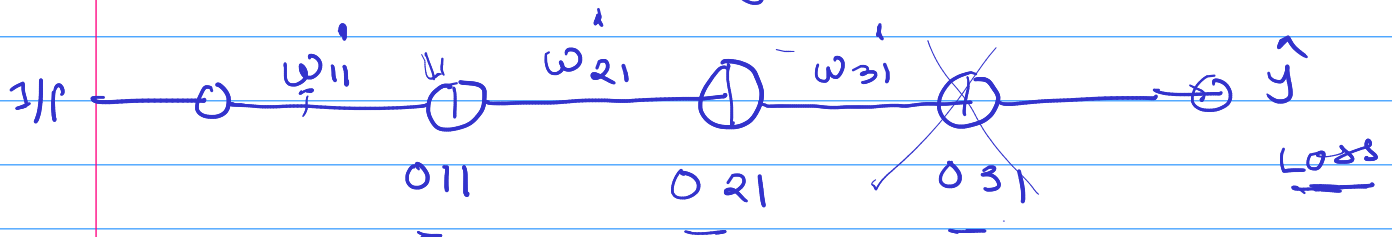
Optimizer (Reduce the loss)

1) Gradient Descent

y & \hat{y} similar



② Exploding Gradient Descent (Problem)
 \downarrow because of weight.



$$w'_{11, \text{new}} = w'_{11, \text{old}} - \frac{dL}{dw'_{11, \text{old}}}$$

475

500

$$\frac{dL}{dw'_{11, \text{old}}} = \frac{dL}{d031} \cdot \frac{d031}{d021} \cdot \frac{d021}{d11} \cdot \frac{d11}{dw'_{11}}$$

$$\text{Sigmoid} = \frac{1}{1 + e^{-x}} = 0 \text{ to } 1$$



$$= 0 \leq 0.25$$

$$\frac{dL}{d03}$$

ex
Loss = 500

$$= 0.25 \times 500$$

= 125 \rightarrow very high

old &
new.

③ Vanishing Gradient Problem

(when Rule

derivative of Sigmoid is always 0 to 0.25 → Exploding

so derivative always smaller

$$x \times 500 = 0$$

$$w_{\text{new}} \approx w_{\text{old}}$$

$$\underline{\text{Relu}} = (0 \text{ to } 2)$$

$$\frac{d}{dz}(\text{zero}) = \text{not define.}$$

mean = dead neuron

$$\text{Sig} = 0 \text{ to } 1$$

$$\text{tanh} = -1 \text{ to } 1$$

$$\text{Relu} = 0 \text{ to } \infty$$

$$\text{Leaky Relu} = 0.01 \text{ to } 2$$

$$\text{Softmax} = \text{mult}$$

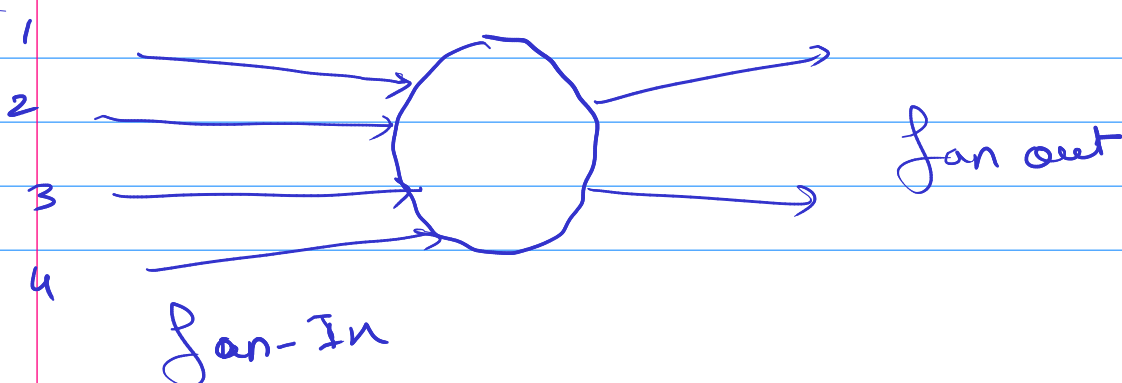
— x —

Weight Initialization

① weight should be small ✓

② weight should not be same ✓

③ should have good variance. ✓



Uniform Distribution

$$w_{ij} \sim \text{Uniform} \left[\frac{-1}{\sqrt{f_{\text{in}} - 1}}, \frac{1}{\sqrt{f_{\text{in}} - 1}} \right]$$

② Xavier / Glorot Distribution

$$w_{ij} \sim N(0, \sigma) \quad \begin{array}{l} \text{std. deviation} \\ \text{mean} \end{array} \quad \text{Normal}$$

(Sigma) $\sigma = \sqrt{\frac{2}{(f_{\text{in}} + f_{\text{out}})}}$

$$w_{ij} \sim U \left[\frac{-\sqrt{6}}{\sqrt{f_{\text{in}} + f_{\text{out}}}}, \frac{\sqrt{6}}{\sqrt{f_{\text{in}} + f_{\text{out}}}} \right]$$

Six $\sqrt{6}$

Uniform

④ He init

He uniform

$$w_{ij} \sim U \left[-\sqrt{\frac{6}{f_{\text{in}}}}, \sqrt{\frac{6}{f_{\text{in}}}} \right]$$

He init Normal

$$w_{ij} \sim N(0, \sigma) \quad \begin{array}{l} \text{sigma} \\ \text{sigma} = \sqrt{\frac{2}{f_{\text{in}}}} \end{array}$$

③ SGD (Stochastic Gradient Descent)

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$$

$\frac{dL}{dw_{\text{old}}}$ → When we consider all data point called as Gradient descent optimizer.

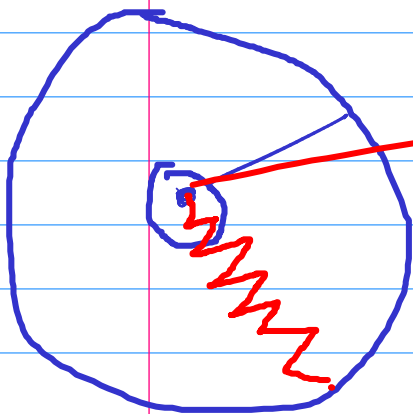
$$\Downarrow$$

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

SGD → Consider only 1 point

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

→ it take time to reach global miniming



mini batch gradient decent = k is batch size

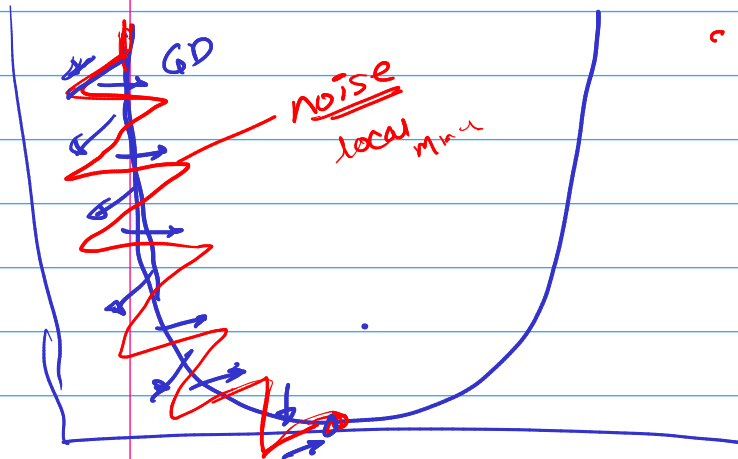
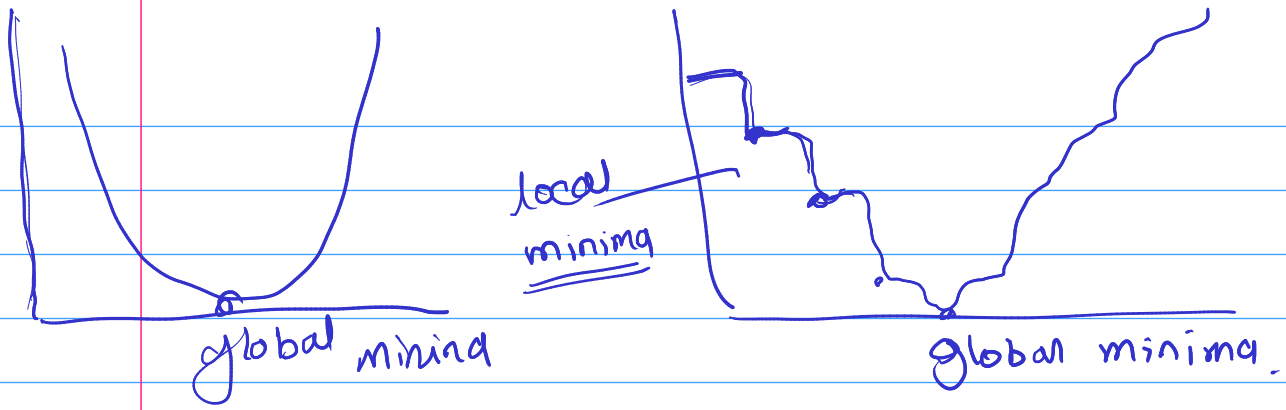
k < n n- data point
k batch

dL/dw_{old} (mini batch) \approx Gradient Decent

Loss in
mini batch

$$Loss = \sum_{i=1}^k (y - \hat{y})^2$$

↑
neg-mean - Square - error.



• 1 batch in

④ SGD with momentum →

Since $d_1, d_2, d_3, \dots, d_n \rightarrow$ momentum add in loss function

noise Reduce

↓
minimize loss not possible by SGD.

⑤ Adagrad optimizer (minimize the loss)

$$\omega_{\text{new}} = \omega_{\text{old}} - \eta \frac{dL}{d\omega_{\text{old}}}$$

$$W_{t, \text{new}} = W_{t-1} - \eta \frac{dL}{dW_{t-1}}$$

Can we use different learning rate for each & every neuron of hidden layer.

why we required?

in ANN \rightarrow 2 factor DENSE & SPARSE

DENSE \rightarrow most factor is non zero

SPARSE \rightarrow most factor is zero.

for Dense we use different learning rate
Sparse layer we use different learning rate for each.

$$\eta_t = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$$

epsilon always small. use for not to zero

Constant Learning rate

$$\alpha_t = \epsilon \sum_{i=1}^t \left[\frac{dL}{dW_i} \right]^2$$

$t =$ time / epoch

$\epsilon =$

Disadvantage = When iteration \uparrow α_t is also high \uparrow
 η is very small

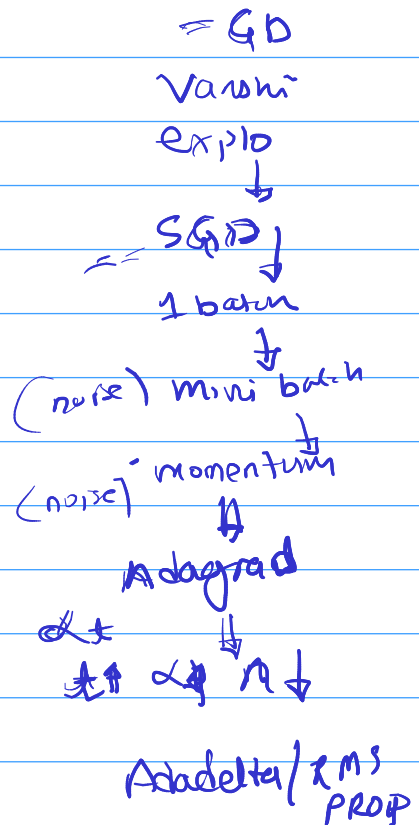
⑥ Adadelta & RMSPROP

$$\eta_i = \eta / \sqrt{w_{Arg_t} + \epsilon}$$

⑦ Adam

(Adaptive movement estimation)

Momentum \neq RMSPROP



$$\eta_i = \frac{\eta}{\sqrt{w_{Arg_t} + \epsilon}}$$

$\neq M \leftarrow \text{momentum}$

Loss function (cost function)

$$1$$
$$\text{Loss} = (y - \hat{y})^2$$

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

ANN - Tabular data

Regression

Square error loss

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

< neg-mean-sq-error.

Absolute

$$L = |y - \hat{y}|$$

$$\text{RMS} = \sqrt{|y - \hat{y}|^2}$$

classification
entropy

(cross-entropy)

(binary-cross-entropy)

✓ \hat{y} = Sigmoid for binary

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

multiclass-cross-entropy

one-hot encoding

Categorical Cross Entropy

Softmax

Sparse Categorical Cross Entropy
Label-encoding.

ANN