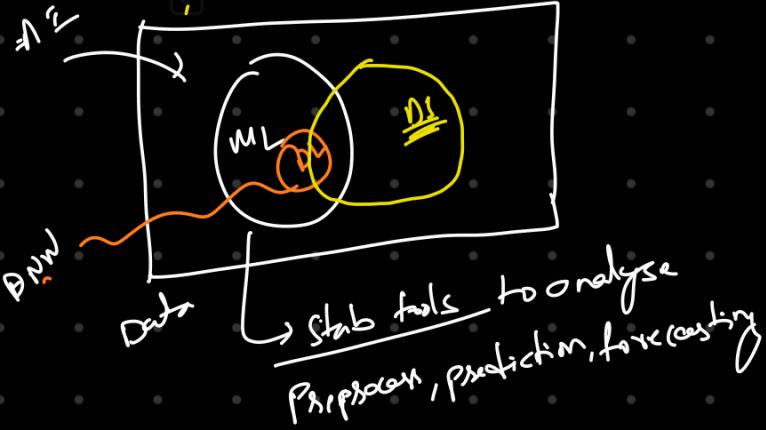


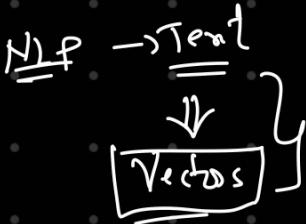
AI VS ML VS DL VS DS



AI application

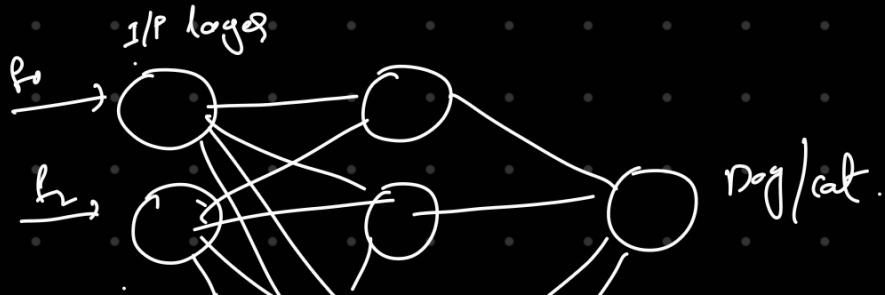
- ① NLP
- ② Self Driving
- ③ Amazon
- ④ YouTube

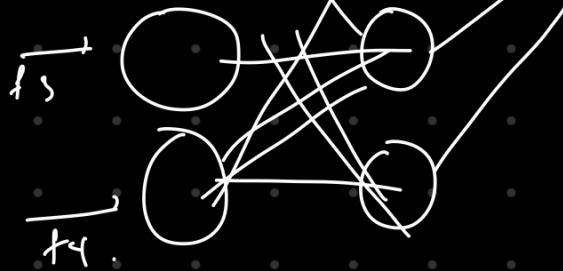
- ① Deep learning
if ANN \leftarrow Artificial Neural Network \leftarrow Tabular kind of dataset.
- ② CNN \leftarrow Convolutional neural network \leftarrow Images, Videos.
- a) Image classification \leftarrow CNN, Transfer Learning
 - b) Object detection \leftarrow RCNN, fAST RCNN, fASTER RCNN, SSD, Yolo, Detectron
 - c) Segmentation
 - d) Tracking
 - e) GAN
- ③ RNN \leftarrow Recurrent Neural Network \leftarrow I/P \div Text, Time series, sequential Data
- \rightarrow RNN
 - \rightarrow LSTM RNN
 - \rightarrow Bidirectional LSTM NN
 - \rightarrow Encoder Decoder
 - \rightarrow Transformer
 - \rightarrow BERT
 - \rightarrow GPT, GPT2, GPT3.



\rightarrow first N/W \rightarrow Perceptron \rightarrow not effective

\rightarrow Jeffrey Hinton \rightarrow proposed backtracking \rightarrow become effective.

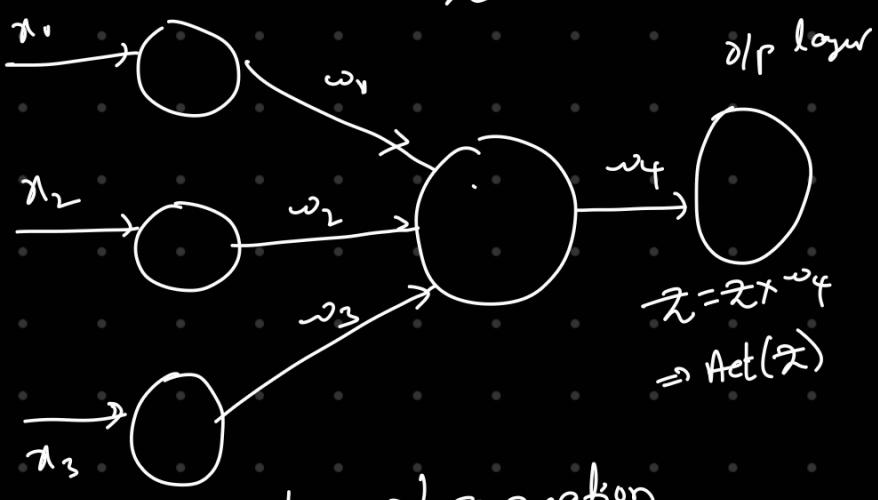




How ReLU works?

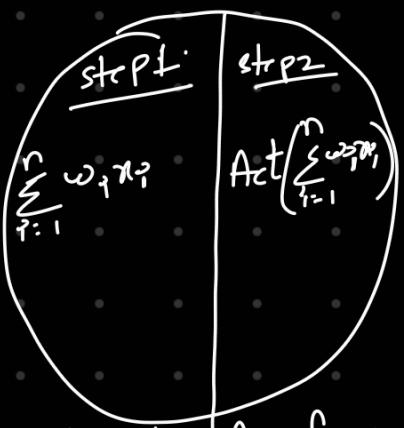
$$y = w_1x_1 + w_2x_2 + w_3x_3 + \text{bias}$$

$$\bar{x} = \text{Act}(y)$$



forward propagation.

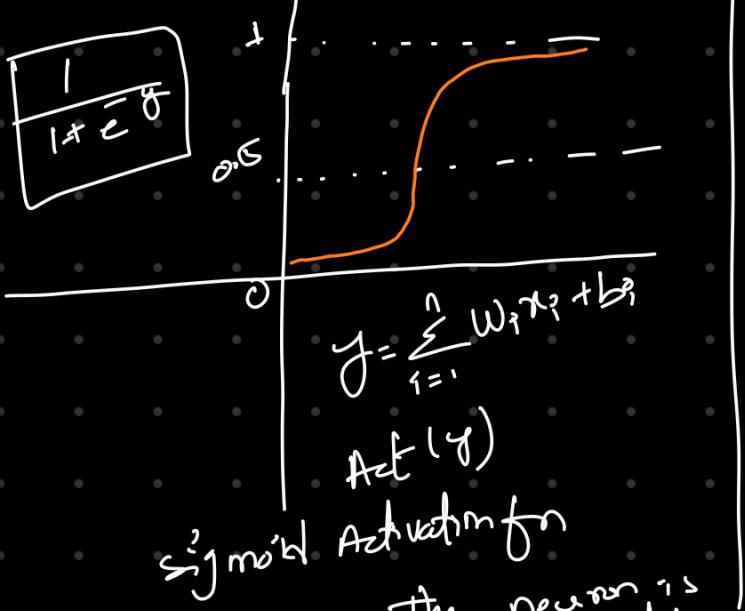
+ Heaviside



Sigmoid Activation fn =

$$\frac{1}{1 + e^{-x}}$$

Activation functions



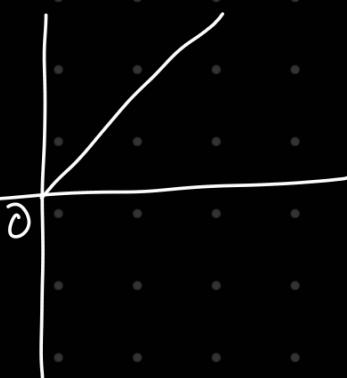
Sigmoid Activation fn

$$y = \sum_{i=1}^n w_i x_i + b$$

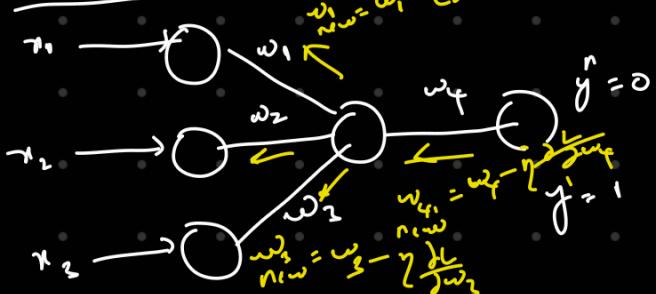
$$\text{Act}(y)$$

(1) ReLU AF

$$\max(y, 0)$$



→ output 1 means the neuron is activated.

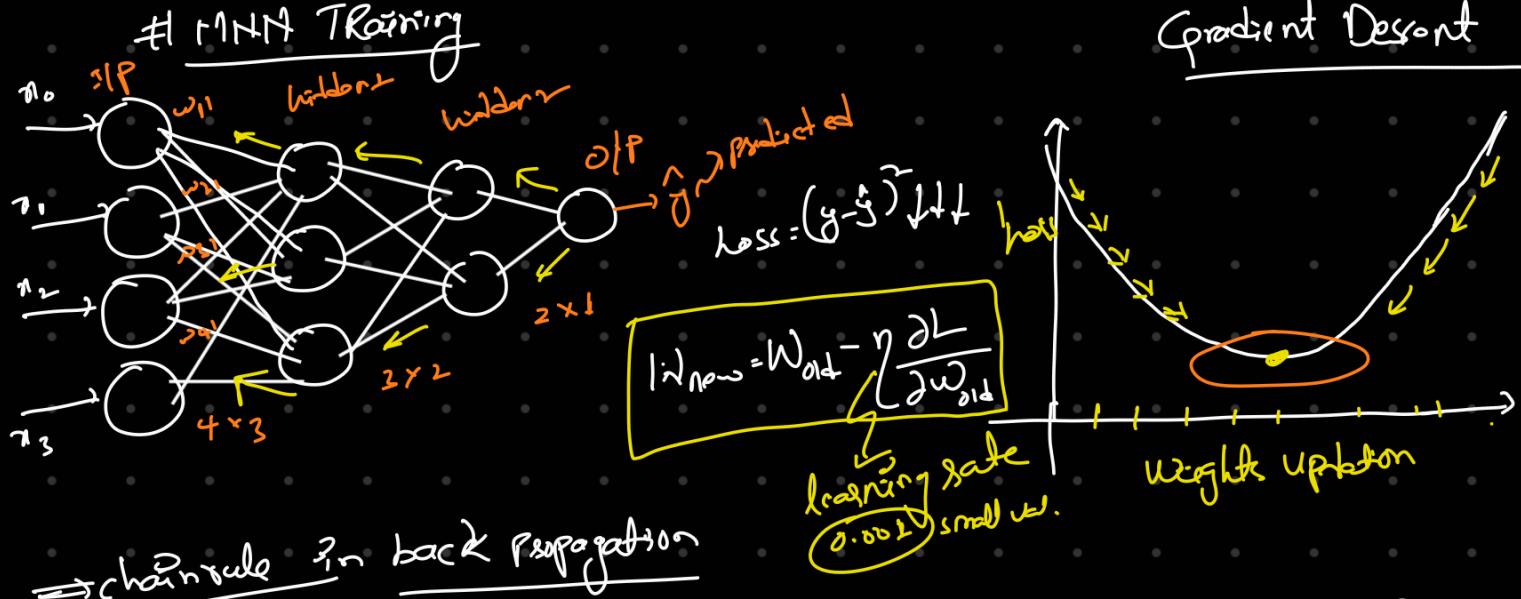


$$\text{loss} = (y - \hat{y})^2 \downarrow \text{need to decrease this}$$

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

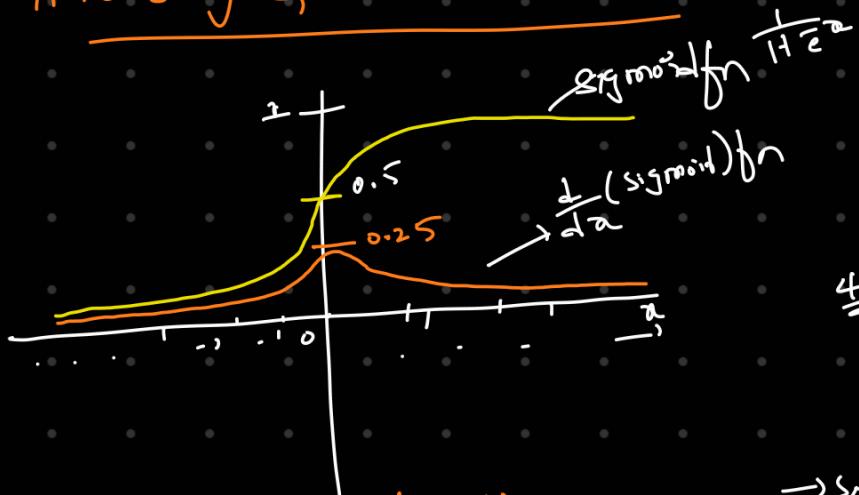
$$\text{cost fn} = \sum_{i=1}^n (y - \hat{y})^2$$

minimize losses → Optimizer.



$$w_{i,j}^{new} = w_{i,j}^{old} - \eta \frac{\partial L}{\partial w_{i,j}^{old}} \rightarrow \frac{\partial L}{\partial w_{i,j}^{old}} = \frac{\partial L}{\partial O_{j,i}} \times \frac{\partial O_{j,i}}{\partial O_{i,i}} + \frac{\partial L}{\partial O_{j,i}} \times \frac{\partial O_{j,i}}{\partial O_{i,i}} + \dots$$

Vanishing Gradient Problem



$$w_{i,j}^{new} = w_{i,j}^{old} - \eta \frac{\partial L}{\partial w_{i,j}^{old}}$$

$$\frac{\partial L}{\partial w_{i,j}^{old}} = \frac{\partial L}{\partial O_{j,i}} \times \frac{\partial O_{j,i}}{\partial w_{i,j}^{old}} \rightarrow \text{chain rule}$$

$$4 \text{ layers: } 0.20 \times 0.1 \times 0.05 \approx 0.001$$

$$= 0.2 \times 1 \times 0.04 = 0.08$$

$$= 2.46 \times \dots \approx 2.49$$

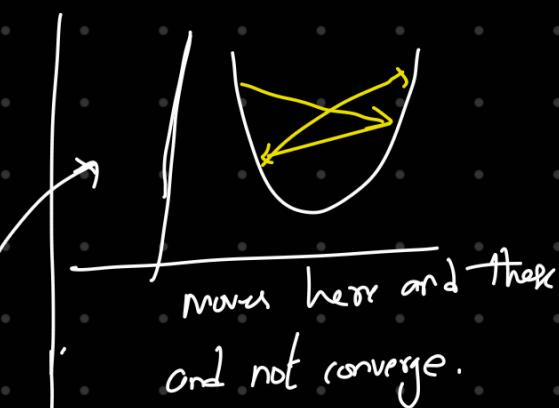
$w_{i,j}^{new} \approx w_{i,j}^{old}$ weights not getting updated.

Exploding Gradient Problem

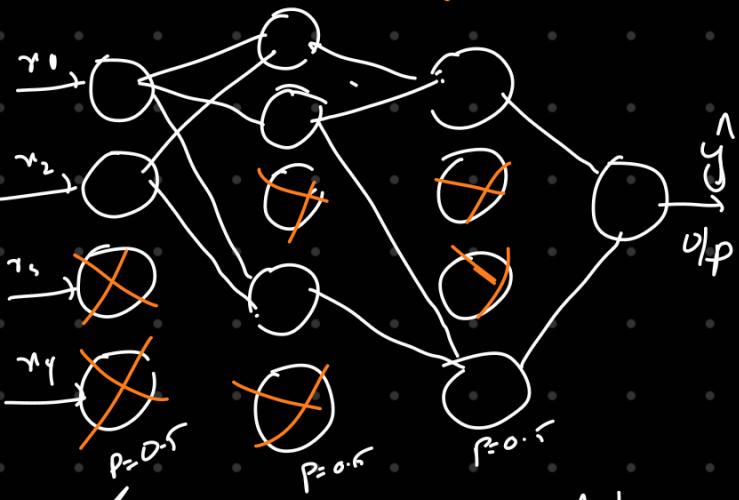


$$\begin{aligned} \frac{\partial O_{21}}{\partial O_{11}} &= \frac{\partial \phi(z)}{\partial z} \times \frac{\partial z}{\partial O_{11}} = \frac{\partial (\phi_1 \cdot z_1 + b_1)}{\partial O_{11}} \\ &\Rightarrow [0 \leq \dots \leq 0.25] * w_{21} \\ &\Rightarrow 0.25 \times 500 \sim 125 \end{aligned}$$

If wts are initialized with higher values, then only exploding gradient problem happen.



Drop Out & Regularization



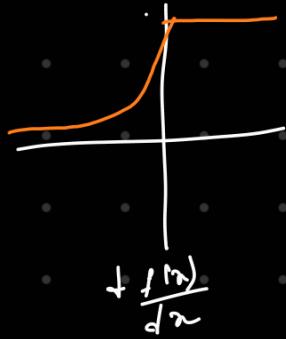
\Leftrightarrow % of neurons are deactivated.

Dropout rate

$$0 \leq p \leq 1$$

\rightarrow to reduce overfitting.

ELU {Exp Linear Units}



$$f(z) = \begin{cases} z, & z > 0 \\ \alpha(e^z - 1), & z \leq 0. \end{cases}$$

\rightarrow advantage is computationally exp.

PReLU (Parametric ReLU)

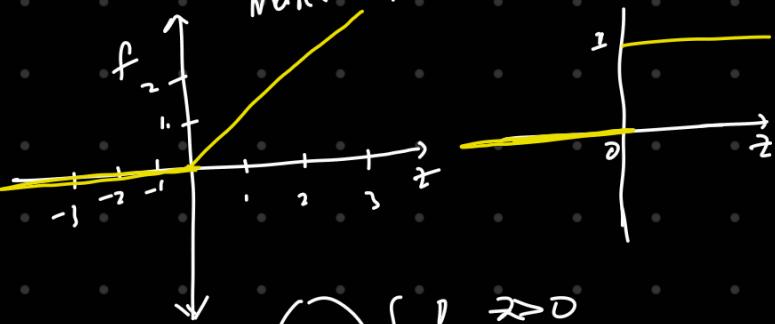
$$f(z) = \begin{cases} z, & z > 0 \\ \alpha z, & z \leq 0. \end{cases}$$

learning param.

$\therefore \alpha = 0.01 \rightarrow$ leaky ReLU

$\alpha = 0 \rightarrow$ ReLU.

Rectified Linear Unit $\frac{\partial}{\partial z} = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0. \end{cases}$

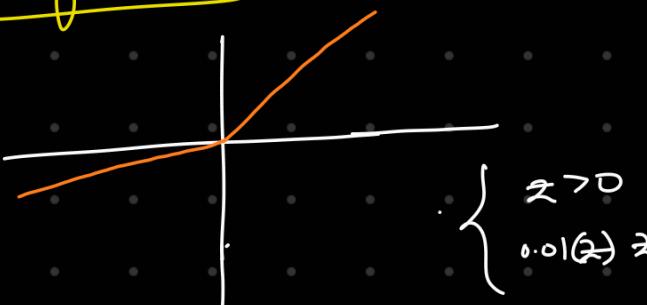


$$\text{In}_{\text{new}} = \text{In}_{\text{old}} - \eta \left(\frac{\partial}{\partial w_{\text{old}}} \right) \begin{cases} 1, & z > 0 \\ 0, & z \leq 0. \end{cases}$$

$$\text{In} = \text{In} - \eta \times 1 \times 1$$

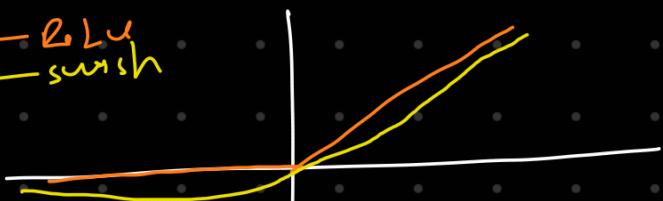
\rightarrow created dead neuron. during back propagation.

leaky Relu



Softplus (A soft-coded) Function

ReLU
softplus



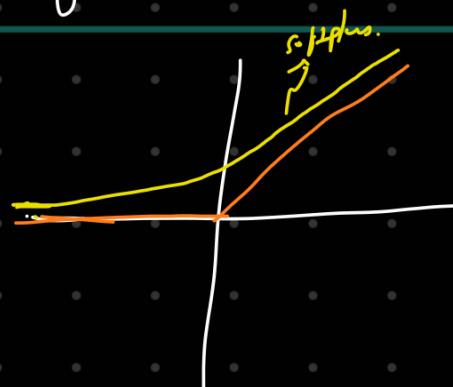
\rightarrow only works when layers are > 40 .

$$\rightarrow y = z \times \text{sigmoid}(z)$$

\rightarrow computationally exp.

Softplus

$$f(z) = \ln(1 + e^z)$$



Weight Initialization

- (1) weights should be small
- (2) weights should not be same
- (3) weights should have good variance

(2) Xavier / Glorot

\rightarrow Xavier Normal

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

$$\sigma = \sqrt{\frac{2}{\text{fan_in} + \text{fan_out}}}$$

\rightarrow Xavier Uniform

$$w_{ij} \sim U\left[\frac{-\sqrt{6}}{\sqrt{\text{fan_in} + \text{fan_out}}}, \frac{\sqrt{6}}{\sqrt{\text{fan_in} + \text{fan_out}}}\right]$$

(1) Uniform Distribution

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

(3) He init

(a) He Normal

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

$$\sigma = \sqrt{\frac{2}{\text{fan_in}}}$$

(b) He Uniform

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

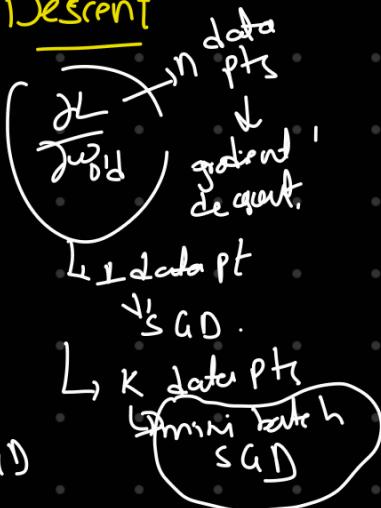
Stochastic Gradient Descent

$$w_{\text{new}} = w_{\text{old}} - \eta \times \frac{\partial L}{\partial w_{\text{old}}}$$

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

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

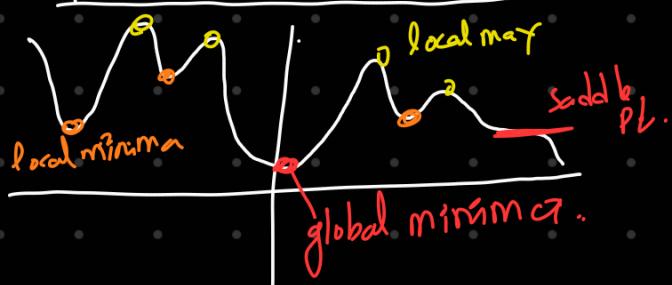
$$(y - \hat{y})^2 \rightarrow \text{SGD}$$



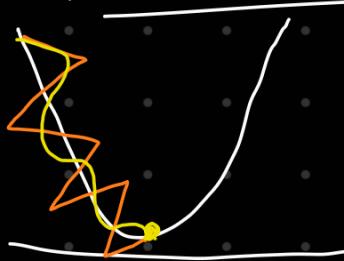
$$\left[\frac{\partial L}{\partial w_{\text{old}}} \right]_{\text{minibatch}} \approx \left[\frac{\partial L}{\partial w_{\text{old}}} \right]_{\text{population}}$$

sample

\Rightarrow Global Minima vs. Local Minima



SGD with momentum



$$\begin{aligned} t &= 1 & v_1 &= b_1 \\ v_2 &= \gamma * v_1 + L_2 \\ v_3 &= \gamma * v_2 + b_3 \end{aligned}$$

$$w_{\text{new}} = w_{\text{old}} - \eta \times \frac{\partial L}{\partial w_{\text{old}}}$$

$$\gamma = 0.9 = w_{\text{old}} - \left[\gamma v_{t-1} + \eta \times \frac{\partial L}{\partial w_{\text{old}}} \right]$$

$$v_{t-1} = \gamma \left[\frac{\partial L}{\partial w_{\text{old}}} \right]_{t-1} + \left[\frac{\partial L}{\partial w_{\text{old}}} \right]_{t-2} \dots$$

AdaGrad Optimizer

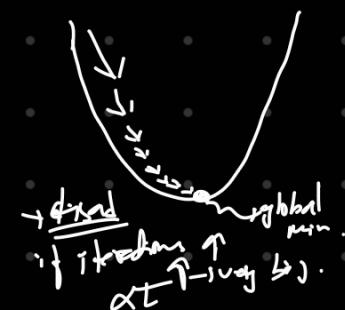
$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$$

Decreasing stepsize after time.

$$\eta_t = \eta / \sqrt{\alpha_t + \sum_{i=1}^{t-1} \text{val}_i^2}$$

$$\alpha_t = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_i} \right)^2$$

high freq no.



AdaDelta of RMS Prop

→ used to restrict η to go small beyond certain values.

$$\text{Adagrad } h_t = h_{t-1} - \eta \frac{\partial L}{\partial w_t}$$

$$\eta_t = \frac{\eta}{\sqrt{k_{avg} + \epsilon}}$$

$$k_{avg} = \beta \times k_{avg} + (1-\beta) \left(\frac{\partial L}{\partial w_t} \right)^2$$

$$k_{avg} = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_t} \right)^2$$

Adam Optimizer (adaptive Moment estimation)

\hookrightarrow momentum → smoothing
 \hookrightarrow RMS prop → variance

$S_w = P_1 h_{t-1} + (1-\beta) \left(\frac{\partial L}{\partial w} \right)^2$

$S_b = P_2 b_{t-1} + (1-\beta) \left(\frac{\partial L}{\partial b} \right)^2$

$RM > PR$

Algorithm

$v_{w,b} = 0$, $s_{w,b} = 0$, $\hat{s}_{w,b} = 0$

On iteration t

- (1) Compute $\frac{\partial L}{\partial w} + \frac{\partial L}{\partial b}$
- $v_w = \beta_1 v_w + (1-\beta_1) \frac{\partial L}{\partial w}$ momentum
- $v_b = \beta_1 v_b + (1-\beta_1) \frac{\partial L}{\partial b}$
- Adam opt

Loss functions

$$\Rightarrow \frac{MSE}{\Rightarrow \frac{1}{m} \sum (y - \hat{y})^2}$$

1) plot the quadratic eqn, we get a gradient descent with only global minima

2) we don't get any local minima

2) The MSE loss penalizes the model for making large errors by squaring them.
 \rightarrow Disadvantages

① It is not robust to outliers.

② Absolute Error Loss

① $L = |y - \hat{y}|$, $J = \frac{1}{t} \sum_{i=1}^t |y_i - \hat{y}_i|$ → imperturbational
 \rightarrow the model is more robust to outliers as compared to MSE

③ Huber Loss

① MSE ② MAE

Loss = $\begin{cases} \frac{1}{2} (y - \hat{y})^2, & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2, & \text{otherwise} \end{cases}$ ← hyperparam

↳ linear eqn.

→ Cross Entropy → Logistic Regression

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

Binary cross entropy, in classification

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

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

Multiclass Cross Entropy Loss

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

↳ categorical cross entropy

o/p → Multiclass → good, Bad, Neutral → one-hot encoded

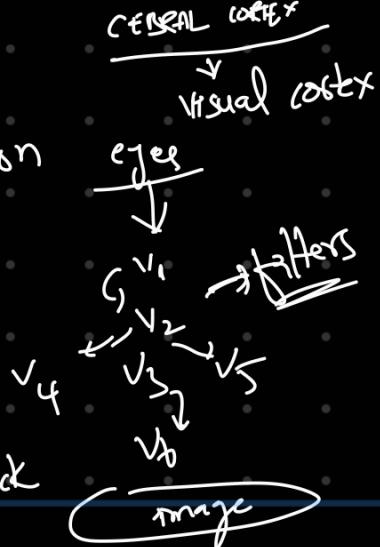
$$\begin{matrix} f_1 & f_2 & f_3 \\ 2 & 3 & 4 \end{matrix} \quad \text{o/p} \quad \begin{matrix} \text{good} \\ \text{Bad} \end{matrix} \rightarrow [1 \ 0 \ 0]$$

$$\begin{matrix} 4 & 5 & 6 \\ 7 & 8 & 9 \end{matrix} \quad \text{Neutral} \rightarrow [0 \ 0 \ 1]$$

$$\begin{matrix} f_1 & f_2 & f_3 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{matrix} \quad \begin{matrix} y_{11} & y_{12} & y_{13} \\ y_{21} & y_{22} & y_{23} \\ y_{31} & y_{32} & y_{33} \end{matrix} \quad \begin{matrix} \text{one hot encoded} \\ \rightarrow \text{softmax fn.} \\ e^{z_i} \\ \sum_{j=1}^3 e^{z_j} \end{matrix}$$

#CNN

① ANN \rightarrow regression
classification



② CNN

skate or Jack

0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1

$$\rightarrow n=6 \rightarrow 4 \times f$$

$f=3 \rightarrow n-f+1$
 $|n-f+1|$ changes for padding

With stride = 2

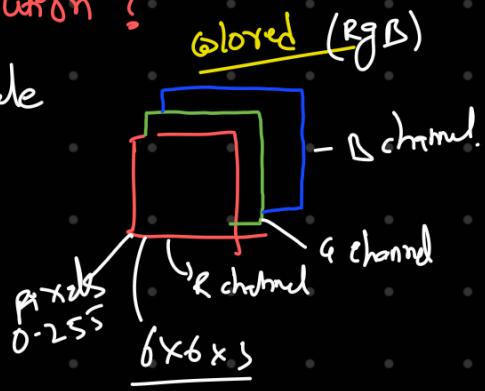
$$* \begin{array}{|c|c|c|} \hline +1 & 0 & -1 \\ \hline +2 & 0 & -2 \\ \hline +1 & 0 & -1 \\ \hline \end{array}$$

Vertical edge filter
(or $1+0x0+0x-1$) +
($0x+2+0x0+0x-2$) +
($1x+1+0x0+0x-1$) = 0

killed us convolution?

0 to 255		
0	0	0
255	255	255
255	255	255
0	0	0

Grey scale
 4×4



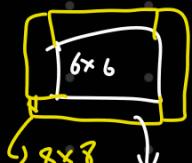
Pooling in CNN

→ like avg pooling some info when filter is applied.

$$8 \times 6 \xrightarrow[3 \times 3]{} 4 \times 4$$

$$n-f+1 = 6 \rightarrow \text{not loosing anything}$$

$n = 8$ so add extra edges



now

$$8 \times 8 \xrightarrow[3 \times 3]{} 4 \times 4$$

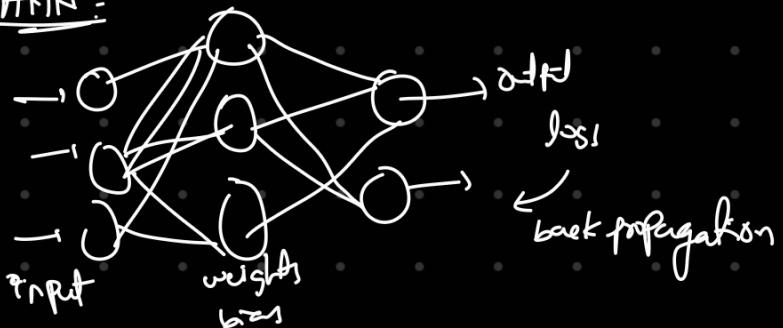
values = (1) 0's of padding
(2) near values

→ new formula:

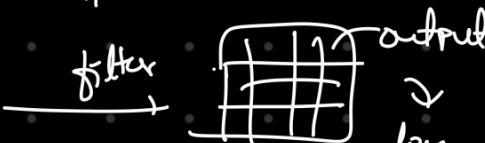
$$n-f+1+2p$$

Operation of CNN

ANN:



↓ Similarly



RNN
back propagation.