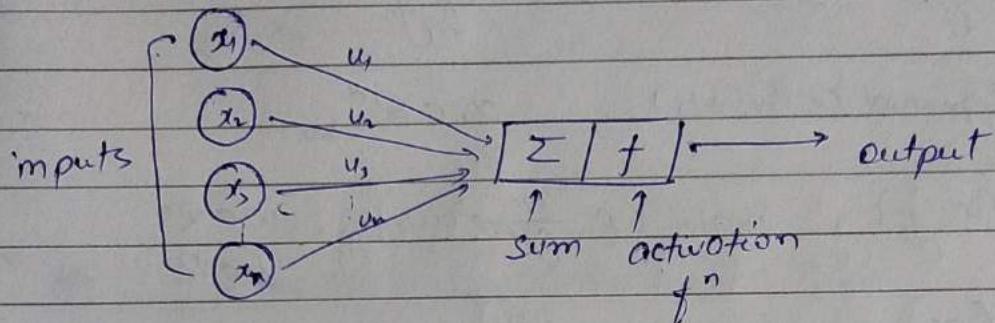


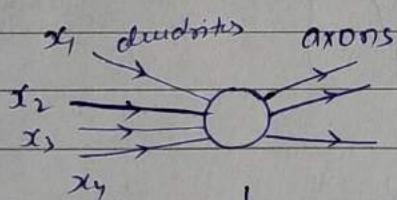
DEEP LEARNING

* History of Neural networks :-



* How Biological Neurons work?

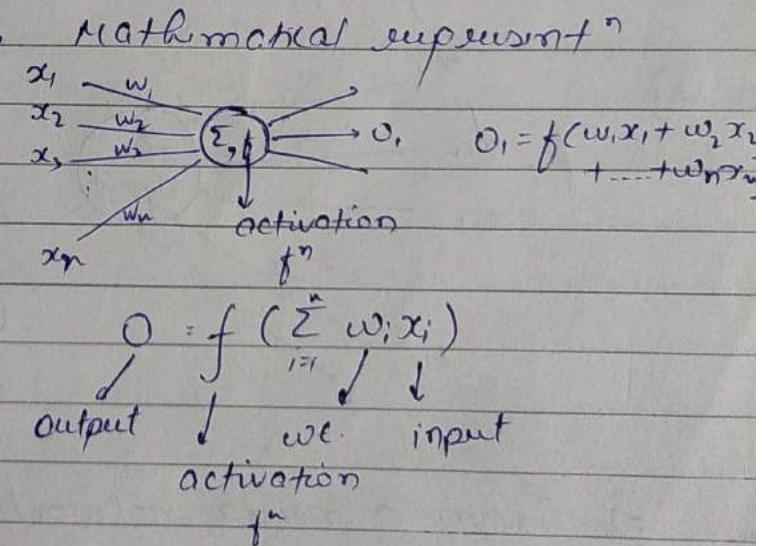
Simplified view of neuron.



— thicker edge \Rightarrow more wt.
— thinner edge \Rightarrow less wt.

* Growth of biological neural networks.

- formation of new branches b/w neurons
- thickness of branch depends on how much important it is.



④ Diagrammatic Representation: Logistic Regression and perceptron (single layered perceptron)

LR: $x_i \rightarrow \hat{y}_i \rightarrow \text{predicted value of } y_i$

$$\hat{y}_i = \text{Sigmoid}(\mathbf{w}^T \mathbf{x}_i + b) \quad \mathbf{x}_i \in \mathbb{R}^d$$

$$D = \{\mathbf{x}_i, y_i\} \quad \text{train LR} \xrightarrow{\text{we get}} \mathbf{w}, b \quad \mathbf{w} \in \mathbb{R}^d \quad b \in \mathbb{R}$$

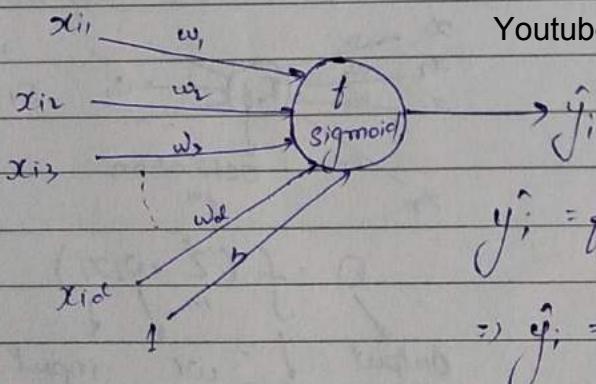
$$\mathbf{x}_i = [x_{i1}, x_{i2}, x_{i3}, x_{i4}, \dots, x_{id}]$$

$$\Rightarrow \hat{y}_i = \text{Sigmoid}\left(\sum_{j=1}^d w_j x_{ij} + b\right) \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{LR}$$

$$\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}]$$

$$\mathbf{w}_i = [w_1, w_2, \dots, w_d]$$

$$\Rightarrow D = f\left(\sum_{j=1}^d w_j x_{ij}\right) \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{NN}$$



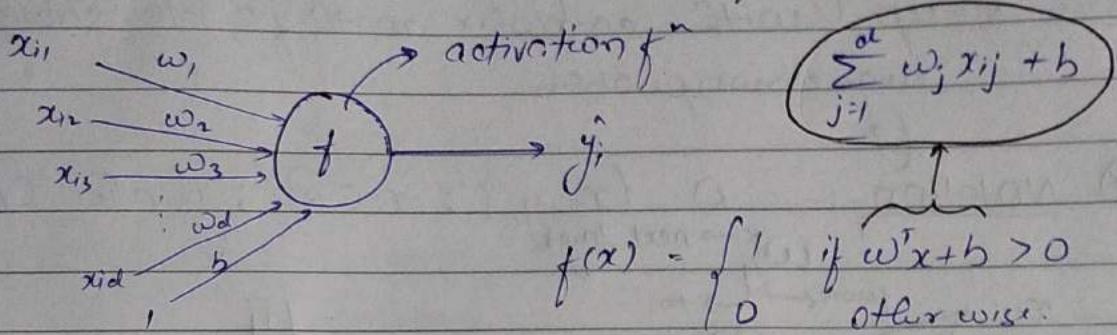
Youtube: Programming Cradle

$$\begin{aligned} \hat{y}_i &= f(w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id} + b) \\ \Rightarrow \hat{y}_i &= f(\mathbf{w}^T \mathbf{x}_i + b) \end{aligned}$$

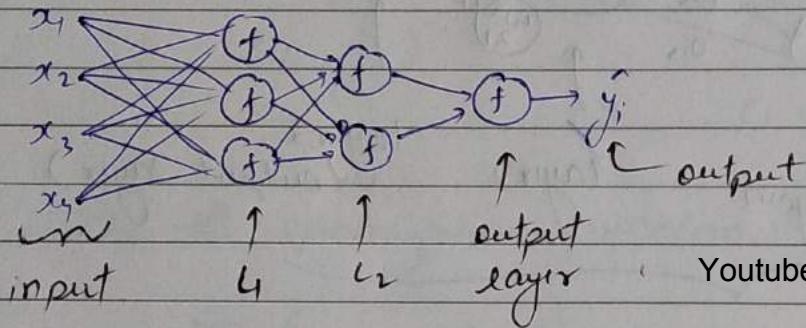
\Rightarrow Train a neural network means finding weights on edges / vertices

\Rightarrow Only diff' b/w perceptron and LR is the loss fn
in LR we do squashing using Sigmoid fn
whereas in perceptron we don't do squashing.

④ Perceptron (1957) → linear classifier



⑤ Multi-layered Perceptron (MLP)



Youtube: Programming Cradle

Q) Why should we care about MLP?

• Biological inspiration : biological neurons are connected in multiple layers (millions)

• Mathematical argument :-

example :- $2 * \sin(x^2) + \sqrt{x^5}$

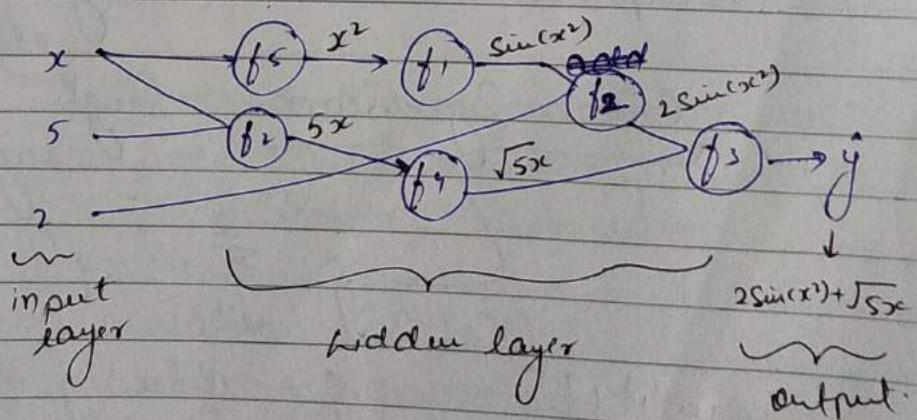
$$f_1 = \sin(x)$$

$$f_2 = \text{mult}()$$

$$f_3 = \text{add}()$$

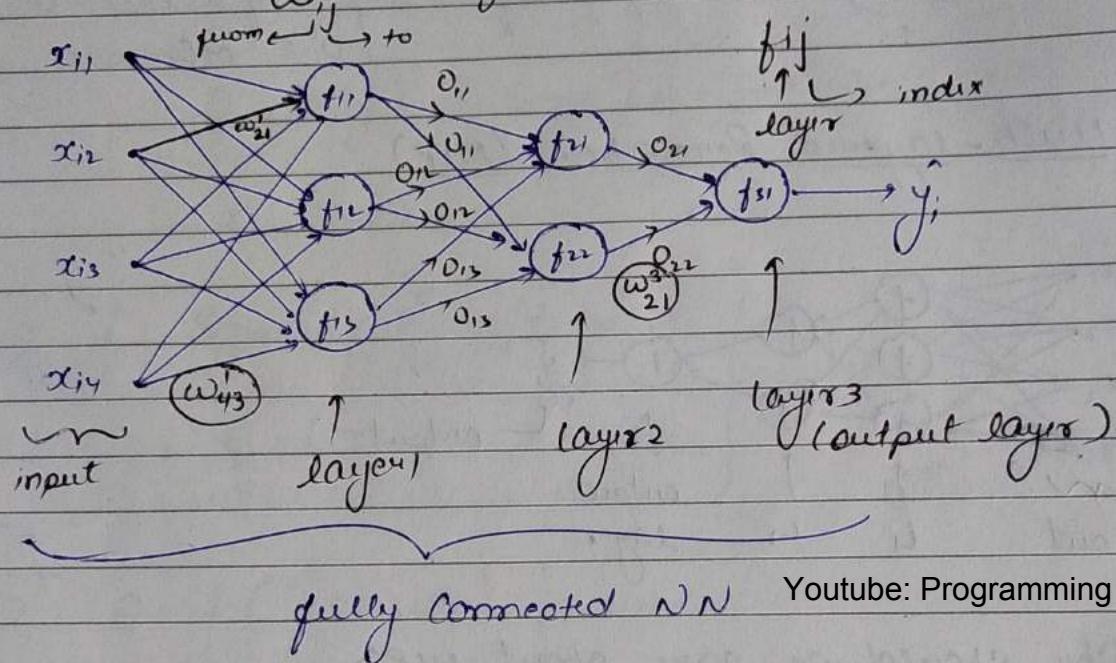
$$f_4 = \text{sqrt}()$$

$$f_5 = \text{squared}()$$



* By using multi layered structure we can come up with complex math fn to solve your regression problem.

* Notation :- $D = \{x_i, y_i\}; x_i \in \mathbb{R}^n; y_i \in \mathbb{R}$ (Regression)



Youtube: Programming Cradle

* Train a Single Neuron model :-

↳ (find the best edge weights using D_{10})

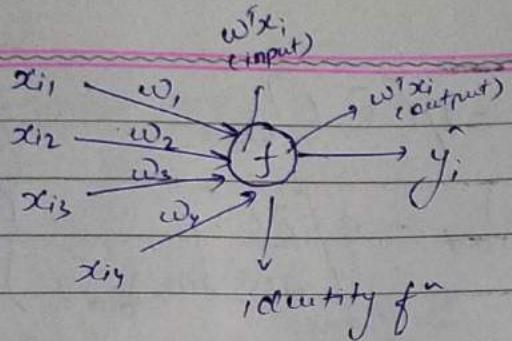
Perception and CR → Single Neuron models for classifⁿ

Linear Regression → Single Neuron models for regression

$$\hat{y}_i = \sum_{j=1}^d w_j x_{ij} \quad \left. \begin{array}{l} \\ y_i = w^T x_i \end{array} \right\} \text{Linear reg. opt}^n$$

identity f^n

$$f(x) = x$$



[in logistic regn
(f - sigmoid f)]

Linear Regression optimization :-

① Define loss fn $\min_{\mathbf{w}; i=1} \sum_{i=1}^{n \rightarrow \text{no. of points.}} (y_i - \hat{y}_i)^2 + \text{reg}$ $\|\mathbf{w}\|_1$ or $\|\mathbf{w}\|_2^2$

$$\hat{y}_i = \mathbf{w}^T \mathbf{x}_i$$

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

\hookrightarrow Sum of all the points

$i \rightarrow$ for a particular point.

② Write the optimization problem.

$$\min_{\mathbf{w}; i=1} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \text{reg} \quad \text{--- (1)}$$

$$\hat{y}_i = f(\mathbf{w}^T \mathbf{x}_i)$$

Youtube: Programming Cradle

\Rightarrow writing eqn (1) in general form.

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - f(\mathbf{w}^T \mathbf{x}_i))^2 + \text{reg.}$$

$f \rightarrow$ identity fn for linear regn

$f \rightarrow$ logistic / sigmoid fn for LR

$f \rightarrow$ threshold fn \rightarrow prediction.

③ Solve the optimization problem (SGD)

a) initialization of w_j 's \rightarrow random.

b) $\nabla_{\mathbf{w}} L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_n} \end{bmatrix}$

$w \in \mathbb{R}^n$

$x_i \in \mathbb{R}^n$

$$\textcircled{c} \quad w_{\text{new}} = w_{\text{old}} - \eta [V_w] w_{\text{old}}$$

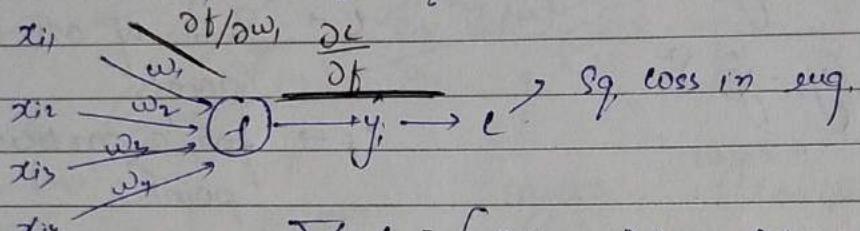
$$(w_i)_{\text{new}} = (w_i)_{\text{old}} - \eta \left[\frac{\partial L}{\partial w_i} \right] (w_i)_{\text{old}}$$

} updating weights.
update grad

GD: $\nabla_w L \rightarrow x_i$'s and y_i 's

SGD: $\nabla_w L \approx$ one pt $\{x_i, y_i\} \leftarrow$
small batch of pts \leftarrow batch SGD

Q) How to compute $\nabla_w L$?



$$\nabla_w L = \left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial w_3}, \frac{\partial L}{\partial w_4} \right]^T$$

$$\frac{\partial L_i}{\partial w_i} = \frac{\partial L_i}{\partial f} \cdot \frac{\partial f}{\partial w_i} \rightarrow \text{chain rule of diff}$$

$$\Rightarrow L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{(reg.)} \rightarrow \text{ignore for now}$$

$$= \sum_{i=1}^n (y_i - f(w^T x_i))^2$$

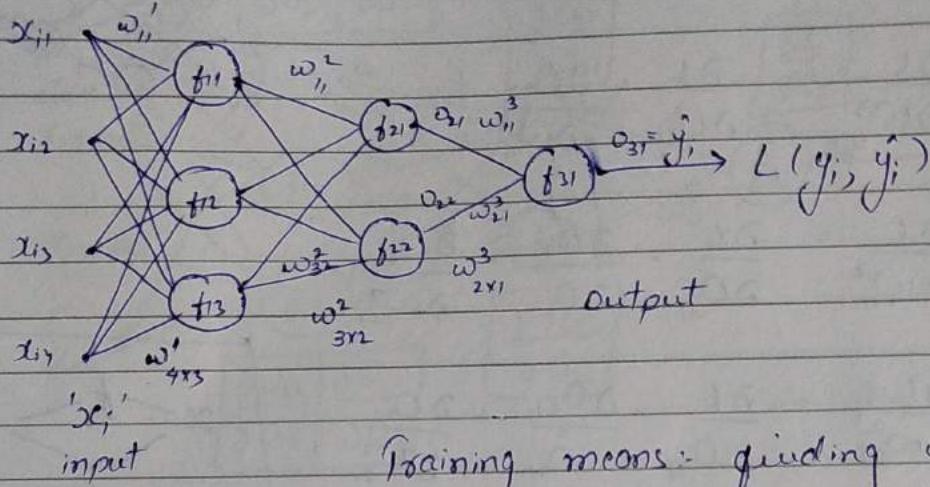
$$\frac{\partial L_i}{\partial f(w^T x_i)} = \frac{\partial L_i}{\partial f} = -\sum_{i=1}^n 2(y_i - f(w^T x_i))$$

Youtube: Programming Cradle

$$\frac{\partial L_i}{\partial w_i} = -2x_i(y_i - \hat{y}_i) : f \rightarrow \text{identity } f(x) = x \\ f(w^T x_i) : w_i x_i$$

$$\frac{\partial L}{\partial w_i} = \sum_{i=1}^n -2x_i(y_i - \hat{y}_i) = -\sum_{i=1}^n 2x_i(y_i - f(w^T x_i))$$

✳️ Training MLP



Training means: finding $\underbrace{w_{11}^1, w_{12}^1, w_{13}^1}_{\text{weights}}, \underbrace{w_{21}^1, w_{22}^1, w_{23}^1}_{\text{weights}}, \underbrace{w_{31}^1}_{\text{weights}}$

① Define loss f^*

$$\Rightarrow L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{(reg)} \rightarrow \sum_{i,j,k} (w_{ij}^k)^2$$

or

$$\Rightarrow L_i = (y_i - \hat{y}_i)^2 \Rightarrow L = \sum_{i=1}^n L_i + \text{reg.} \quad \frac{\sum_{i,j,k} |w_{ij}^k|}{C_{ik}}$$

② Optimizn :-

$$\min_{\substack{w_1, w_2, w_3 \\ w_{ij}^k}} L \rightarrow \text{sq. loss + regul'}$$

Youtube: Programming Cradle

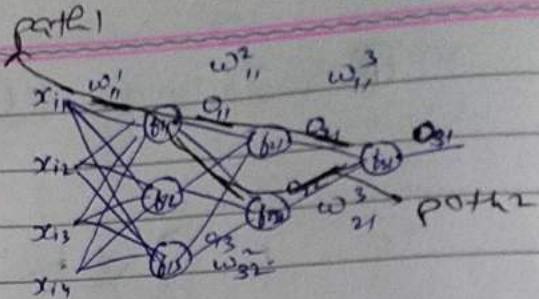
③ SGD or GD

a) initialize of w_{ij}^k randomly.

update eq^a \rightarrow b) $(w_{ij}^k)_{\text{old}} = (w_{ij}^k)_{\text{old}} - \eta \frac{\partial L}{\partial w_{ij}^k}$

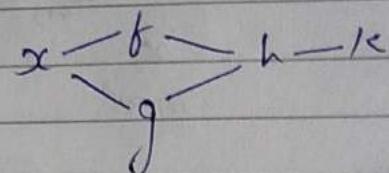
c) perform update till convergence. learning rate

$$\left. \begin{array}{l} \frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial w_{21}^3} \\ \frac{\partial L}{\partial w_{31}^3} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial w_{31}^3} \end{array} \right\}$$



$$\left. \begin{array}{l} \frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}^2} \end{array} \right\}$$

$$\left. \begin{array}{l} \frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{21}^2} \end{array} \right\}$$



$$\left. \begin{array}{l} \frac{\partial L}{\partial w_{31}^2} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{31}^2} \end{array} \right\}$$

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial o_{31}} \cdot \underbrace{\left(\frac{\partial o_{31}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}^1} + \right)}_{\text{path1}}$$

$$\underbrace{\left(\frac{\partial o_{31}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}^1} \right)}_{\text{path2}}$$

$$\frac{\partial k}{\partial x} = \frac{\partial k}{\partial h} \cdot \frac{\partial h}{\partial x}$$

$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x}$$

$$\frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}$$

$$\frac{\partial k}{\partial x} = \frac{\partial k}{\partial h} \left(\frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x} \right)$$

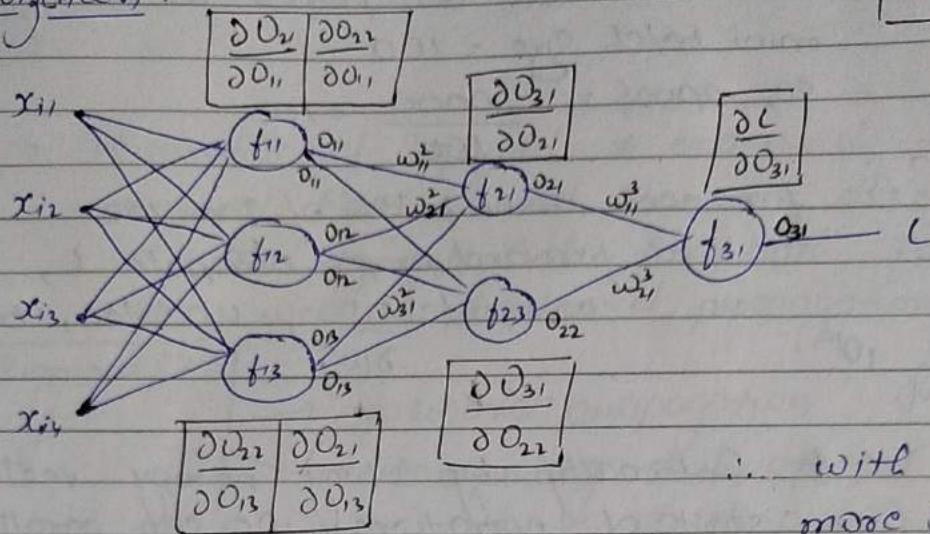
④ Training MLP memoization

memoization :-

if there is any operat'n that is used many times repeatedly, it is a good idea to compute it once and store it.

Epoch :- no. of times whole data is passed through the network

Memoization :-



... with slightly more space, speed can be ↑ drastically.

④ Backpropagation :- chain rule + memoization

$$D = \{x_i, y_i\}$$

① initialize w_{ij}^k 's

② for each x_i in D

a) pass x_i forward through the netw: \rightarrow forward prop

b) Compute $L(\hat{y}_i, y_i)$

c) Compute all the derivatives using chain-rule and memoization

d) update weights from end of the network to start : back propagation

③ Repeat step 2 till convergence

→ Instead of sending one point at a time to the network (SGD) we can send batch of points (mini batch SGD)

→ Keeping all the datapoints in RAM and computing derivatives (d) is hard (next to impossible if data is huge) } GD

→ To overcome above problem mini-batch based back prop is used most of the times

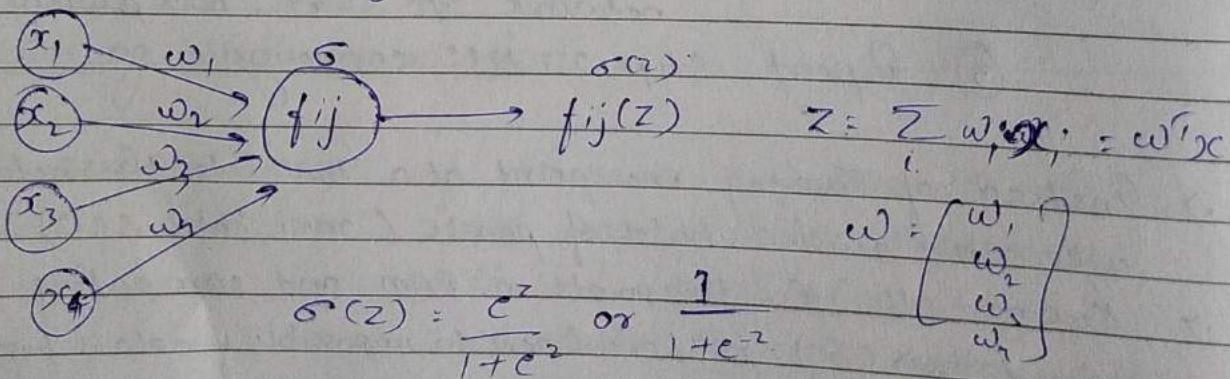
Let there are 10K datapoints in D and
 mini batch size = 100
 so, epoch = $\frac{10000}{100} = 100$

We have to sum this only 100 times instead of 10K
 \Rightarrow for each mini-batch of size 100
 → forward prop, compute L,
 compute $\frac{\partial L}{\partial w_{ij}}$, update, back prop.

- ④ Internally by using numpy vectorizing sort of operations, we can parallelly send multiple inputs to the nw. think of them as a matrix with datapoints one stacked on top of another. While calc. the gradient, gradient for all the pts. is calculated and avg. value is taken which is used for updating the wt.

Youtube: Programming Cradle

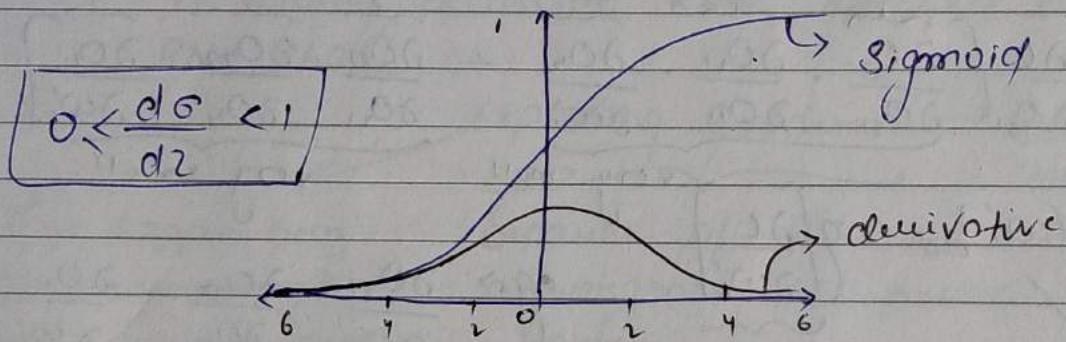
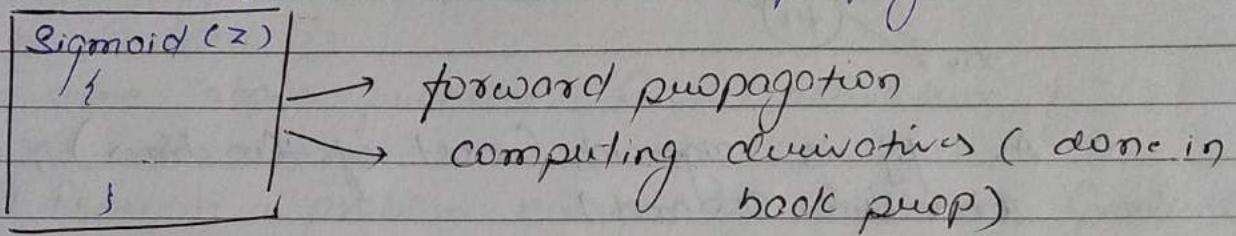
- ⑤ Activation function :- f_{ij} (Sigmoid, tanh) \rightarrow population $(80, 90)$
 activation f' should be differentiable and easy to differentiate.



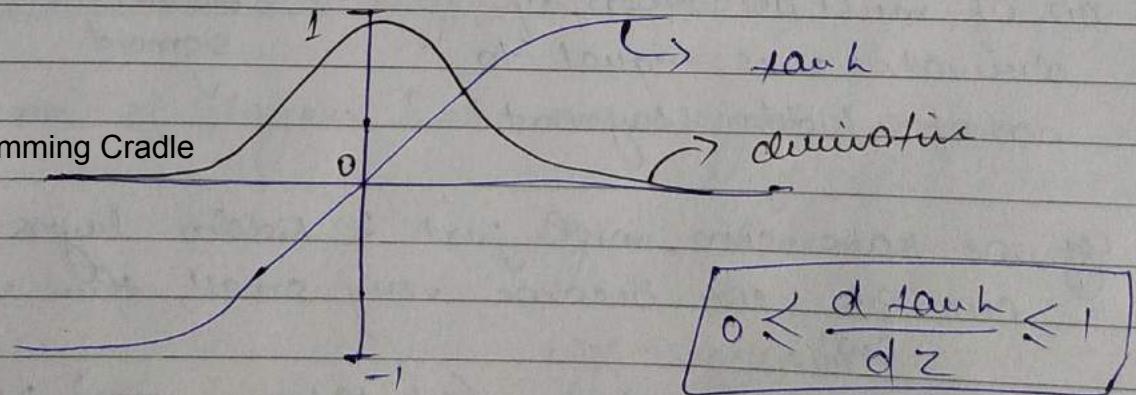
$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

$$\frac{d\sigma}{dz} = \sigma(z)(1-\sigma(z))$$

→ derivative of Sigmoid ^{can be} is expressed in terms of sigmoid itself and hence we can use sigmoid calc. for both forward as well as backward propagation.

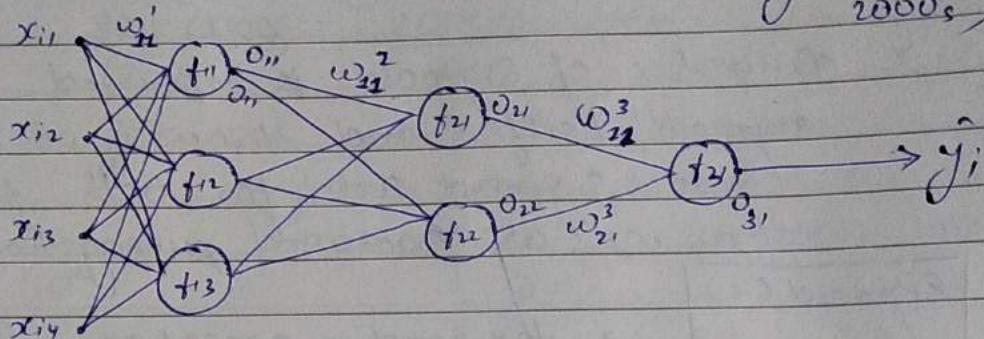


* $\tanh = \frac{e^z - e^{-z}}{e^z + e^{-z}} = a \Rightarrow \frac{da}{dz} = 1 - a^2$



Modern activation f': ReLU

⊗ Vanishing gradients :- (Biggest problem in NN during 80's, 90's and 2000's)



→ As f_{ij} is sigmoid (most of the times), value lies b/w 0 and 1

$$\frac{\partial C}{\partial w_{ii}} = \frac{\partial C}{\partial O_{31}} \left[\underbrace{\frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{ii}}}_{\text{very small}} + \underbrace{\frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{ii}}}_{\text{very small}} \right]$$

$$(w'')_{\text{new}} = (w'')_{\text{old}} - \eta \left[\frac{\partial C}{\partial w_{ii}} \right]$$

will be almost
same

e.g. $\frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{ii}}$

↓ ↓ ↓

$0.2 * 0.1 + 0.05 = 0.005$

Please all small
cause etc. because of
sigmoid

⊗ no. of multiplication of derivatives is equal to no. of hidden layers + 1

⊗ If we have sun, with just 2 hidden layers (3 multip) our $\frac{\partial C}{\partial w_{ii}}$ can become very small, then with more hidden layers (let 10 or even more) $\frac{\partial C}{\partial w_{ii}}$ will tend to 0.

⊗ As $\frac{\partial C}{\partial w_{ii}}$ getting smaller and smaller it is called a) Vanishing gradient

* Sigmoid / Tanh $\rightarrow 0 < < 1 \rightarrow$ chain rule and multiplication

$$(w_{ij}^k)_{\text{new}} \approx (w_{ij}^k)_{\text{old}} \leftarrow \text{vanishing gradient problem} \leftarrow \text{v.v small } \frac{\partial L}{\partial w_{ij}^k}$$

* Because of above problem deep NNs could not be trained.

* To solve above problem ReLU activation function was discovered.

* There is a problem known as Exploding Gradient it mainly occurs in RNN. We will see this in RNN chap.

Exploding Gradient $\rightarrow \frac{\partial L}{\partial w_{ij}^k} \rightarrow$ very long.

This happens coz because of this convergence will never happen.
 $(\frac{\partial L}{\partial w_{ij}^k})$ components will be greater than 1 and when we multiply values greater than 1 it becomes v. large

* Bias-Variance Tradeoff :- (MLP)

* no. of layers $\uparrow \Rightarrow$ more weights / params

\Downarrow
higher chance of overfitting
high variance.

* no. of layers $\downarrow \Rightarrow$ less wts \Rightarrow high chance of underfit

\Downarrow
High bias

- ④ To get rid of overfitting and underfitting we can use ℓ_1 or ℓ_2 regularization.
- ⑤ To get rid of over/underfitting we fine tune hyperparameters:
 - ① $L = \text{loss} + \lambda \ell_2 \text{ reg}$
 - ② $\lambda \uparrow \Rightarrow \text{var} \downarrow$
 - ③ no. of layers :-
no. of layers $\uparrow \Rightarrow \text{var} \uparrow$

⑥ Deep multi-layer perceptrons (1980s to 2010s)

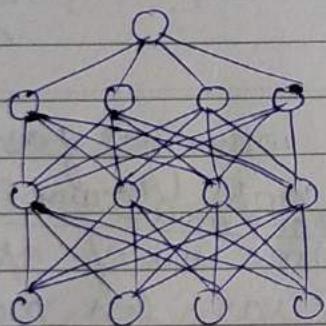
failed during (80-2010) because

→ vanishing gradient
 → too little data -> becz of which overfitting is possible common
 → too little computation power. hence it took too much time to train.

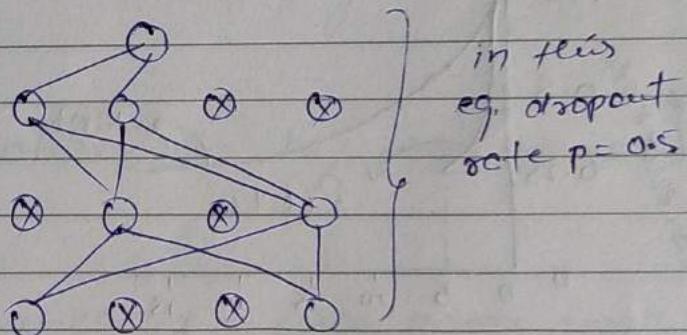
- helped in modern DL*
- ⑦ But by the end of 2010 → ① we had lots of data (also labelled data like imagenet)
 - ② computation power → powerful CPUs and GPUs
 - ③ New ideas and algorithms.
 - ⑧ classical ML - SVM / 1st theory was \neq experimental - RF / developed them \Rightarrow proofs
 - ⑨ modern DL → experiment first then find theoretical proofs.

④ Dropout layers and Regularization

- ① \Rightarrow It is easy to overfit in Deep NN because there are lots of weights.
- ② One way to avoid overfitting is by using L1 and L2 regularizers. Other way can be "dropout" layers.
- ③ In Random Forest we use randomization of features to create regularization
 - \rightarrow in Random forest we create multiple fully grown trees on subset (randomly selected features) of features and then we take majority vote
- ④ Dropout \approx random subset of features in RF



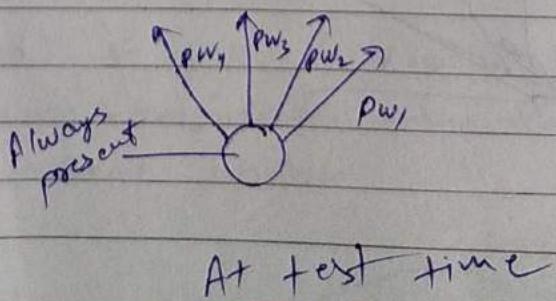
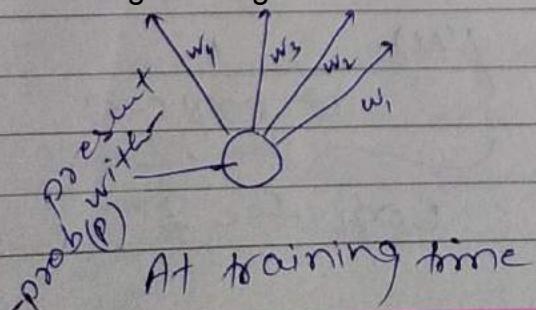
\Rightarrow Standard NN



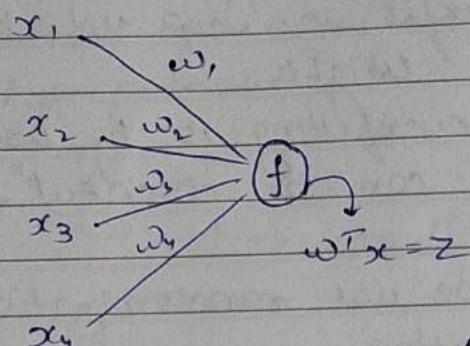
\Rightarrow After applying dropout.

in this
eg. dropout
rate $p=0.5$

- ⑤ Dropout Rate:- In a particular layer what percent of neurons/nodes will be inactive/dropped out.



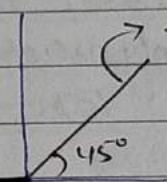
★ Rectified Linear Units (ReLU)



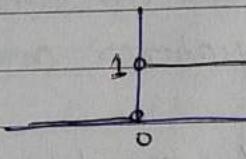
$$f(z) = z^+ = \max(0, z)$$

ReLU

not differentiable
at zero.

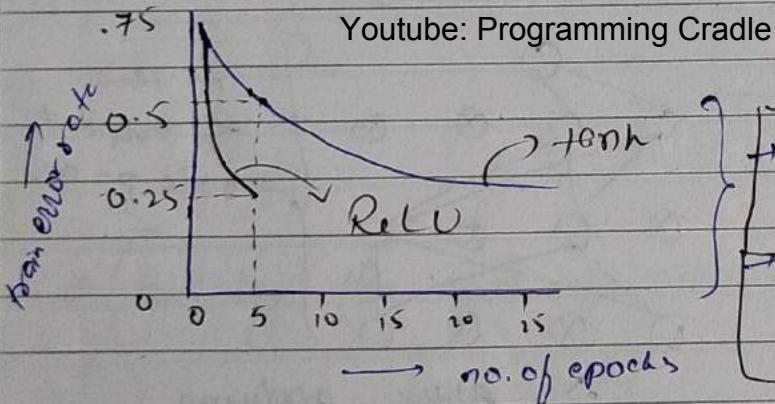


$$f(z) = \max(0, z) = \begin{cases} 0 & \text{if } z \leq 0 \\ z & \text{otherwise} \end{cases}$$



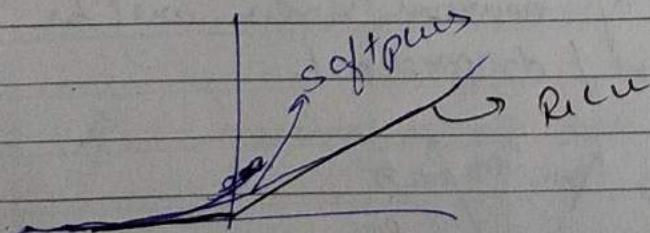
bcz of zero, dead
activation problem may
occur.

$\frac{df}{dz} \in \{0, 1\} \rightarrow$ no exploding
problem \rightarrow no vanishing
problem



- ReLU converges faster than tanh because ReLU doesn't have vanishing gradient problem
- for same no. of epochs ReLU has less error than tanh

★ Smooth approximation of rectifier $f(x) = \log(1 + \exp x)$



$$f'(x) = \frac{1}{1 + \exp(-x)}$$

Softplus

logistic fn

④ problem with ReLU :-

when z is negative $f(z) = 0$ and $\frac{df}{dz} = 0$

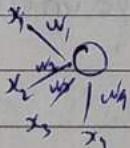
and $\frac{\partial L}{\partial w_{ij}^k}$ will become 0 { $\frac{\partial L}{\partial w_{ij}^k} = \{ \text{multiplication of several intermediate derivatives} \}$ }

we know that $(w_{ij}^k)_{\text{new}} = (w_{ij}^k)_{\text{old}} - \eta \frac{\partial L}{\partial w_{ij}^k}$ $\underbrace{1 \times 1 \times 1 \times \dots \times 0}_{\text{even if 1 deriv.}}$
 \downarrow even if 1 deriv.
 eqn ① $\quad \quad \quad$ is 0 whole mult. will become zero

if in eqn ① $\frac{\partial L}{\partial w_{ij}^k} = 0$ convergence will not happen.

⑤ what is z ?

$$z = w^T x$$



$x_i \rightarrow$ always normalized

⑥ One sol'n for dead activation can be, we can use leaky ReLU.

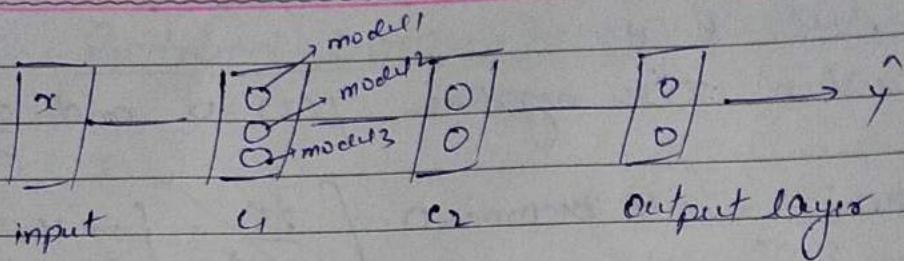
→ In leaky ReLU when z is -ve we use some small value a which is > 0 and 1 instead of using '0' but this can again lead to vanishing prob.

→ $a \rightarrow$ hyperparam. [Youtube: Programming Cradle](#)

⑦ Weight initialization :- (Deep MLP)

① init $w_{ij}^k = 0$ (or any no.) $\forall i, j, k$ } very bad idea bcz:
 Symmetry but no want asymm. $\left\{ \begin{array}{l} \textcircled{*} \text{ all neurons compute the same thing} \\ \textcircled{*} \text{ Same gradient update} \end{array} \right\}$

↳ of this
 each node $\quad \quad \quad$ our model / MLP isn't learning anything



- ① More different each models are, the better will be the performance. (some concept from RF and GBDT where more diff base models results in better performance)
- ② init $w_{ij}^k = \text{large -ve number } + ij, k$

$$w^T x = z = \underset{\substack{\uparrow \\ \text{normalized}}}{\text{large -ve values}} \Rightarrow f(z) = 0 \quad \underset{\substack{\uparrow \\ \text{ReLU}}}{f(z)}$$

Solution :-

Youtube: Programming Cradle

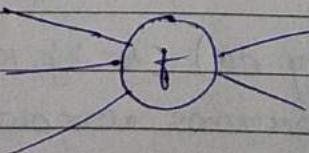
- idea 1:-
- weights should be init. with small no. (not very small)
 - not all zero.
 - good variance \Rightarrow (diff model/modulus)

$$w_{ij}^k \sim N(0, \sigma) \quad \sigma: \text{small}$$

∴ weights are normally distributed.

- ② Better init strategies:-

→ use fan-in and fan-out for init.



$$\text{fan-in} = 3$$

$$\text{fan-out} = 2$$

idea ② uniform init - (works fairly well for sigmoid) Pd t

$$\omega_{ij} \sim \text{Uniform}\left(\frac{-1}{\sqrt{\text{fan-in}}}, \frac{1}{\sqrt{\text{fan-in}}}\right) \Rightarrow$$

NOTE :- No concrete agreement with in all researchers.

Every researcher comes with their own init method
and it works fairly well.

idea ③ Xavier/Glorot init

$$\begin{array}{lll} \text{Normal} & a) \quad \omega_{ij} \sim N(0, \sigma) & \sigma_i = \frac{\sqrt{2}}{\text{fan-in} + \text{fan-out}} \\ \text{Xavier/Glorot} & & \end{array}$$

$$\begin{array}{ll} \text{Uniform} & b) \quad \omega_{ij} \sim \text{Uniform}\left(\frac{-\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}}, \frac{\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}}\right) \\ \text{Xavier/Glorot} & \end{array}$$

idea ④ He-init (2015) (works well for ReLU and leaky ReLU)

$$\text{normal} \rightarrow a) \quad \omega_{ij} \sim N(0, \sigma) \quad \sigma = \frac{\sqrt{2}}{\sqrt{\text{fan-in}}}$$

$$\text{uniform} \rightarrow b) \quad \omega_{ij} \sim \text{Uniform}\left(-\frac{\sqrt{6}}{\sqrt{\text{fan-in}}}, +\frac{\sqrt{6}}{\sqrt{\text{fan-in}}}\right)$$

⊗ Batch normalization :-

(2015)

this problem is known as internal covariance shift

↓ layers

for layer 5 data in diff batches may have diff dist. being normalized is not done at each layer.

$x \rightarrow l_1 \rightarrow l_2 \rightarrow l_3 \rightarrow l_4 \rightarrow l_5 \rightarrow \text{Output} \rightarrow y$

mini batch for layer 1

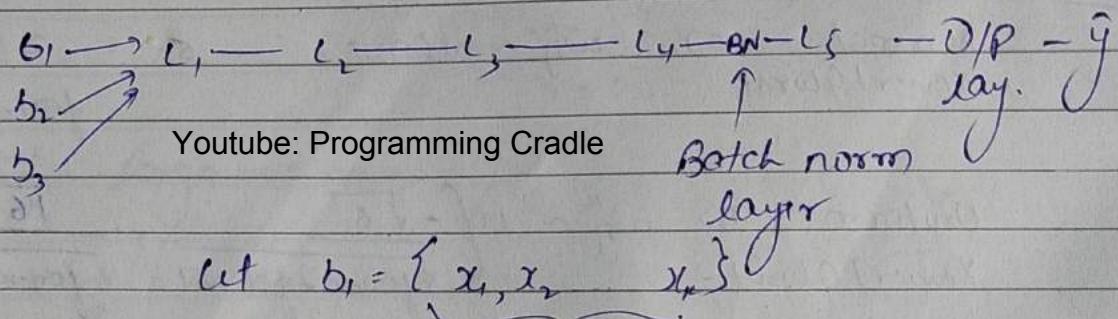
data is properly normalized since b_1 and b_2 will be same

also data is normalized since

also batch size is small

Solution for internal cov. shift :-

- We add addition layer known as batch normalization layer. In this layer only current batch is normalized based on batch's mean and std-dev.
- Usually batch layer is used deep inside the network.

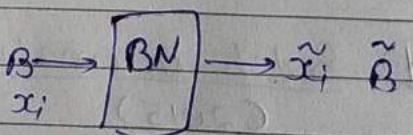


NOTE:- BN can be added to any layers.

norm will be done for only these pts. not on whole data.

Input:- Values of x over a mini-batch : $B = \{x_1, \dots, x_n\}$
parameters to be learned γ, β

Output:- $\{\tilde{x}_i : BN_{\gamma, \beta}(x_i)\}$



All our backpropagation will be done w.r.t γ, β

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad \| \text{mini-batch mean}$$

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad \| \text{mini-batch variance}$$

$$x_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad \| \text{normalize}$$

small value to avoid divide by zero

$$\tilde{x}_i \leftarrow \gamma x_i + \beta \equiv BN_{\gamma, \beta}(x_i) \quad \| \text{scale and shift}$$

→ Advantages of BN:-

•→ faster convergence because it is guaranteed that all the batches will have same distribution and because of it we can afford to have larger learning rate (η)

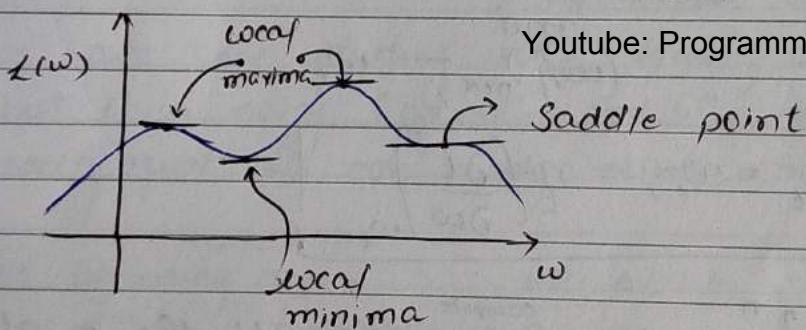
$$(\omega_{ij}^n)_{\text{new}} = (\omega_{ij}^n)_{\text{old}} - \eta \left(\frac{\partial L}{\partial \omega_{ij}^n} \right)$$

•→ BN also works as weak batch regularization so we should use BN in addition to dropouts.

•→ Avoids internal cov. shift. and hence deep NN can be trained.



Optimizers :- Hill descent analogy in 2D

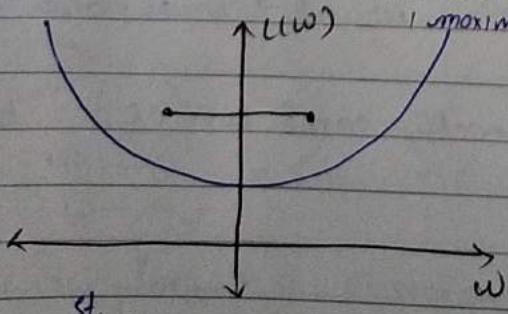


Youtube: Programming Cradle

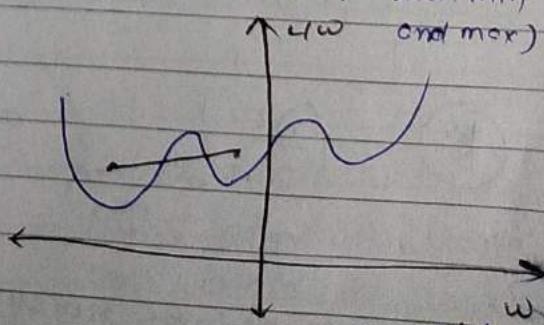
NOTE :- SGD or mini-batch SGD could get stuck at saddle point.

convex f'' (will have 1 local

minima or
1 maxima)



non-convex f''
(can have multiple local min and max)



if line connecting 2 points in the same region lies completely in one region then it is convex f'' else non-convex.

NOTE :- → for Logistic Regression, linear Regression and SVM loss f'' can be shown as convex f'' and hence SGD or batch SGD works fine.

→ Whereas in MLP loss f'' is non convex and can have multiple local minimum and local maxima. Also it can have saddle point and hence SGD or both SGD won't always work.

* SGD and momentum

$$\text{in each iter. } (\omega_{ij}^k)_{\text{new}} = (\omega_{ij}^k)_{\text{old}} - \eta \left[\frac{\partial L}{\partial \omega_{ij}^k} \right] \quad \left. \begin{array}{l} \downarrow \\ (\omega_{ij}^k)_{\text{old}} \end{array} \right\} \text{update } f^n$$

for simplicity $(\omega_{ij}^k)_{\text{new}} \rightarrow \omega_t$

and
 $(\omega_{ij}^k)_{\text{old}} \rightarrow \omega_{t-1}$

$$\boxed{\omega_t = \omega_{t-1} - \eta \left[\frac{\partial L}{\partial \omega} \right]_{\omega_{t-1}}}$$

$\mathcal{D} = \{x_i, y_i\}_{i=1}^n$

$\frac{\partial L}{\partial \omega} \xrightarrow{\text{compute}} \text{using all the } n \text{ pts in } \mathcal{D}, \text{ it}$
 is called as GD

$\xrightarrow{\text{compute using one pt chosen at random, known as SGD}}$
 mostly used
 $\xrightarrow{\text{using random subset of } k \text{ pts in } \mathcal{D} \rightarrow \text{mini-batch SGD}}$



$$\left(\frac{\partial L}{\partial \omega} \right)_{\text{mini-SGD}}$$

↑
or
SGD

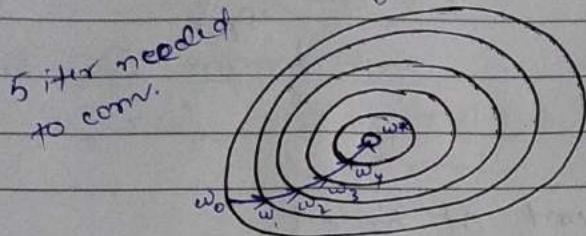
tends to be noisy.

are not exactly same as $\left(\frac{\partial L}{\partial \omega} \right)_{\text{GD}}$ but close

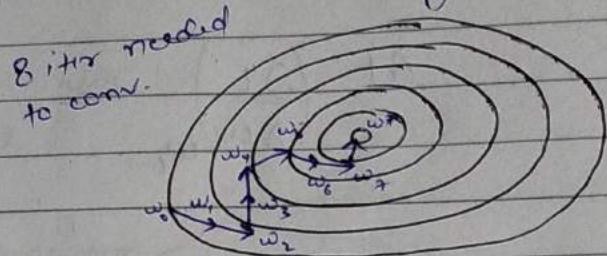
↑
estimates

ideally we want

④ GD convergence.



* SGD convergence



NOTE:- To get the "almost" same w^* as in GD, SGD requires more no. of iterations.

$$w_{GD}^* \xrightarrow{\text{requires}} 5 \text{ iter.}$$

similar \rightarrow SS

$$w_{SGD}^* \xrightarrow{\text{requires}} 20 \text{ iter may be, or even more}$$

⑤ Batch SGD with momentum :-

↳ exponential weighting

As we have seen finding w^* using SGD work well but it is noisy. (as w_{GD}^* and w_{SGD}^* are approx. equal not exactly equal) We can apply denoise methods.

⑥ Denoise by using augs:-

t_1, t_2, t_3, \dots

date pts. $\rightarrow a_1, a_2, a_3, \dots$

$$t=1 \quad v_1 = a_1$$

$$t=2 \quad v_2 = \sqrt{v_1} + a_2$$

$$t=3 \quad v_3 = \sqrt{v_2} + a_3 = \sqrt{(\sqrt{v_1} + a_2)} + a_3$$

$$\text{let } r = 0.5$$

$$\begin{aligned} &= \sqrt{v_1} + \sqrt{a_2} + a_3 \\ &= 0.25a_1 + 0.5a_2 + 1a_3 \end{aligned}$$

$$0 < r \leq 1$$

$$\Rightarrow \left\{ \begin{array}{l} v_1 = a_1 \\ v_t = \sqrt{v_{t-1}} + a_t \end{array} \right\}$$

$v_t \rightarrow \approx$ denoised estimate.

$$\begin{matrix} \uparrow & \uparrow & \uparrow \\ t_1 & t_2 & t_3 \end{matrix}$$

\therefore Giving more wt. to latest pt. and lesser and lesser to previous pts.

$$\text{mini-batch SGD} \rightarrow \omega_t = \omega_{t-1} - \eta \left(\frac{\partial L}{\partial \omega} \right)_{\omega_{t-1}}$$

$$\text{at } \left(\frac{\partial L}{\partial \omega} \right)_{\omega_{t-1}} = g_t$$

↑ gradient at time t

$$\text{std. update} \Rightarrow \omega_t = \omega_{t-1} - \eta g_t \quad \text{①} \rightarrow \text{mini-batch SGD}$$

eqn

eqn ① modified using expon. weighted. $[0 < \gamma < 1]$

$$\begin{cases} v_t = \gamma v_{t-1} + \eta g_t \\ \omega_t = \omega_{t-1} - v_t \end{cases}$$

good value of $\gamma = 0.9$
(thumb rule)

$$\text{case 1 } \gamma = 0 ; v_t = \eta g_t \Rightarrow \omega_t = \omega_{t-1} - \eta g_t$$

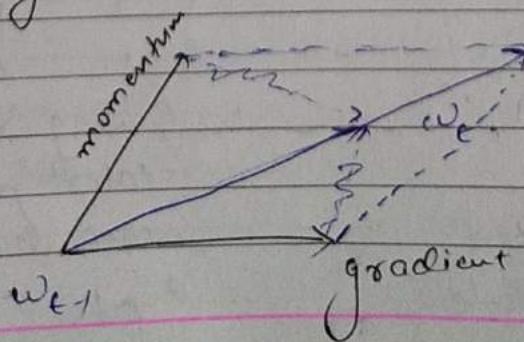
$$\text{case 2 } \gamma = 0.9 ; v_t = 0.9 v_{t-1} + \eta g_t ; \omega_t = \omega_{t-1} - (0.9 v_{t-1} + \eta g_t)$$

⊕ SGD + momentum = exp. weighting to denoise your SGD gradients.

$$\omega_t = \omega_{t-1} - [\gamma v_{t-1} + \eta g_t]$$

↑ momentum ↑ gradient

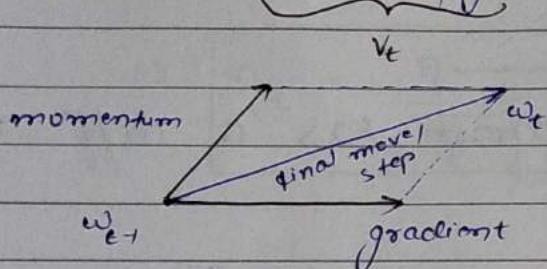
Speeds up
convergence



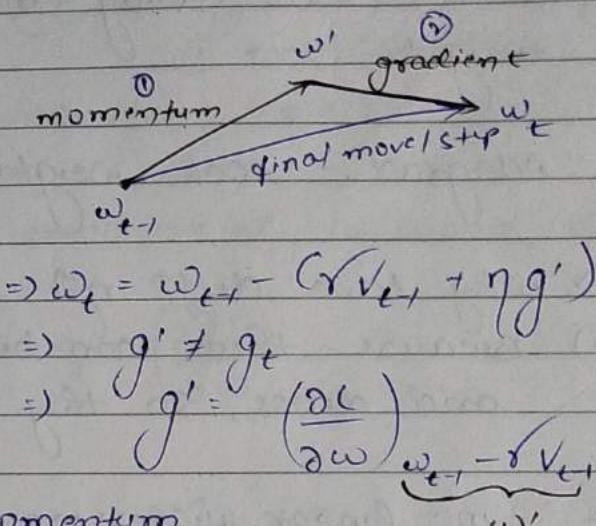
(*) Nesterov Accelerated Gradient (NAG)

④ SGD+ momentum

$$\omega_t = \omega_{t-1} - (\gamma v_{t-1} + \eta g_t)$$



④ NAG

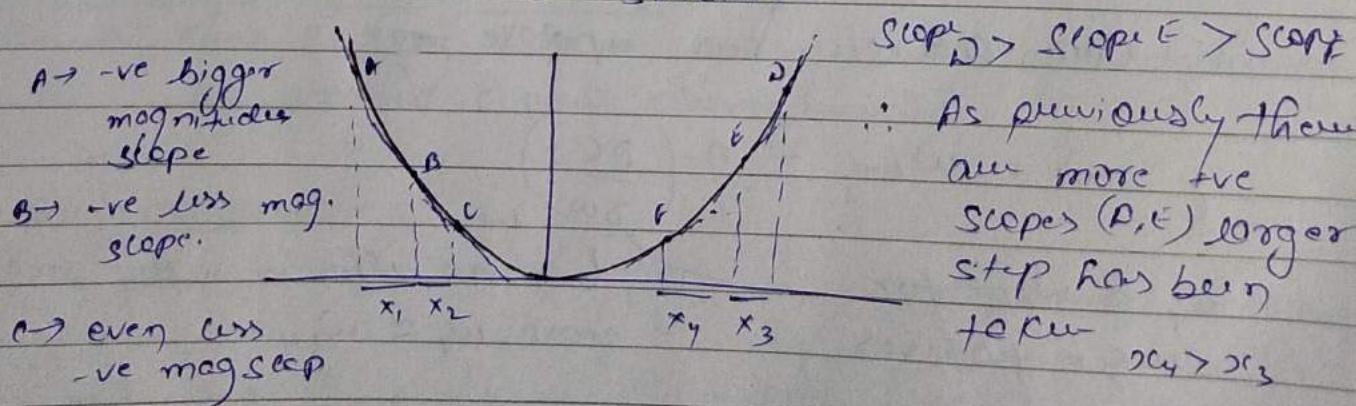


Q) Why gradient and momentum has diff. directions?

Youtube: Programming Cradle

Ans) → Momentum gives direction based on previous gradient / slope, so if there are many +ve slopes steps in -ve direction will be bigger and if no. of -ve slopes are more, steps in +ve dirn will be bigger.

→ Gradient gives direction based on magnitude and polarity of the slope. If slope is large -ve value, step in +ve direction will be bigger and if slope is small -ve value, step in +ve direction will be smaller.



∴ As we can see x4 > x2

Adagrad (adaptive gradients)

- ④ For SGD and SGD + momentum learning rate (η) is constant (very small no. 0.01) for each weight.
- ⑤ Adagrad :- each weight / param has a diff η

Q) Why have diff η ?

A) Because there can be features which are sparse and dense, so they need diff learning rates.

e.g. in linear regression (loss = sqrd loss)

$$\left(\frac{\partial L}{\partial w} \right)_{w_{t-1}} = \sum_{i=1}^n (-2x_i)(y_i - w_{t-1} \cdot x_i)$$

sparse dense

most of times
will be '0' and
hence $\frac{\partial L}{\partial w}$ will be
very small

most of times values
will be non zero
hence $\frac{\partial L}{\partial w}$ will be
large.

Now lets see our update eqn

$$w_t = w_{t-1} - \eta \left(\frac{\partial L}{\partial w} \right)_{t-1}$$

$w_t \neq w_{t-1}$ for large features \rightarrow as this is very small for sparse, $w_t \approx w_{t-1}$

SGD: $\omega_t = \omega_{t-1} - \eta \widehat{g_t}$ same for all wts.

Adagrad: $\omega_t = \omega_{t-1} - \eta' g_t$ diff' η' for each feature
@ each iteration.

$$\Rightarrow \eta'_t = \frac{\eta}{\sqrt{d_t + \epsilon}} \rightarrow \text{constant (0.01)}$$

small +ve no. to avoid divide by

$$d_t = \sum_{i=1}^{t-1} g_i^2$$

↑ +ve

$$g_i = \left(\frac{\partial L}{\partial \omega} \right)_{\omega_{i-1}}$$

and $d_t \geq d_{t-1} \Rightarrow t \uparrow \Rightarrow \alpha \downarrow \Rightarrow \eta'_t \downarrow$
iteration

* Advantages:-

- no need of manually tune η
- chooses right value of η based on previous grad for sparse and dense features

* Disadvantages:-

- d_t can become very large as $t \uparrow$ and because of this convergence will be slower.

* AdaDelta and RMS Prop.

- These 2 algos help in getting rid of disadvantage of adagrad (slow convergence)

=) Ada delta:-

$$\omega_t = \omega_{t-1} - \eta'_t g_t \rightarrow \text{exponential decay avg.}$$

$$\eta'_t = \frac{\eta}{\sqrt{\text{cda}_t + \epsilon}} ; \text{cda}_t = \gamma * \text{cda}_{t-1} + (1-\gamma) g_t^2$$

Adadelta: exp. wt avg. of g_i^2 instead of sum of g_i^2 (add grad) so as to avoid large denominator in η_i hence slow conv. can be avoided.

★ Adam: (Adaptive Moment Estimation)

$$\text{eda } g_t : m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t \quad 0 < \beta_1 < 1$$

$$\text{eda } g^2 : v_t = \beta_2 v_{t-1} + (1-\beta_2) g^2 \quad 0 < \beta_2 < 1$$

$$\hat{m}_t = \frac{m_t}{1-(\beta_1)^t} ; \hat{v}_t = \frac{v_t}{1-(\beta_2)^t} \quad \left. \right\} \text{ bias correction}$$

$$w_t = w_{t-1} - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \quad \begin{array}{l} \text{mean} \rightarrow 1^{\text{st}} \text{ order moment} \\ \text{var} \rightarrow 2^{\text{nd}} \text{ order moment} \end{array}$$

typically $\alpha = 0.001$ Youtube: Programming Cradle

Q> which optimizer / algorithm to choose?

Mini Batch SGD \rightarrow Small / shallow NN



Momentum based \rightarrow works well in most cases
NAG
↓

Adagrad \rightarrow sparse f

Adadelta & RMS prop \rightarrow sparse



Adam \rightarrow Best of all mentioned above

* Gradient Monitoring and clipping...

- * To get rid of exploding gradient we can use clipping

It is a good habit to monitor gradients and updates for each epoch and weights and layers. It will help in knowing if vanishing grad. prob. is occurring or not, exploding prob or any other prob.

$$\text{Grad}_{\text{new}} = \frac{\text{Grad} * \tau}{\|\text{Grad}\|_2}$$

threshold.

$$\text{At } \text{Grad} = \begin{bmatrix} g_1 & g_2 & g_3 & g_4 & g_5 \end{bmatrix} / 1000000$$

$$\frac{\text{Grad}}{\|\text{Grad}\|_2} = \underbrace{\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix}}_{< 1}$$

$$\|\text{Grad}\|_2 = \sqrt{g_1^2 + g_2^2 + g_3^2 + g_4^2 + g_5^2}$$

Youtube: Programming Cradle

$\Rightarrow \text{Grad}_{\text{new}}$ will always be less than τ (threshold)

- * Softmax and cross-entropy for multi-class classifier

Softmax classifier = logistic func + multi-class

$\Rightarrow LR$

$$p(y_i=1|x_i) = \hat{y}_i = \sigma(z) = \frac{1}{1+e^{-z}}$$

Softmax classifier

$$D = \{(x_i, y_i)\} \text{ multi-class}$$

$$y_i \in \{1, 2, 3, \dots, K\}$$

$$x_i \rightarrow M \rightarrow \begin{cases} p(y_i=1|x_i) \\ p(y_i=2|x_i) \\ p(y_i=3|x_i) \\ \vdots \\ p(y_i=K|x_i) \end{cases}$$

Sum 1

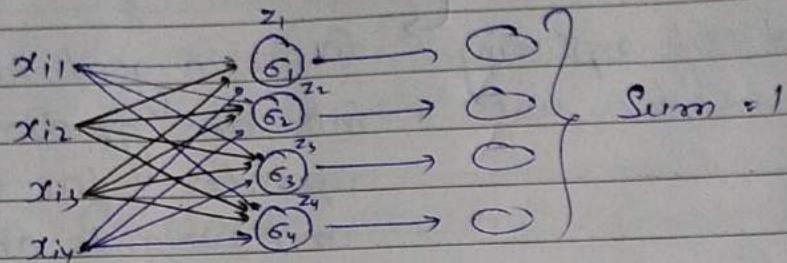
- * Cross-entropy - it can be used to define a loss

f^n in ML and optimizat^n, it is

defined on prob. dists, not single values.

* Soft max

wt $K=4$



$$z_1 = \sum_{j=1}^d x_{ij} w_{j1}; z_2 = \sum_{j=1}^d x_{ij} w_{j2} \dots$$

$$\sigma_1(z_1) = \frac{e^{z_1}}{\sum_{i=1}^K e^{z_i}}; \sigma_2(z_2) = \frac{e^{z_2}}{\sum_{i=1}^K e^{z_i}} \dots$$

$$\Rightarrow \sigma_1(z_1) + \sigma_2(z_2) + \dots + \sigma_K(z_K) = 1$$

$$\Rightarrow \frac{e^{z_1}}{\sum_{i=1}^K e^{z_i}} + \frac{e^{z_2}}{\sum_{i=1}^K e^{z_i}} + \dots + \frac{e^{z_K}}{\sum_{i=1}^K e^{z_i}}$$

$$\Leftrightarrow \frac{e^{z_1} + e^{z_2} + \dots + e^{z_K}}{\sum_{i=1}^K e^{z_i}}$$

Youtube: Programming Cradle

$$\Rightarrow \frac{\sum_{i=1}^K e^{z_i}}{\sum_{i=1}^K e^{z_i}} = 1$$

* Softmax \rightarrow generalization of LR to multi-class setting.

LR :- $\sigma(z) = \frac{1}{1+e^{-z}} = \frac{e^z}{e^z + 1}$

minimizes log-loss

Softmax :- $\sigma_1(z_1) = \frac{e^{z_1}}{\sum_{i=1}^K e^{z_i}}; \sigma_1, \sigma_2, \dots, \sigma_K$

* Multi class log-loss

$$-\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K y_{ij} \log p_{ij}$$

* 2 class log-loss

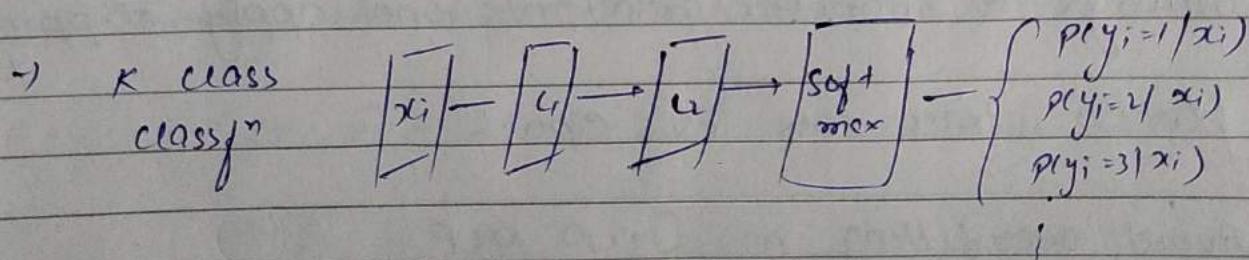
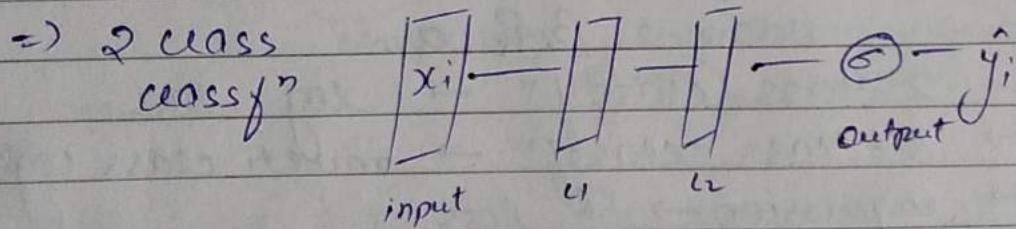
$$-\frac{1}{N} \sum_{i=1}^N [y_i \log p_i + (1-y_i) \log (1-p_i)]$$

$$x_i, y_i = 3$$

$$\begin{matrix} 1 & 2 & 3 & \dots & K \\ | & | & | & \dots & | \\ 0 & 0 & 1 & \dots & 0 \end{matrix}$$

$$y_{i1} \quad y_{i2} \quad y_{i3} \quad y_{i4}$$

$$x_i \rightarrow \boxed{M} \rightarrow \begin{cases} 0.2 \\ 0.1 \\ 0.7 \\ 0 \end{cases} \begin{array}{l} \rightarrow p(y_i=1/x_i) \\ \rightarrow p(y_i=2/x_i) \\ \rightarrow p(y_i=3/x_i) \\ \vdots \end{array}$$



* In case of regression we can use linear unit instead of sigmoid

$$f(z) = z$$

linear unit

(*) "How to" train an MLP?

- ① Preprocess :- Data Normalization
- ② Weight initialization :- ✓ Xavier / Glorot \rightarrow Sigmoid / tanh
:- $H_i \rightarrow R_{i, j}$
:- Gaussian (small value of σ)
- ③ Choose activation f^n :- ReLU
- ④ Add Batch normalization especially for deep MLP (later layers)
dropout (regularization)
- ⑤ Optimizer \rightarrow ADAM
- ⑥ hyperparameters :- Architecture \rightarrow no. of layers
 \rightarrow no. of neurons
:- Dropout rate \rightarrow
:- Adam $\rightarrow \beta_1, \beta_2, \alpha$
- ⑦ loss f^n \rightarrow 2 class classif \rightarrow log loss
 \rightarrow K-class classif \rightarrow multi class log loss
 \rightarrow regression \rightarrow SQ. loss
- ⑧ Always monitor gradients and apply clipping.
- ⑨ plots : train loss vs epoch
- ⑩ Avoid overfitting in Deep MLP

NOTE :- It's good to have fewer neurons in other layers as compared to earlier layers.

* Type of NN:- Auto encoder

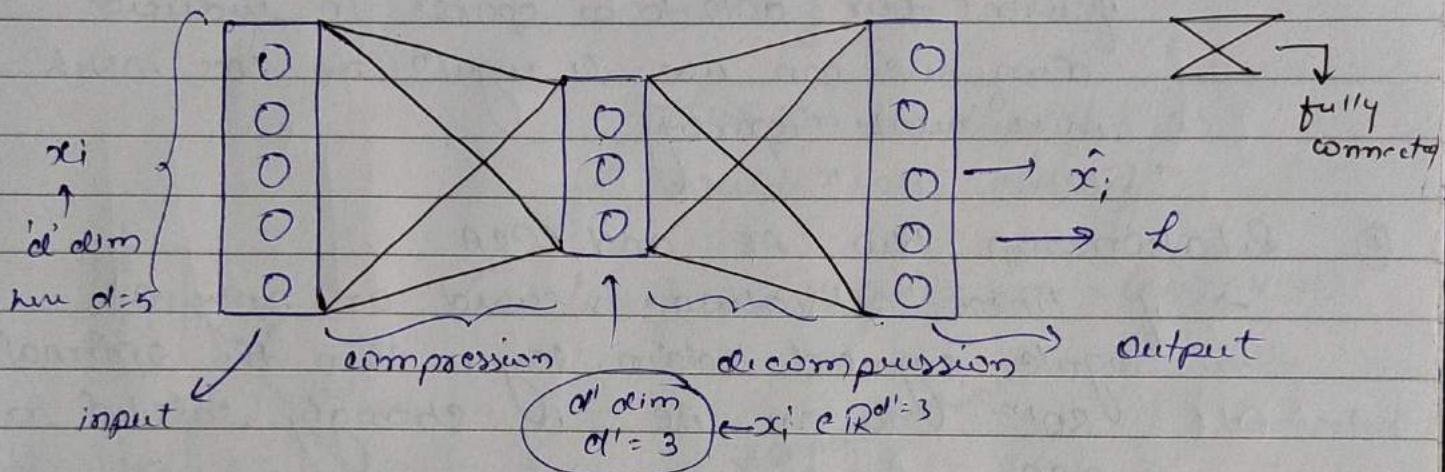
→ Auto encoder is a NN which performs dim. reduction

→ Objective :- $D = \{x_i\}_{i=1}^n, x_i \in \mathbb{R}^d$



$$d' < d$$

$$D' = \{x'_i\}_{i=1}^n, x'_i \in \mathbb{R}^{d'}$$



* If $\hat{x}_i \approx x_i$ we can say that d' (3) is sufficient enough to represent data in d (5) dim

$$L(x_i, \hat{x}_i) = \|x_i - \hat{x}_i\|^2$$

↑ loss f

Youtube: Programming Cradle

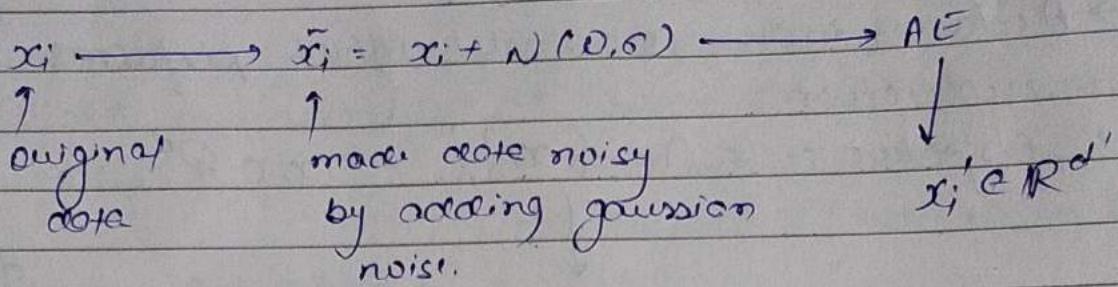
NOTE :- We can have multiple hidden layers instead of just 1.

* Denoising Auto encoder

$D = \{x_1, x_2, x_3, \dots, x_n\} \leftarrow$ actual data
we get $\tilde{D} = \{\tilde{x}_1, \tilde{x}_2, \tilde{x}_3, \dots, \tilde{x}_n\}$

let \tilde{x} is noisy

Steps in denoising AE :-



④ Sparse Encoders:-

if want our data to be sparse in reduced dim we can use 'l1 regul' on loss which will create sparsity.

⑤ Relationship b/w AE and PCA

→ if linear activation is used or only a single sigmoid hidden layer, then the optimal soln to an AE is strongly related to PCA

⑥ Word2Vec CBOW (continuous bag of words)

focus word :- word for which vec is being generated.

context word :- All the other words in the sentence

e.g.

The cat sat on the wall.

~~~~~                    ~~~~~

context word            focus word            context word.

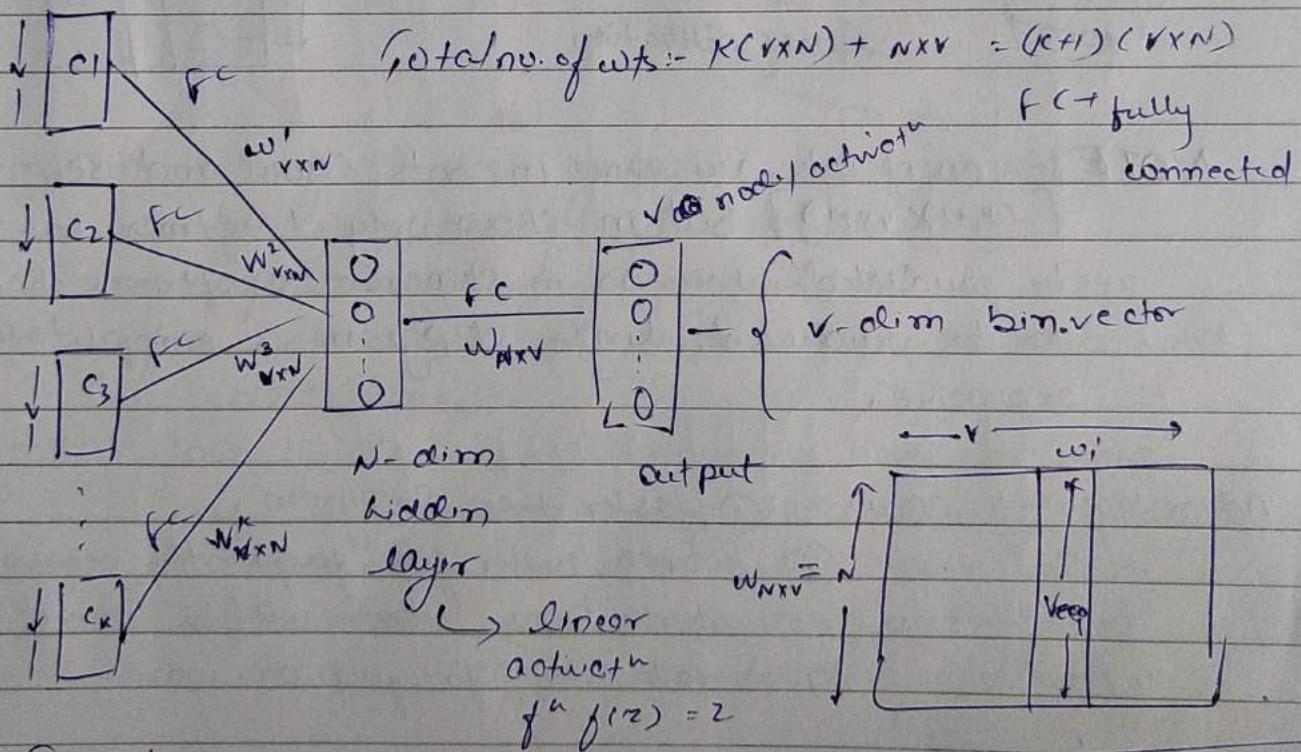
↓  
sat 'sat' is a focus word.

①  $\rightarrow V \rightarrow$  dictionary / vocabulary of words  
 $\Rightarrow v \rightarrow \text{dim}(V)$

② One hot encoding each word  
 $w_i \in \mathbb{R}^v$ : binary vector of  $V$ -dim.

$\Rightarrow$  Core Idea of CBOW :- Given Context words, can we predict focus word

This can be thought as multi class classifier



dataset will look something like this:

focus word

$w_1$

$w_2$

context word

$w_2, w_3, w_4, w_5$

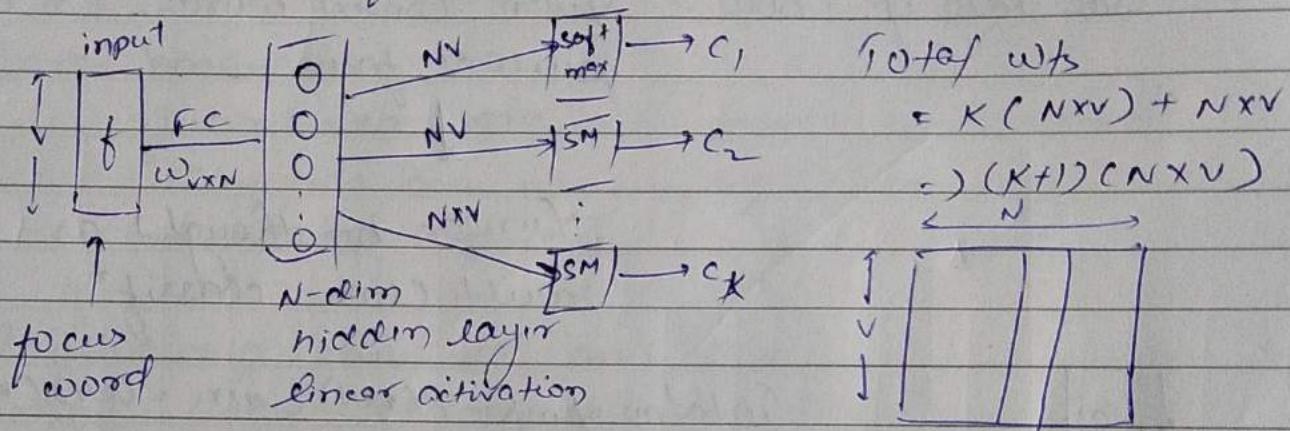
$w_1, w_3, w_5, w_4$

sent:  $w_1, w_2, w_4, w_3, w_5$

## \* Word2Vec: skipgram

CBOW :- predict focus word given context words

Skipgram:- predict context words given focus word.



**NOTE :-** no. of wts are same in both CBOW and skipgram  $[ (K+1)(VxN) ]$  but in CBOW only 1 softmax has to be predicted whereas in skipgram  $K$  softmax's has to be predicted. hence skipgram is computationally expensive.

Advantages of CBOW :-

- ① faster than Skipgram
- ② works better for frequently occurring words because there will be more data points for freq. occ. words.

Advantages of skipgram -

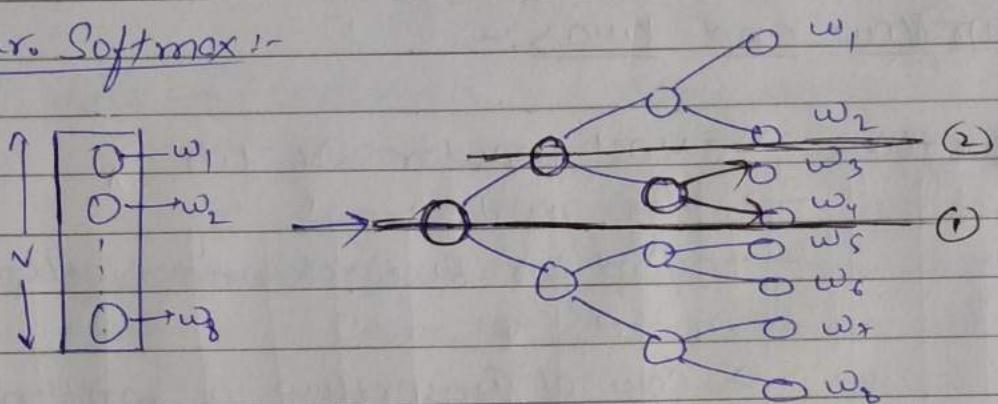
- ① can work well with smaller amount of data
- ② works well for infrequent words.

Major problem with CBOW and skipgram is they have so many wts.  $(k+1)(v \times N)$

### Algorithmic Optimizations:-

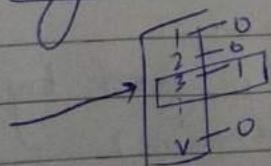
- ↳ hierarchical Softmax
- ↳ Negative Sampling.

### Hier. Softmax :-



- instead of linearly going through all the  $V$ -class/wts we can use binary tree data structure
- at first node if  $p(y_i = w_3 | x_i)$  is more ~~less~~ then we can skip all the nodes below block line ①.
- Now if  $p(y_i = w_3 | x_i)$  is more we can skip all the nodes above block line ②
- at last node we can select our class.
- Here we are computing on  $\log(V)$  instead of  $V$  wts.

Negative Sampling :- update only a sample of words per iteration



- Sampling :- ① always keep the target word. ( $word_3$ )
- ② amongst all the non-target words don't ~~update~~ all, update sample ~~of~~ <sup>from</sup> target.

Sample is selected based on below prob.

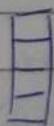
$$p(w_i) = 1 - \frac{1}{f_{\text{eng}}(w_i)} \rightarrow \text{threshold typically } 10^{-5}$$

Probability of a word getting selected in the sample will be more if it occurs more frequently.

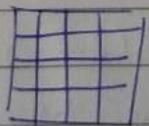
## ⊗ Tensorflow and Keras:-

- Tensorflow → most popular DL lib
- By google
- Good for research and development as well.
- Core of Tensorflow is written in c/c++
- Has interface with Py., Java, JS
- Tensorflow ~~etc~~ → can run on android
- learning curve is steep.

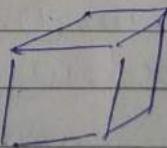
## ⊗ Tensor :- n-dim array.



1D array/  
tensor



2D array/tensor 3D-tensor

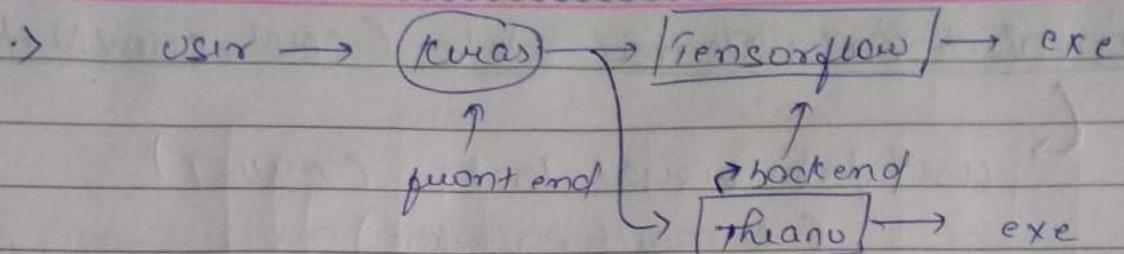


- ⊗ Flow → Tensors are flowing from i/p layer to o/p layer.

## ⊗ Keras :-

- High-level NN-lib (mostly used by developers)
- easy to learn
- fewer lines of code.

- learning.tensorflow.com
- cloud.google.com learn deeper without phd.



✳ Other alternatives of TF :- Theano, coffee, PyTorch, MXNET

✳ Keras is deeply integrated with TF

↳ Some low level control.

✳ TF has → constants } same as in  
                   → variable      } python  
                   → placeholders. ⇒ stores mini-batches,  
 $\Rightarrow x = \text{tf.placeholder}(\text{tf.float32},$       once the processing is  
 $[None, 784])$       done on one batch,  

 next batch will be stored  
 in placeholder

↳ don't know, can be copy + clip

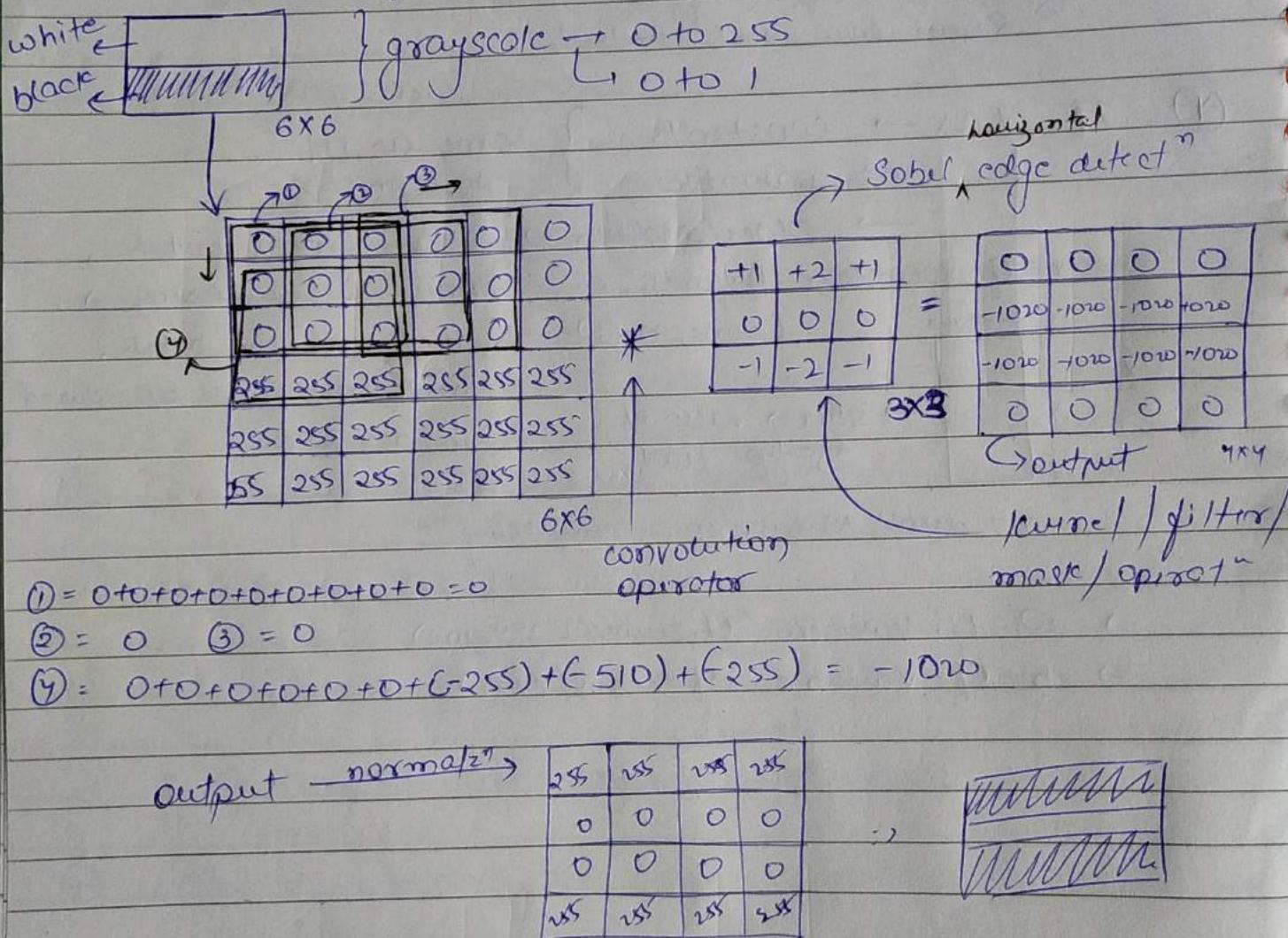
$\Rightarrow w = \text{tf.Variable}(\text{tf.zeros}([784, 10]))$

$\Rightarrow b = \text{tf.Variable}(\text{tf.zeros}([10]))$

## ④ Convolutional Neural Nets (CNN, convNets)

↳ used for visual tasks (MNIST)  
(computer vision)

### ⑤ Convolution Edge Detection on Images:



### ⑥ Sobel vertical edge det

|    |   |    |
|----|---|----|
| +1 | 0 | -1 |
| +2 | 0 | -2 |
| +1 | 0 | -1 |

$$\text{input } (n \times n) \xrightarrow[\text{turn} \cdot 1]{(K \times K)} (n - k + 1) \times (n - k + 1) \quad \text{output} \quad \left| \begin{array}{l} \text{eg: } 6 \times 6 * 3 \times 3 = 4 \times 4 \\ \uparrow \text{memory size mg.} \end{array} \right.$$

## ✳️ Convolution Padding:

→ If we don't want to reduce the size of the output image we can apply padding on the input image.

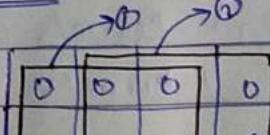
① Zero padding (most often used)    ② Some value padding

|   |     |     |     |     |     |     |   |
|---|-----|-----|-----|-----|-----|-----|---|
| 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0 |
| 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0 |
| 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0 |
| 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0 |
| 0 | 255 | 255 | 255 | 255 | 255 | 255 | 0 |
| 0 | 255 | 255 | 255 | 255 | 255 | 255 | 0 |
| 0 | 255 | 255 | 255 | 255 | 255 | 255 | 0 |
| 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0 |

↳ pad with some value as nearest value.  
(this tech. is not used much)

$$\Rightarrow n \times n \xrightarrow[\substack{\text{with} \\ \text{padding} \\ p}]{K \times K} (n - k + 2p + 1) \times (n - k + 2p + 1)$$

## ✳️ Strides :-

|       |                                                                                     |   |
|-------|-------------------------------------------------------------------------------------|---|
| - 255 |  | - |
|       | 0 0 0   0 0 0                                                                       |   |
|       | 0 0 0   0 0 0                                                                       |   |
|       | 0 0 0   0 0 0                                                                       |   |
|       | - - -   - - -                                                                       |   |
|       | - - -   - - -                                                                       |   |
|       | - - -   - - -                                                                       |   |

from ① to ② stride of 1 pixel  
it can be stride of 2 or more pixels as well.

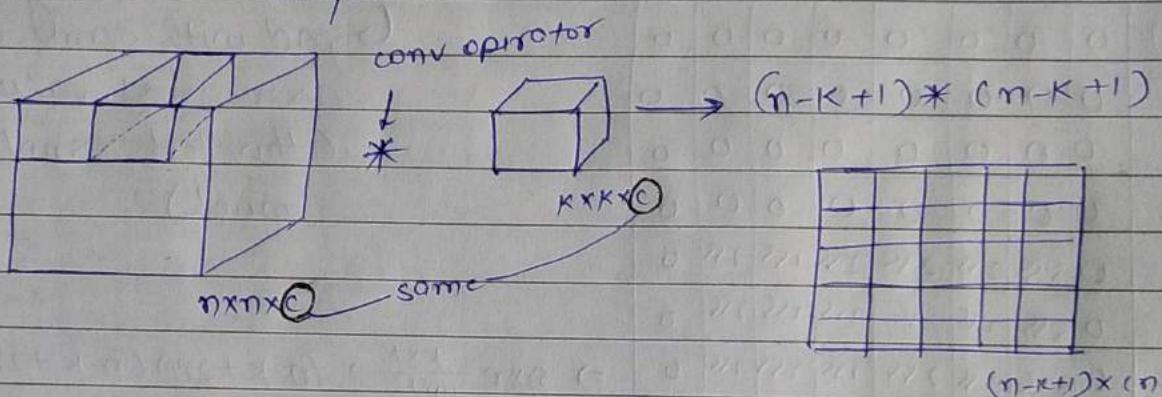
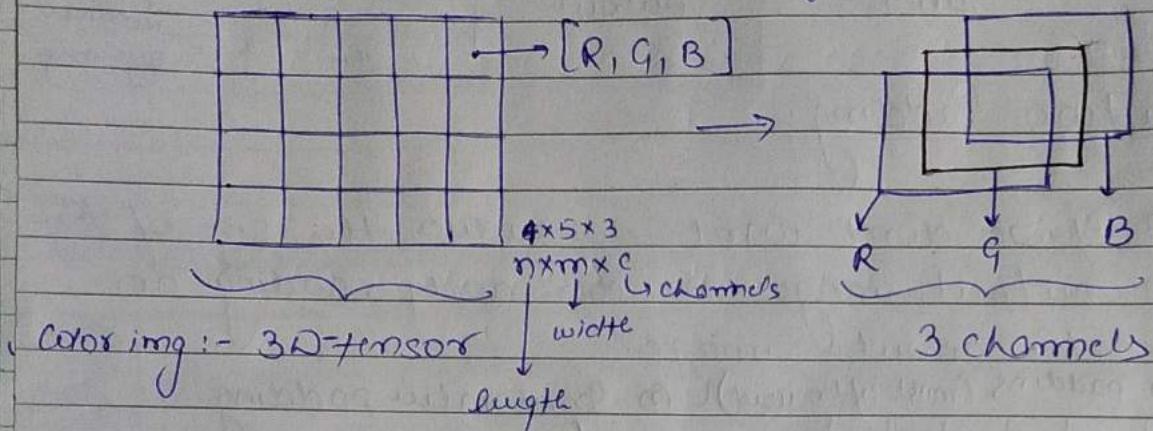
$$n \times n \xrightarrow[\substack{\text{Stride: } s \\ K \times K}]{\text{Stride: } s} \left( \left\lfloor \frac{n - K}{s} \right\rfloor + 1 \right) \times \left( \left\lfloor \frac{n - K}{s} \right\rfloor + 1 \right)$$

$n \times n$   
 $n=6$

eg:-

$$6 \times 6 \xrightarrow[\substack{K=3 \\ s=2}]{\text{Stride: } s} 2 \times 2$$

## \* Convolution on color-images :-



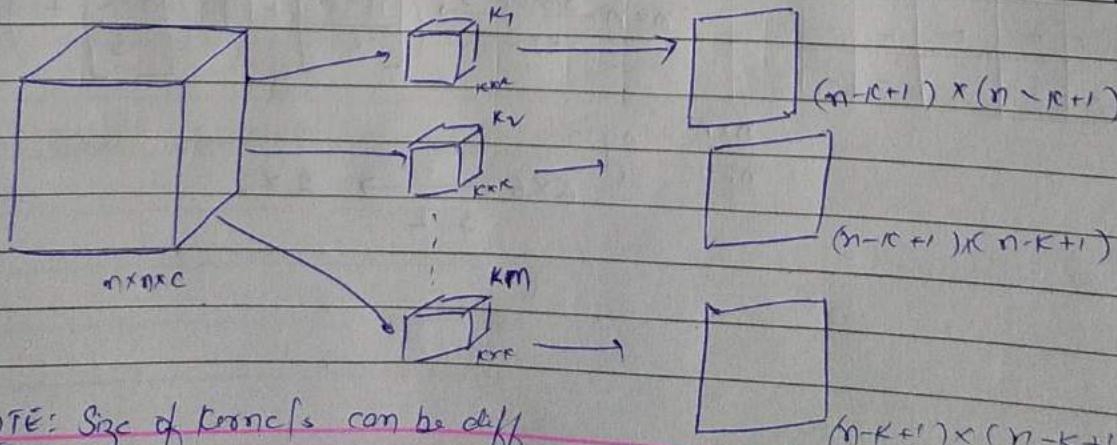
## \* Convolution layer in NN

- ① Conv. layers are inspired by biological brain
- ② for multiple edge detection  $\rightarrow$  multiple kernels are needed  
1 edge detected by 1 kernel.

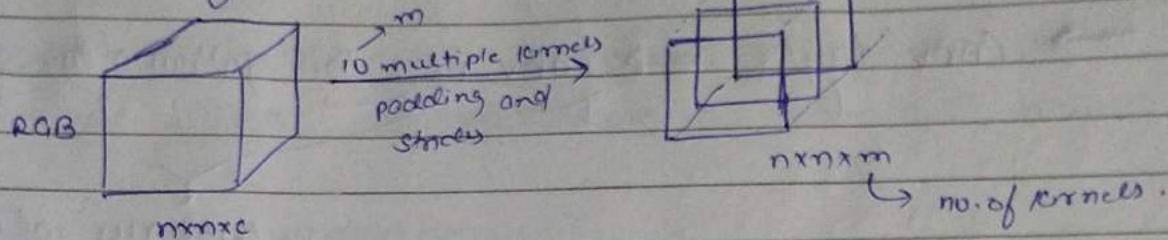
in CNN kernels are learned using backpropagation

in MLP :- w's are learned  $\rightarrow$  w.x

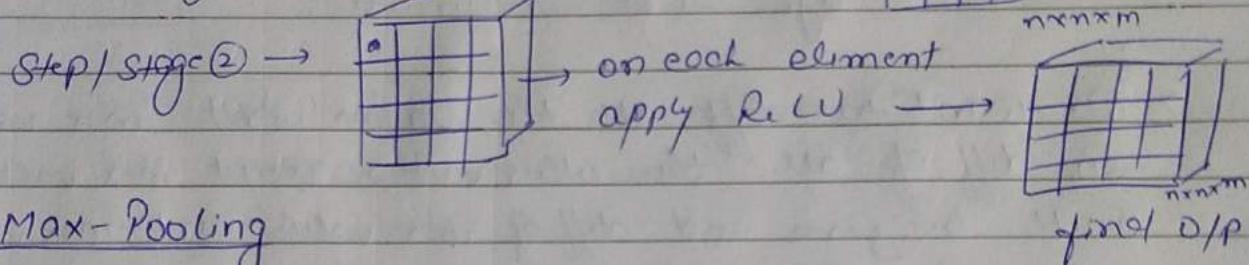
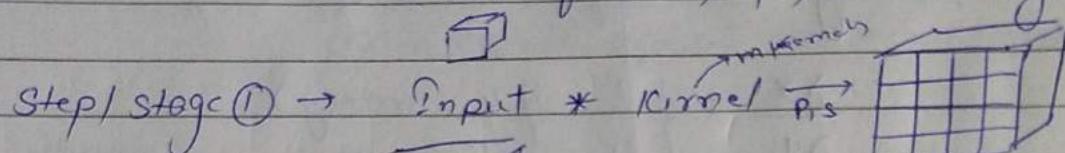
in CNN :- kernel matrices are learned  $\rightarrow 1 \times \omega_{R \times K}$



NOTE: Size of kernels can be diff

Single Conv layer :-

$m$  (no. of kernels);  $p$ ;  $s$  are hyperparam.

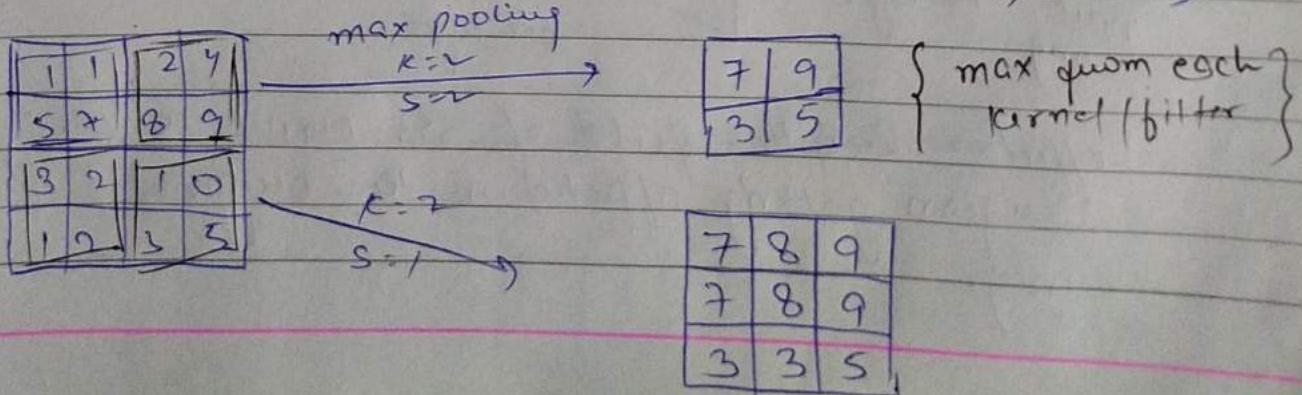
Max-Pooling

$\Rightarrow$  location invariance :- Where the object is in the image, it should be detected

$\Rightarrow$  Scale invariance :- whatever is the size of the obj. it should be detected.

$\Rightarrow$  rotation invariance :- even if the obj is rotated, it should be detected.

$\Rightarrow$  Max-pooling helps in making the model invariant to above 3 invariance. (locn<sup>h</sup>, rot<sup>h</sup>, scale)



## \* CNN training Optimization :-

→ Conv-layer  $\Rightarrow$  conv operation followed by ReLU

both are differentiable

→ Max-pooling  $\Rightarrow$  differentiation of max<sup>n</sup>

$$\frac{\partial x}{\partial y_i} = \begin{cases} 1 & \text{if } y_i = \max(y) \\ 0 & \text{otherwise.} \end{cases}$$

$\Rightarrow$  We want to diff all the values w.r.t max value  
 so, diff of all the other values except max value  
 will be zero and diff of max value w.r.t max value  
 will be one.

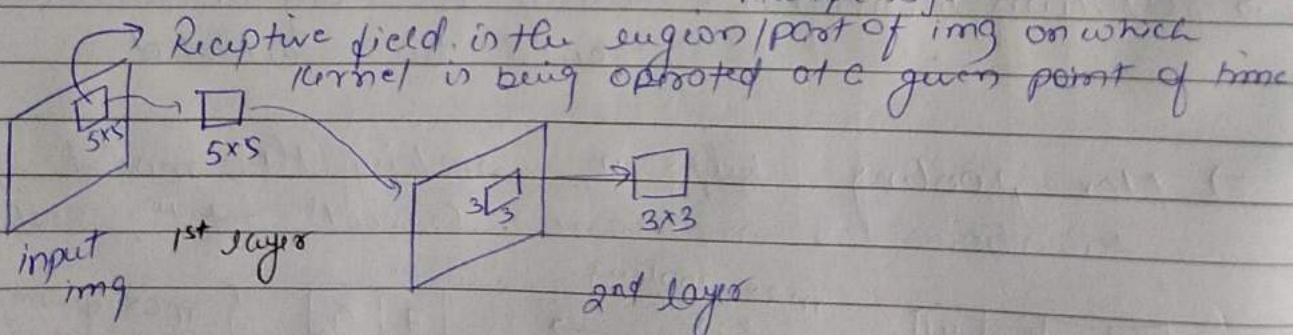
Let  $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$

Let  $b$  is max value

$$\frac{\partial a}{\partial b} = 0, \frac{\partial b}{\partial b} = 1, \frac{\partial c}{\partial b} = 0, \frac{\partial d}{\partial b} = 0$$

## \* Receptive field and effective Receptive field

multiple layer



effective Receptive field of  $3 \times 3$  kernel is the region getting affected in the original / i/p img.

## ✳️ Data Augmentation :-

adding several diff variations of original data  
 e.g. flipping, rotation, zoom in, zoom out version of the data

## ✳️ Why data augmentation?

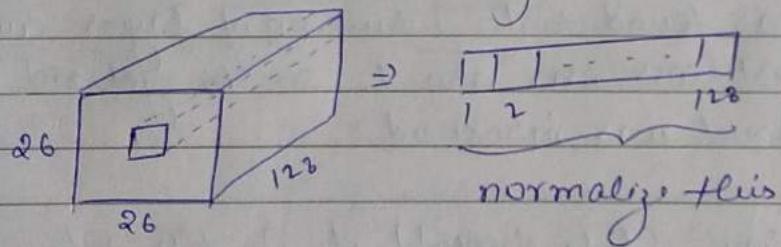
→ invariance (any 3 invariance)

→ from small dataset to large dataset.

## ✳️ AlexNet

- ↳ used ReLU as activation f<sup>n</sup>
- ↳ dropouts were used
- ↳ GPUs were used

→ local response normalization (LRN)

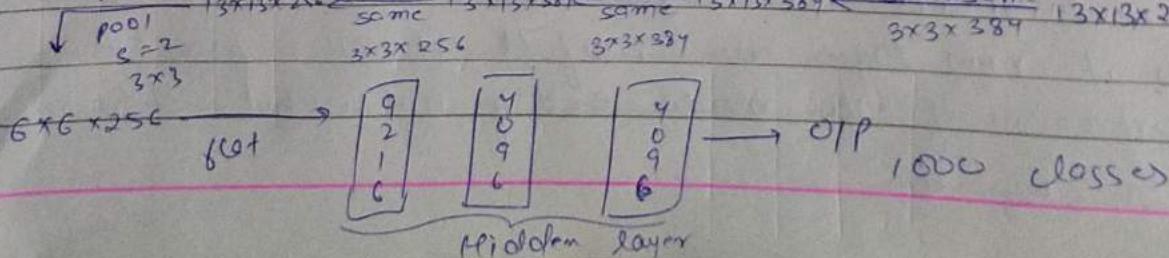
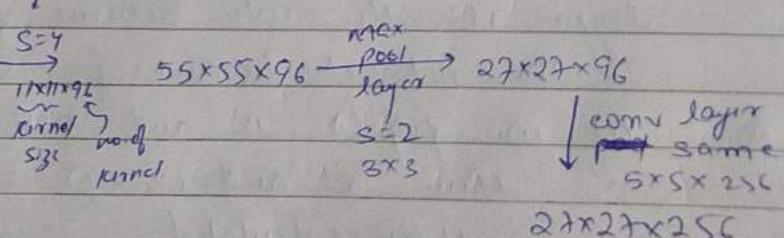


(but not used much, batch norm. is the way to go)

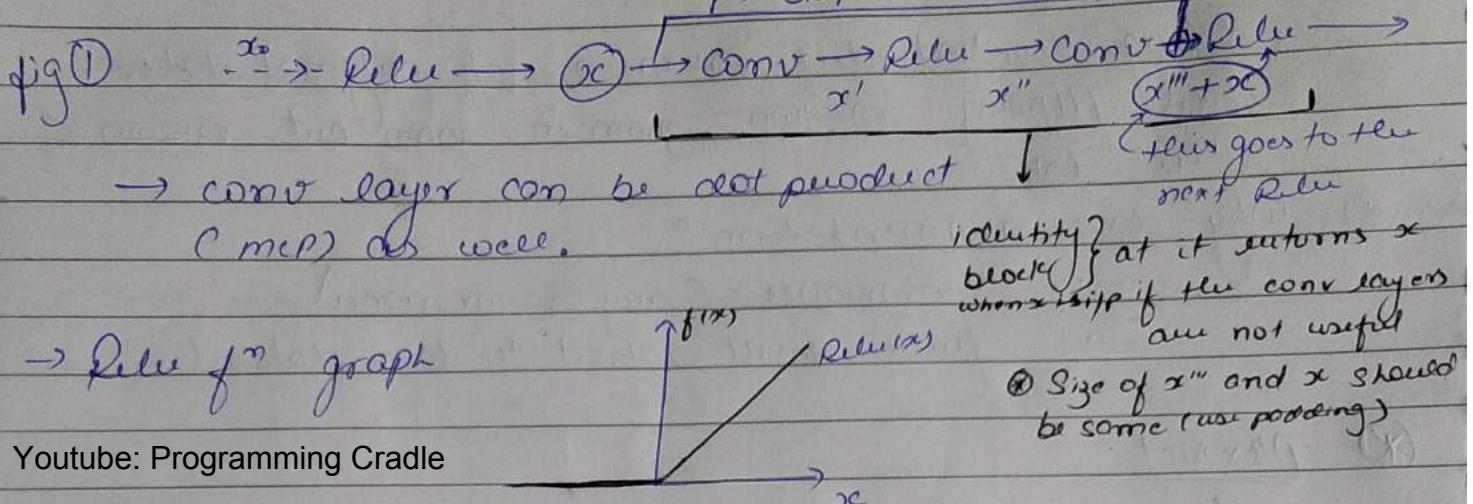
input

$224 \times 224 \times 3 \rightarrow \text{pool}(3,3)$

Youtube: Programming Cradle



\* Residual Network (ResNet) \*



$$\begin{aligned} \textcircled{1} & \left\{ x: +ve \rightarrow \text{ReLU}(x) = x ; \text{ReLU}(\text{ReLU}(x)) = x \right\} \\ \textcircled{2} & \left\{ x: -ve \rightarrow \text{ReLU}(x) = 0 ; \text{ReLU}(\text{ReLU}(x)) = 0 \right\} \end{aligned}$$

from above 2 eq, we say:

$$\text{ReLU}(\text{ReLU}(x)) = \text{ReLU}(x)$$

- ⊕ In Research it was found that when no. of layers are ↑ train and test loss was also ↑. so to get rid of this problem Resnet was introduced
- ⊕ When we use regular (L1, L2, dropout) it is the wts which are not important. and lets suppose in some layers all the wts become zero then it won't affect the network because of skip connection
- ⊕ Let in fig ① 1<sup>st</sup> and 2<sup>nd</sup> conv layers is 0 then  $x'$  and  $x'' = 0 \Rightarrow x''' = 0$  hence  $\text{ReLU}(x)$  will be calculated.
- ⊕ Adding additional / new layers would not hurt performance as regular will skip over them
- ⊕ So if new layers are useful it will improve the performance of your model and if not useful it will be skipped

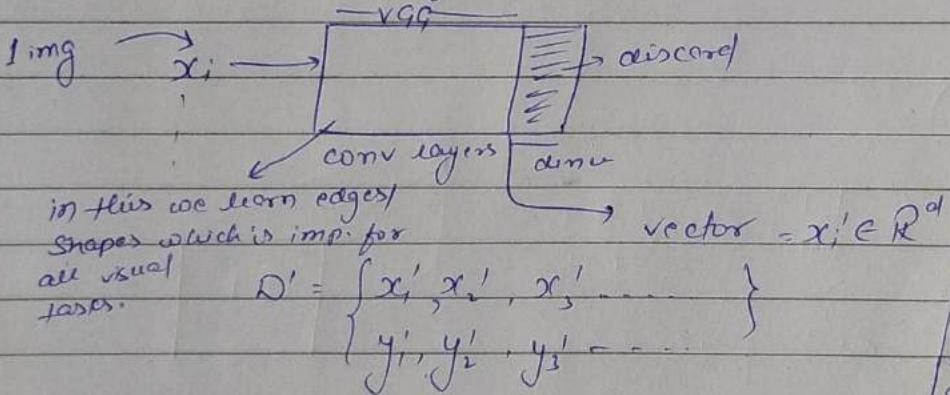
\* Inception Network :- (<http://csfir.ashukumar27.io/cnn-incept^n.html>)

\* instead of deciding whether to use  $1 \times 1$  or  $3 \times 3$  or  $5 \times 5$  conv and whether to use max pool layer, why not use all of them?

\* Transfer Learning :-

→ idea :- Instead of building a NN from scratch to solve our task, we can reuse existing models (VGG16) trained on diff dataset.

\* Case 1 VGG16 + ImageNet



- \* use case 1 when your dataset is very similar to imagenet and size of your dataset is small. ( $\sim 1000$  img.)
- \* use case 2 when your dataset is similar to imagenet and large
- \* use case 3 when your dataset is dissimilar to imagenet (use imnet layers) and dataset is large
- \* when dataset is small and dissimilar use initial few layers off already trained model.

\* Case 2 We fine-tune only all the  $10$  conv. layer but last  $2$  conv. layers, and fine tune last layer using new dataset by keeping the learning rate very small.

\* Case 3 use vgg16 imagenet as initialization and train the whole model using new dataset  $D = \{(x_i, y_i)\}$  by putting small learning rate.