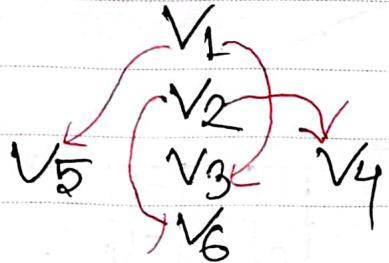


Difference type of layer in CNN

CNN is a part of deep learning which comes into the picture inspiring from, the how the brain works.

In the back part of our head we have Cerebral Cortex and an important portion of this is called Visual Cortex responsible for creating the image in our head.



After getting data through our sensory organ eye, it pass through various layer presented in the visual cortex area. And all those interconnected layers do some operation on the image data and, pass and share information among all of the layers and finally add the end image get created. And all the layers are responsible for detecting

specific part of the image.

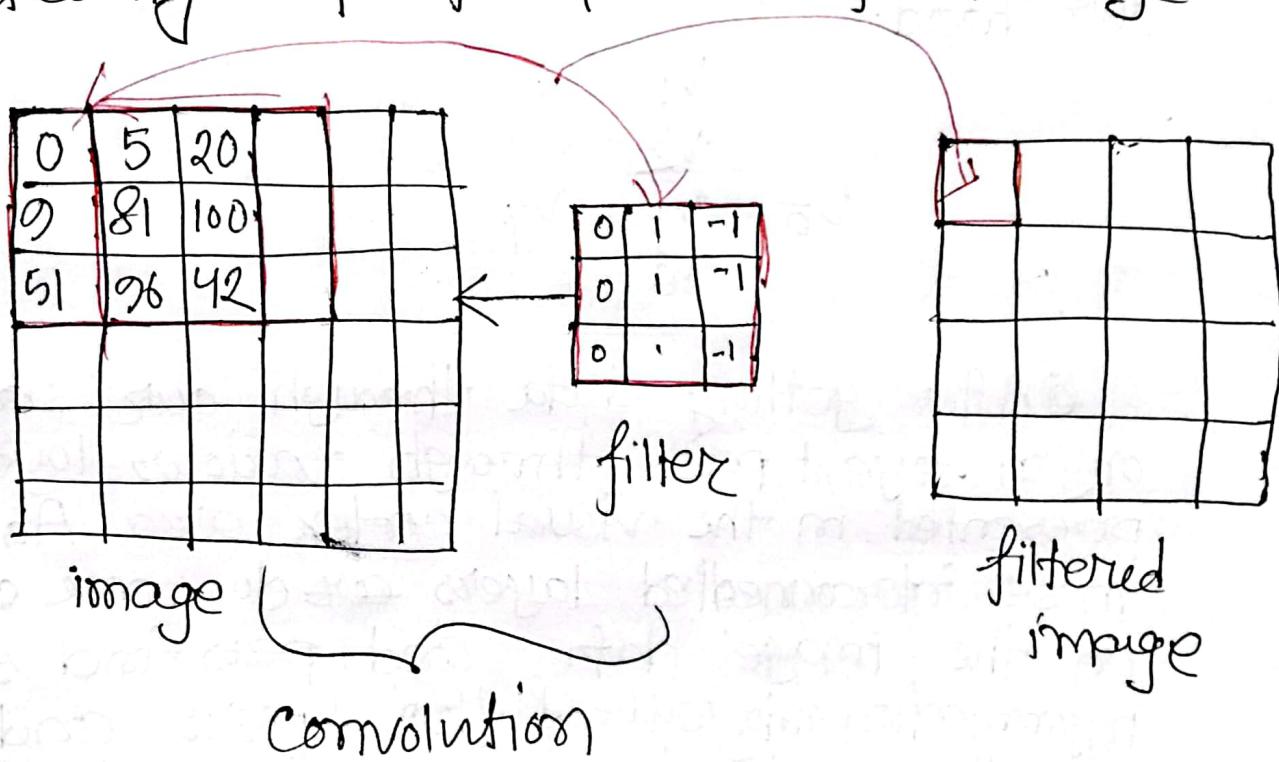
* Convolution

Basically we work with two types of

1) Grey scale (6×6) →

2) RGB scale $(6 \times 6 \times 3)$ →

Filter: are $m \times n$ matrix responsible for detecting specific portion of an image.



* Padding: It is used to retrieve the corner information

* zero padding

* Neighbour padding.

* Stride: It means how much steps we will take both vertically & horizontally.

⇒ Equation for getting output shape of an image after filtering.

$$(\text{Width}, \text{height}) = \left\{ \left(\frac{W-f_w+2P+1}{S} \right), \left(\frac{H-f_h+2P+1}{S} \right) \right\}$$

N.B if we want the output size will be same as input,

$$W-f+2P+1 = W$$

$$\text{hence, } P = (f-1)/2$$

And here, height & width of the filter will be same.

* How CNN is different from ANN.

* Parameter sharing.

* All the weights and (Input value or Neuron's value) do not get multiplied to define the value of a neuron in next layers.

* Why do we CNN over ANN for image classification:

- * Large numbers of trainable parameters.
- * Loss of spatial features.
- * Vanishing & exploding gradient descent.

* What is vanishing Gradient descent:

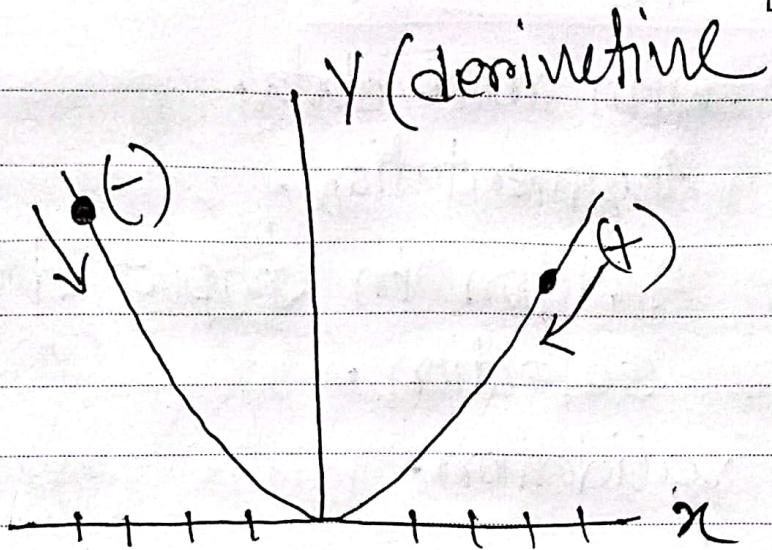
If occurs when weight updation is less, ~~can be~~ and ignorable amount

$$\begin{aligned}\text{new} &= \text{old} - \eta * \frac{\partial}{\partial w} \\ &= 5 - \boxed{.000123} \leftarrow \text{can be ignored}\end{aligned}$$

* What is Exploding Gradient descent:

If occurs when, the summation multiplication of derivative of loss function w.r.t. weights is too large.

$$\begin{aligned}\text{new} &= \text{old} - \eta * \frac{\partial}{\partial w} \frac{\partial}{\partial w} \frac{6}{6L} \\ &= 5 - \boxed{4.123} \leftarrow \text{way too much}\end{aligned}$$



Gradient descent Graph.

- * Why do we need bias? Go to ML Khatua
- * Which pooling technique will you use? same

* Overfitting & Underfitting:

Model get underfitted, when in single layer or very less num Neural Network or which has less number of ~~Neura~~ layer. Underfitted model has high variance bias low variance.

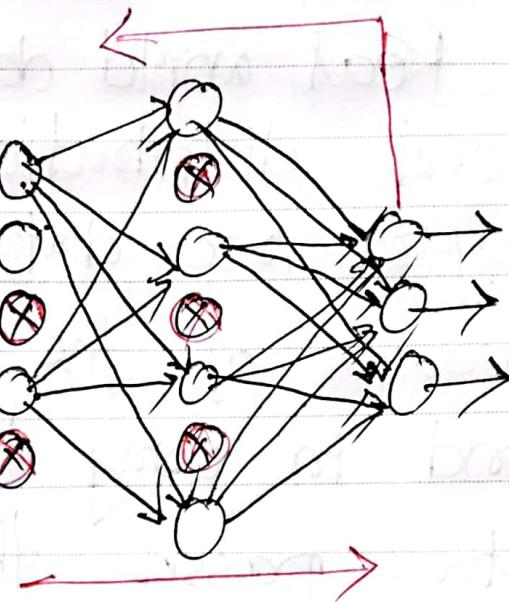
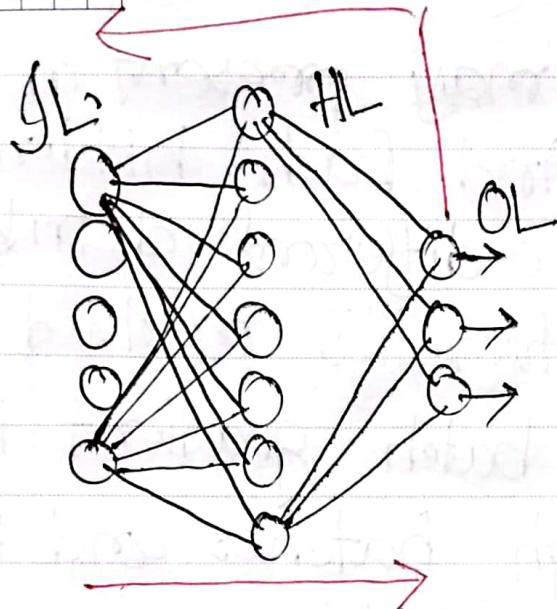
Overfitting simply defines the memorization of training data. And model will have high bias & low variance.

* To reduce the overfit:

- ① Train with more data.
- ② Data Augmentation.
- ③ Noise Addition in Input data.
- ④ Feature selection.
- ⑤ Cross-validation.
- ⑥ Model tuning.
 - * Pruning of Decision tree .(Post & Pre)
 - * L₁, L₂ Regularisation.
 - * Using Dropout in the Model's layers.
- ⑦ Ensembling.
- ⑧ Early stopping.

* Dropout layers:

Dropout means randomly removing Neuron from the hidden layers and the Input layer, which contains ~~at~~ the features of data. In a word random feature and neuron selection in each epoch.



In the forward propagation randomly a proportion of Neuron will be deactivated and in backpropagation weight corresponding to those neuron will not be updated. And this will continue for the n numbers of epochs.

And while backpropagation to testing model, weights of each layer will be multiplied by the corresponding dropout ration of that layer.

Look like Random Forest Algorithm.

* Batch Normalization: (Internal covariant shift)

Real world data may present in different distribution. And each hidden layers see data in a different distribution because of the mathematics involved in it. And in using Mini batch gradient descent data passes through batches and back propagate for each batch and update the weights. But just because of variety in distribution, it's hard to recognise pattern and takes time to converge. And that's why Batch normalization comes into the picture for normalizing data in each layer and keep same distribution throughout the process.

~~Basic Mathematics of Batch Normalization~~

- ① Calculate mean value \bar{x} for all the data points in a batch, as well as the standard deviation.

$$(1) \mu = \frac{1}{n} \sum z^{(i)}$$

$$(2) \sigma^2 = \frac{1}{n} \sum (z^{(i)} - \mu)^2$$

(2) Then it normalizes the activation vector, in this way all the neuron's output follows a normal distribution

(3) Then finally it calculate the layer's output by applying a linear transformation with by adding two new variables (trainable) α, β to introduce non linearity.

$$(3) z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad (4) Z = \alpha * z_{\text{norm}}^{(i)} + \beta$$

* Loss function or Cost function

Loss function is used to define individual losses & Cost function for the cumulative loss errors.

→ Different types loss function are there

(1) Logloss or Binary Cross Entropy.

It's used in logistic Regression

Disopan®
Clonazepam tablet & injection

- (2) Sparse-Categorical cross entropy.
- (3) Categorical crossentropy.
- (4) Mean Absolute Error.
- (5) Mean Squared Error.
- (6) Root Mean Square Error.
- (7) R² squared Error.
- (8) Adjusted R² squared error.

* Python Interview Question.

what is Zip function.

zip method takes multiple iterable or containers and return a tuple(iterator object) having mapped values from all the containers

* Few confusing terms in Computer Vision
Image Classification, Object Recognition.
Object detection, Object Segmentation,
Object verification.

⇒ **Image Classification/ Image Recognition**
Means classifying whole image into a category.

⇒ **Image or Object localization:**
Means Localizing the main object in the Picture image alongside the classification.
When Image will contain only 1 object.

⇒ **Object Detection**

If specifies the location of Multiple object along with it's class.

⇒ **Object Segmentation**

It's all about segmenting the pixels of an image, having similar attributes. And another important point is it's gives us the shape of the object.

Image Recognition Project

- ① Loading Necessary library.
- ② Loading the Image
- ③ Finding the feature location from Actual image.
- ④ Getting feature or the encoding of Image.
- ⑤ Repeat ③ & ④ for the test image also
 ^{cont}
- ⑥ And compare those image using the
 face-recognition library.
- ⑦ Also can find the distance between those.

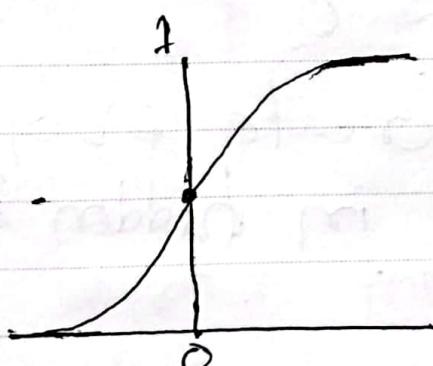
Object detection Algorithm

R-CNN, Fast-CNN, Faster-CNN, YOLO (Have to Read)

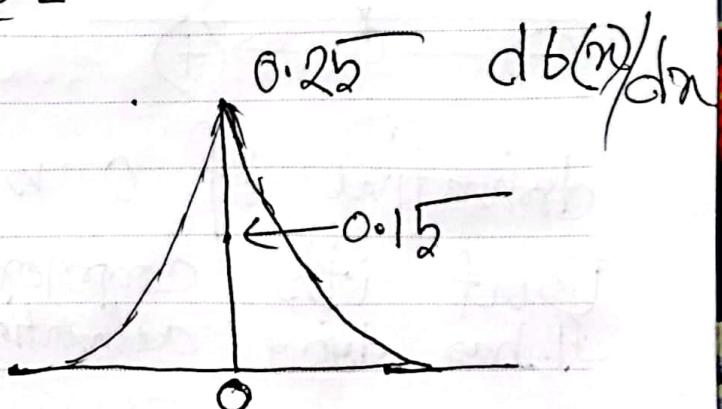
Activation Function

① Sigmoid

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



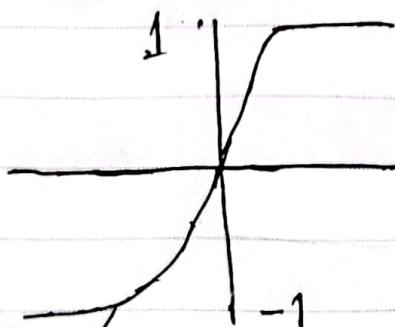
Value (0, 1)



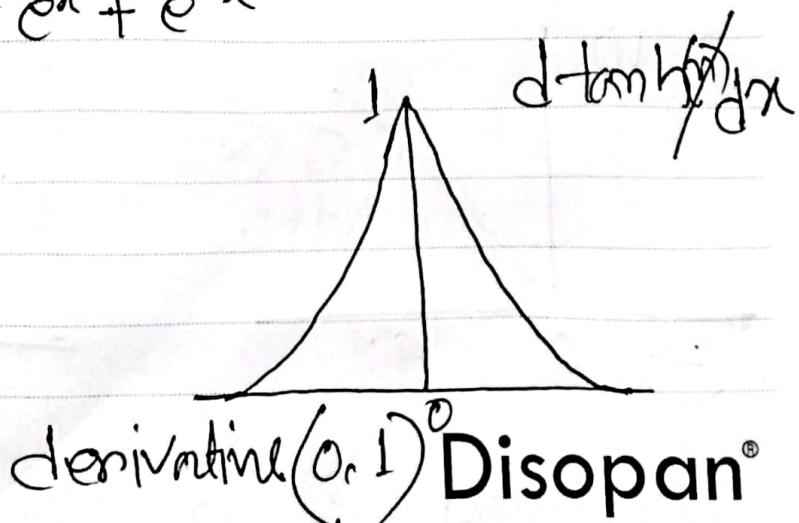
derivative (0, 0.25)

② tanh

$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



inceptia

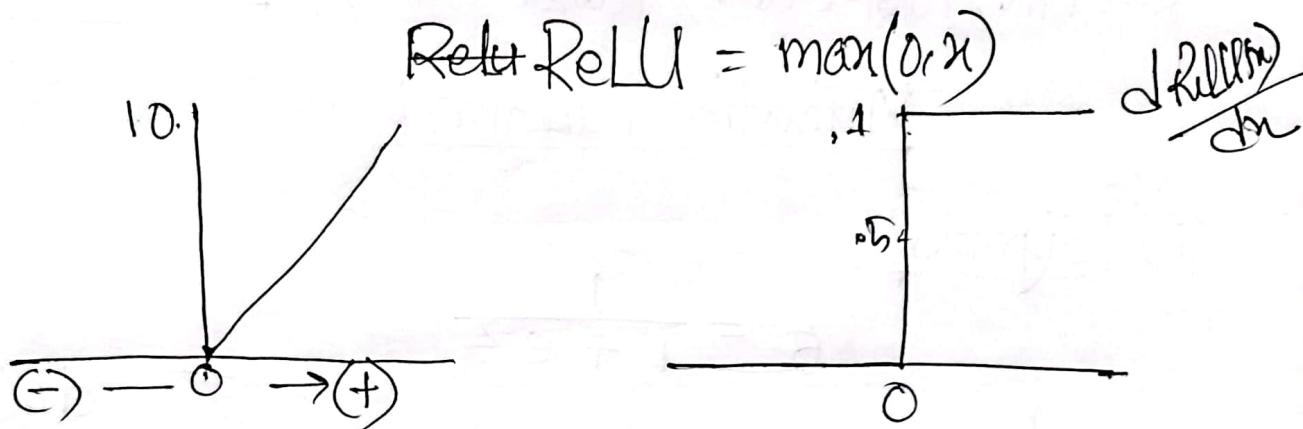


derivative (0, 1) Disopan®

Clonazepam tablet & injection

Adove both as vanishing Gradient descent problem,

③ Relu

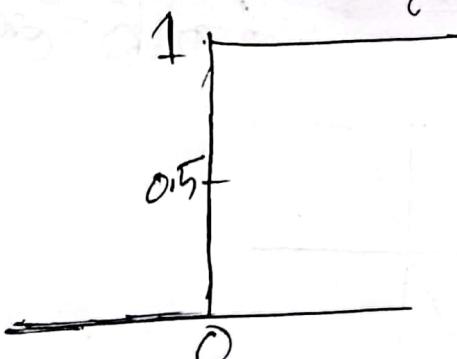
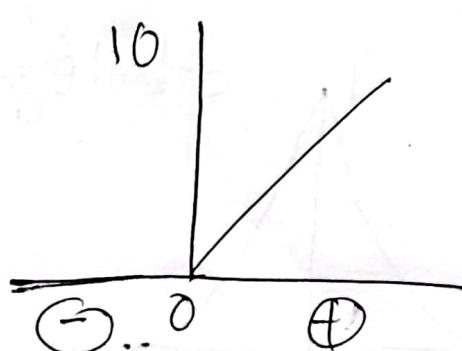


derivative of 0 is 0 and create $\partial V \partial y$.
But it's commonly used in hidden layer.
It has dying activation problem.

④ Leaky Relu

$$f(x) = \max(0.1x, x)$$

$$\frac{df}{dx}$$



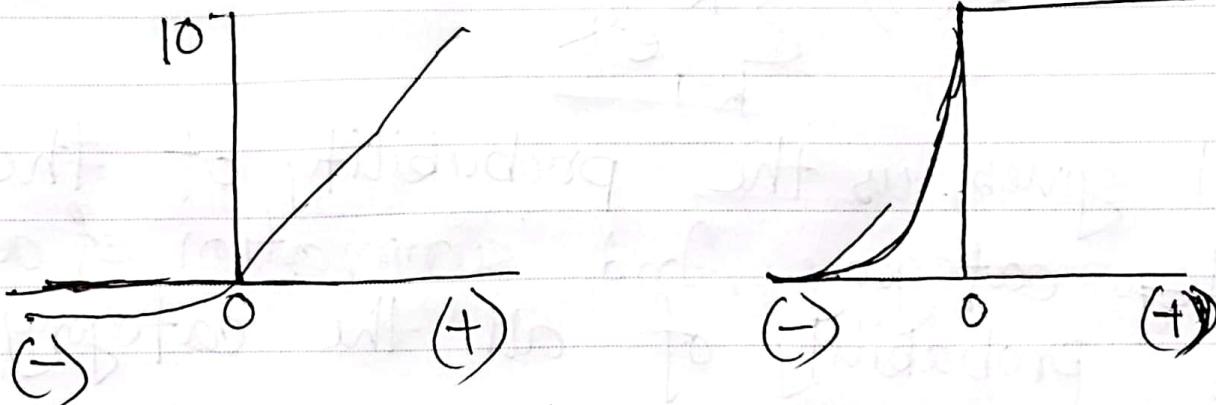
It add a smaller value to negative value to eradicate dying activation problem.

But Leaky ReLU create the vanishing gradient descent problem again

⑤ ELU:

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{otherwise} \end{cases}$$

$d f(x)/dx$



$\alpha(e^x - 1)$ this gives a little higher value than leaky ReLU and derivative of it also higher. thus handle the over GD problem.

But it takes time computationally expensive

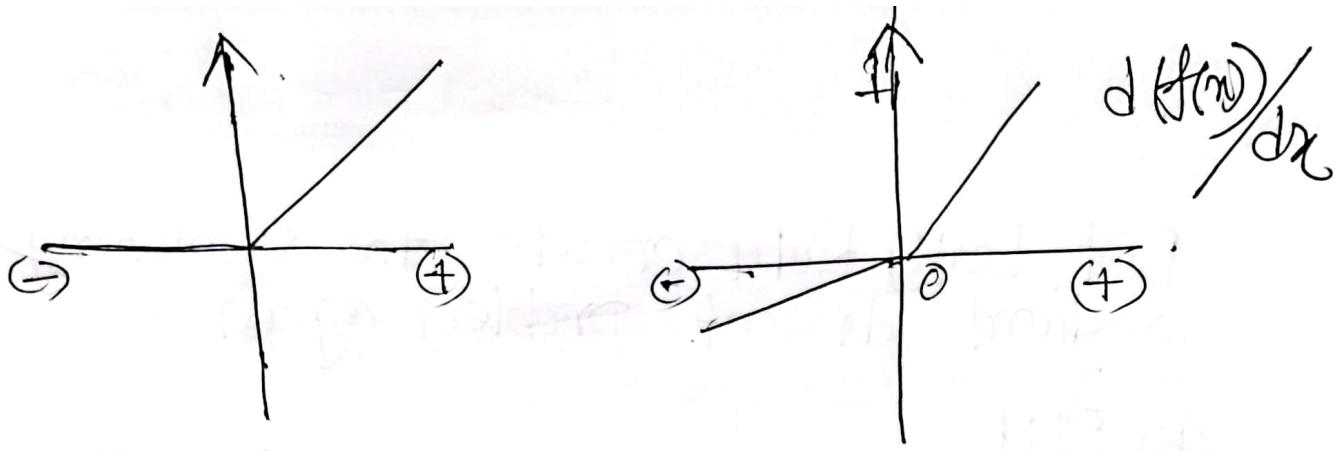
⑥

⑥ PReLU (Parametric ReLU):

$$f(x) = \begin{cases} x & x > 0 \\ \alpha x & \text{otherwise} \end{cases}$$

if $\alpha = 0.01$ it becomes Leaky ReLU

Here α is learning Parameter **Disopan®**



⑦ Softmax

$$S(x_j) = \frac{e^{x_j}}{\sum_{k=1}^K e^{x_k}}, j=1, 2, 3 \dots K$$

It gives us the probability of the category. And summation of all the probability of all the category is 1.

Assume last layer has value something like this.

$$[40, 30, 20, 10]$$

$$P(40) = \frac{e^{40}}{e^{40} + e^{30} + e^{20} + e^{10}}, P(30) = \frac{e^{30}}{e^{40} + e^{30} + e^{20} + e^{10}}$$

To To

✓

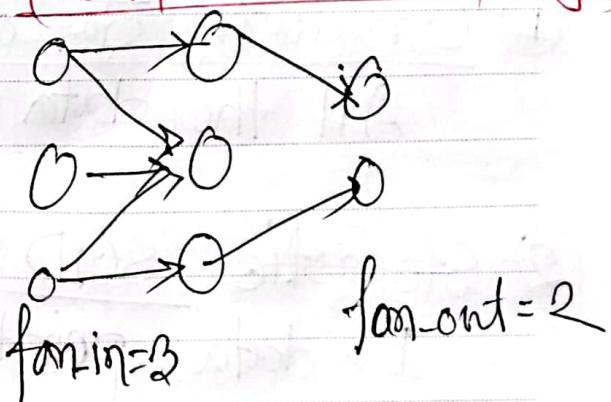
Weight Initialization Technique

- ① Weights must be balanced.
- ② Weights should not be same.
- ③ Weights should have good variance.

Various Techniques

① Uniform Distribution: (Works well with Sigmoid)

$$W_{ij} = \left[\frac{-1}{\sqrt{fan_in}}, \frac{1}{\sqrt{fan_out}} \right]$$



② Xavier / Glorot Normal (Works well with Sigmoid)

- ① Xavier Normal
- ② Xavier Uniform

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

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

$$W_{ij} = \left[\frac{-6}{\sqrt{fan_in + fan_out}}, \frac{6}{\sqrt{fan_in + fan_out}} \right]$$

(3) He init (works well with ReLU activation)

(i) He uniform

$$w_{ij} \approx \left[-\frac{6}{\sqrt{f_{min}}}, \frac{6}{\sqrt{f_{max}}} \right]$$

(ii) He. Normal

$$w_{ij} \approx N(0, \sigma^2)$$

$$\sigma = \sqrt{\frac{2}{f_{min} + f_{max}}}$$

Optimizers

① Gradient descent:

All the data at a time & time killing to take
take so much resources

② Stochastic SGD:

1 data point at a time.
take too much time

③ Mini batch SGD:

Batch of data, take less time and
converge faster but noisy.

④ SGD+Momentum:

SGD+momentum to remove noise.

⑤ Adadelta and RMSprop \Rightarrow Adamgrad

6 Adadelta and RMSprop.

7 Adam (RMSprop + Momentum)

Adam is most used Optimizer.

Tensor

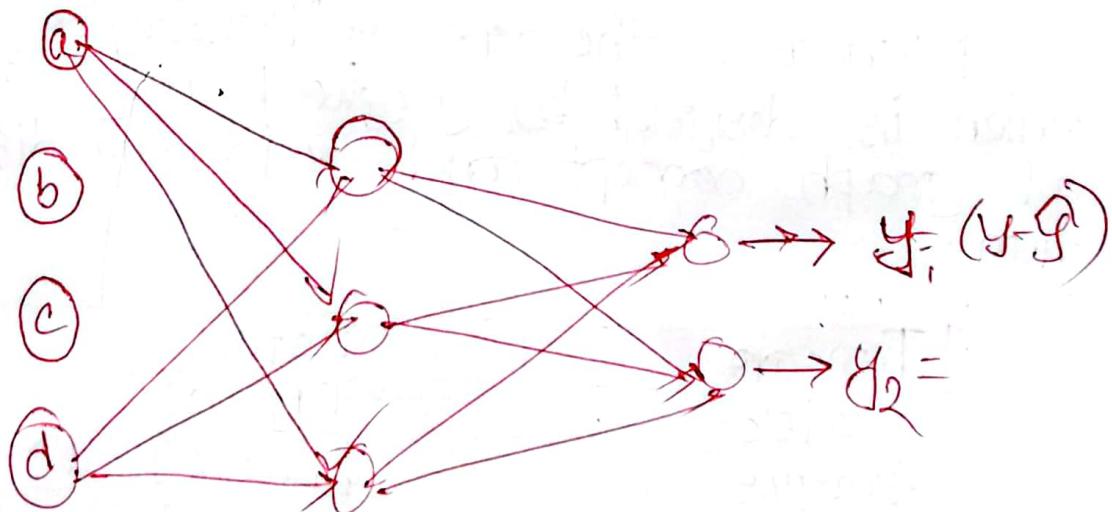
Tensor: Tensor is a vector or matrix of n-dimensional that represent all types of data. All values in a tensor hold identical data type with a known shape. The shape of the data is the dimensionality of the matrix or array.

* Tensor is the core of Tensorflow library, which is designed by google having computational graph concept and automatic differentiation.

Tensor	Data
Scalar	$[.]$ (1D)
Arrays	$\{[1, 2, 2]\}$ (1D)
Matrix	$[[1, 2, 3], [4, 5, 6]]$ (2D)
Multidimensional or ndarray	$[[[1], [2], [3]], [[1], [2], [3]]]$ (3D)
N dimensional tensor	

?Tensor contains:

- ① Input Data.
- ② Store the output of set of computation take place inside a graph (known as Computational graph)
- ③ Keep the values of weights and bias biases update while backpropagation.
- ④ And it also take care of back propagation with automatic differentiation also known as chain rule.



Computation graph

Most Important (Chain Rule)

Interpreted differentiation.]

Cosine Similarity & distance

① Minkowski distance:

$$D = \left(\sum_{i=1}^N |P_i - P_2|^P \right)^{1/P}$$

And from this Equation, can derive Manhattan and Euclidian distance.

Manhattan: When $P=1 = |x_1 - x_2| + |y_1 - y_2|$

Euclidian: When $P=2 = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$

Both are Good but when features increases may create complexity.

② Hamming distance:

~~If finds the similarity between two equal size strings. And check characters by character and do bit wise operation~~

Sudha

Sudda

11101

$$\text{So, HD} = (1+1+1+1) = 4$$

Disopan®
Clonazepam tablet & injection

$$(\sin \theta)_{\text{ref}} = \cos \theta \Leftrightarrow$$

$$\sin \theta = \cos \theta \Leftrightarrow$$

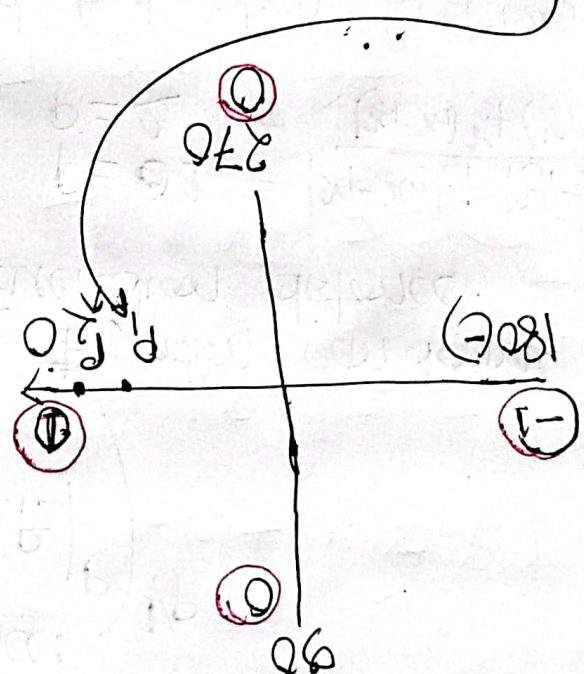
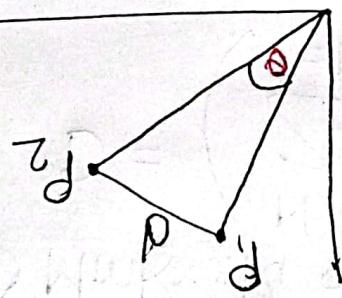
$$\cos \theta = \frac{\|A\| * \|B\|}{A \cdot B} = \boxed{\text{Cosine Similarity}(A, B)}$$

So, this two points are more similar.

$$\cos \theta = \cos 0^\circ \therefore$$

$$\theta = 0^\circ$$

between (x_1, y_1)
and the origin
with respect to
the x-axis



$$1 - \cos \theta = \sin \theta = \cos(180^\circ - \theta)$$

③ Cosine Similarity & Distance:

Recurrent Neural Network (RNN)

⇒ RNN is a special type of NN deals with Sequential type of data by remembering the previous history, which includes a hidden state to remember the past data.

⇒ But it failed to remembering long data, and also vanishing gradient descent problem.

⇒ And then came the LSTM which resolves this issue by adding a cell state alongside hidden state to keep important past information, containing three main gates. And hidden state will be for output the result and cell state is for keeping important information.

⇒ Forget Gate: Responsible for deciding which information we need to be removed.

$$f_t = \sigma(x_t w_f + h_{t-1} w_f + b_f)$$

And sigmoid activation function is used on top it.

Disopan
Incepta
tablet & injection

between (0~1). Or mean remove information.

⇒ Input Gate: It's responsible for what ^{new} information will be stored in the cell state, with the help of candidate state,

$$I_f = \sigma(u_i + h_{t-1}u_i + b_i)$$

$$g_f = \tanh(u_i w_g + h_{t-1} w_g + b_g)$$

$$\boxed{\text{Input}} = I_f * g_f$$

⇒ Output Gate: For deciding which information of cell state will go to output, by using tanh on Input

$$Cell_{in} = \tanh(\text{Input})$$

And finally updating hidden state

$$h_t = \sigma(u_h + h_{t-1}u_h + b_t) * Cell_{in}$$

$$\text{And } \gamma = \text{SoftSigmoid}(h_t)$$

$$\text{And loss} = (\gamma - \hat{y}) \text{ then backpropagation}$$

⇒ BiLSTM Depends on the Knowledge of both the left and right side for prediction

⇒ GRU, BiGRU → (ମାତ୍ର ଦେଖିଲୁ ଅନ୍ତର୍ଭୟ)

Seq to Seq

- ① Encoder-Decoder (BLEU Score)
- ② Attention
- ③ Transformer
- ④ BERT