

Deep Learning for Computer Vision

NPTEL - DL - Dr. Mitesh Khapra

Deep Learning is a subset of AI but fundamentally different from ML

McCulloch Pitts Neuron

- simplest neuron
- input and output data is binary

1. Aggregation
2. Threshold on aggregation

A neuron ~~has both~~
is mathematically
 $f(g(x))$
threshold aggregation

Excitatory Neuron

$$\text{Aggregation} \Rightarrow x_1 + x_2$$

$$\Rightarrow x_1 + x_2 + \dots + x_n$$

$$\Rightarrow g(x) = \sum_{i=1}^n x_i$$

for an AND gate,



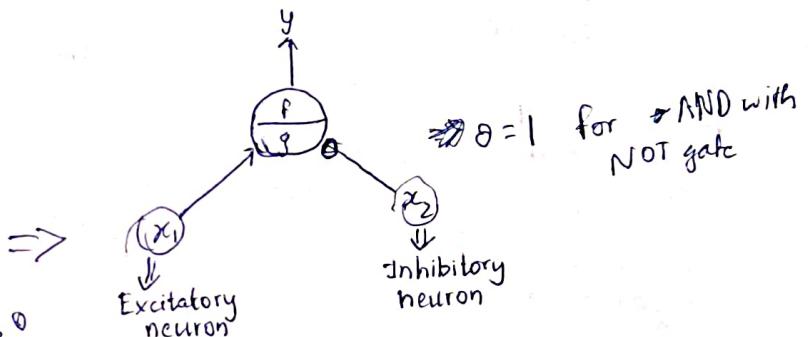
$$f(x) \Rightarrow \begin{cases} 1 & x \geq \theta \\ 0 & x < \theta \end{cases}$$

where θ is the threshold.

for an AND^{NOT} gate $\Rightarrow x_1 \cdot \bar{x}_2$

x_1	x_2	$x_1 \cdot \bar{x}_2$
0	0	0
0	1	0
1	0	1
1	1	0

whenever x_2 is 1, the output is 0
 $\Rightarrow x_2$ is inhibitory neuron



Inhibitory Neuron

- ⇒ if the inhibitory neuron returns an output of 1, the result is always 0.
- it doesn't participate in aggregation

The type of neuron is identified by its associated weight.

for excitatory neuron, weight is 1

for inhibitory neuron weight is $\theta - 1$

Perceptron

For an XNOR gate,

x_1	x_2	$x_1 \otimes x_2$
0	0	0
0	1	0
1	0	0
1	1	0

$$y(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n = \sum_{i=1}^n w_i x_i$$

$$\sum_{i=1}^n w_i x_i - \theta = g(x) \Rightarrow$$

$$f(x) = \begin{cases} 1 & x \geq \theta \\ 0 & x < \theta \end{cases}$$

$$g(x) = \sum_{i=0}^n w_i x_i \quad \text{where } w_0 = \theta \text{ and } x_0 = 1$$

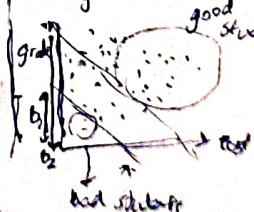
in a perceptron, x_2 inputs can be any real value.

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}, x = \begin{bmatrix} x_0 & x_1 & \dots & x_n \end{bmatrix} \rightarrow n+1 \text{ matrix}$$

$\rightarrow 1 \times n \text{ matrix}$

understanding of bias

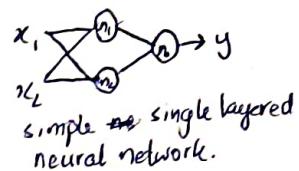
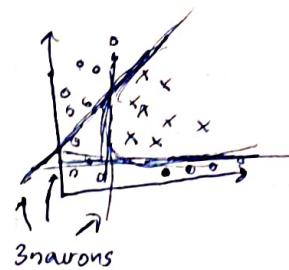
Suppose a college is taking students for admission based on grades and test results, new colleges may be alright with lower values than the good ones



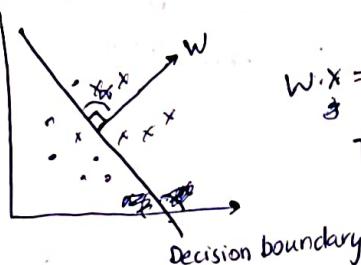
Linear separability

- A single neuron can only separate 2 classes with a straight line
- in order to separate 2 classes with a curve, multiple lines are used by combining multiple neurons and forming a neural network

$$\begin{aligned} w_0 &= -1 \\ w_1 &\propto 1 \\ w_2 &\propto 1 \\ w_1 + w_2 &\geq 1 \end{aligned}$$



Learning Logic



$w \cdot x = 0 \Rightarrow$ decision boundary. (w^\perp is orthogonal to decision boundary)
The decision boundary is the line (imaginary) which classifies points into class 1 or class 2.

if P is positive but classified as negative then $w = w + x$ (repositioning the D.B. to make the correct classification)

$\cos\alpha = \frac{x \cdot w}{\|w\| \|x\|}$ \Rightarrow new $\angle \alpha'$
angle b/w w and w . $\angle \alpha'$ is less than the previous angle.
~~D.B.~~ on updating w . $\alpha_{new} > \alpha$ for when P is -ve but classified as positive.

Perceptron learning algorithm

$$\text{Aggregate} = \sum_{i=1}^n w_i x_i + \text{bias} = y_{in}$$

Activation function \Rightarrow Step activation $\begin{cases} 1 & y_{in} \geq \theta \\ 0 & -\theta \leq y_{in} \leq \theta \\ -1 & y_{in} < -\theta \end{cases}$

compare if $y \neq t$
 $w_{new} = w_{old} + \alpha x \cdot t$
 $bias_{new} = bias_{old} + \alpha \cdot t$

$\alpha \Rightarrow$ learning rate
 $w \Rightarrow$ randomly initialised weights
 $t \Rightarrow$ target result / expected result.

even though the error is not considered, and simply the expected output is used to adjust the weights, the perceptron will converge as the inputs are meant to be simple t can only be +1 or -1. (bipolar-class classification)

⇒ Implement AND with bipolar inputs and outputs, assume $\alpha=1$, $w_{old} = [0, 0, 0]$

x_1	x_2	$a=1$	t	y_{in}	y	Δw_1	Δw_2	Δb	w_1	w_2	b
1	1	1	1	0	0	1	-1	1	1	1	1
1	-1	1	-1	1	1	-1	1	-1	1	1	-1
-1	1	1	-1	2	1	1	-1	=1	0	2	0
-1	-1	1	-1	-3	-1	0	0	0	1	1	-1
1	1	1	1	1	1	0	0	0	1	1	1
1	-1	1	-1	-1	-1	0	0	0	1	1	-1

Q> AND NOT function ; bipolar inputs and outputs

$x_1 \cdot x_2$

$x_1 \cdot x_2$	x_1	x_2	α	t	y_{in}	y	Δw_1	Δw_2	Δb	w_1	w_2	b
	1	1	1	-1	0	0	-1	-1	-1	0	0	0
$w_{new} = w_0 + \alpha x_1 t$	1	-1	1	1	-1	-1	1	1	1	-1	-1	-1
$\Delta w = \alpha x_1 t$	-1	1	1	-1	-2	-1	0	0	0	0	-2	0
$\Delta b = \alpha t$	-1	-1	1	-1	2	1	1	1	-1	0	-2	0
	1	1	1	-1	-1	-1	0	0	0	1	-1	-1
	1	-1	1	1	1	1	0	0	0	1	-1	-1
	1	1	-1	-3	-1	0	0	0	0	1	-1	-1
	-1	-1	1	-1	-1	-1	0	0	0	1	-1	-1
	-1	-1	-1	-1	-1	-1	0	0	0	1	-1	-1

$$\text{when } = w_0 + \alpha x_1 t$$

$$\Delta w = \alpha x_1 t$$

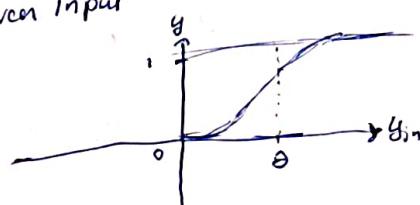
$$\Delta b = \alpha t$$

Drawback of step functions \Rightarrow There is no explainability on why the output is 0 or 1 or -1
 it is difficult to understand

- There is a hard-threshold at a given input

Hence, we begin using Sigmoid input function.

$$\text{sig}(x) = 0 \text{ to } 1 = \frac{1}{1 + e^{-x}} = \frac{1}{1 + e^{-(w_1 x_1 + b)}}$$



Loss function is calculated instead of only checking for $y = t$.
 Another reason for not using step function is the difficulty in convergence.

Vanishing gradient $\Rightarrow \Delta y \approx 0$ for $\Delta y_{in} > 0$

Sigmoid function gives probability of a test case being classified as positive
 i.e. suppose the sigmoid output is 0.7 for a test case \Rightarrow There is a 70% chance of the point being classified as belonging to positive class

Probability of being in negative class is taken as $1 - P$

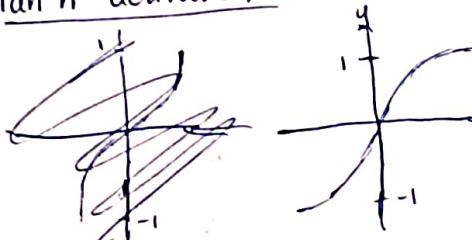
Softmax function is an alternative to sigmoid which is used for multi-class classification

$$\hookrightarrow \frac{e^{z_i}}{e^{z_1} + e^{z_2} + \dots e^{z_n}}$$

where $z_i \Rightarrow$ probability for i^{th} class.

(Sigmoid function is not to be used for multi-class classification)

Tan h activation



This activation function is used when even the neutral (0) is considered in classification

- Doesn't give probability

- computationally intensive

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

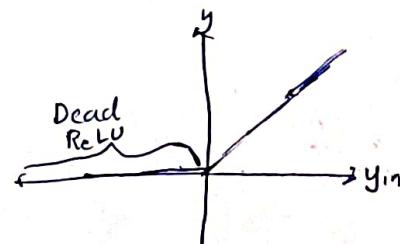
- Vanishing Gradients is still an issue

Rectified Linear Unit (ReLU)

$$\text{ReLU}(x) = \max(0, x)$$

- no vanishing gradient
- computationally simple

ReLU is used when lot of layers are involved due to simplicity in computation.



SiLU \Rightarrow sigmoidal + ReLU

\hookrightarrow simple computation
used in recent YOLOS

Loss Function

- Binary cross-entropy loss / or log loss = $\sum y_i \log p_i + (1-y_i) \log (1-p_i)$
- \hookrightarrow loss is either 0 or 1
- Mean Absolute Error
- Mean Squared Error
- Root Mean Squared Error
- Categorical cross entropy

p_i - class classification

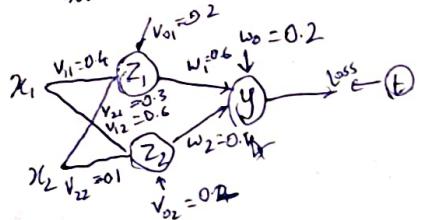
$\sum y_i \log p_i \rightarrow$ multiclass classification

Gradient Descent

$$\Delta w = \frac{\delta L}{\delta w}$$

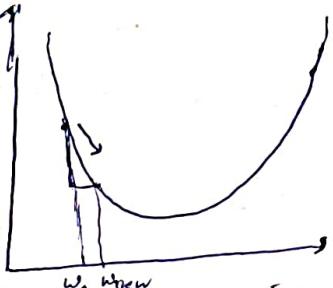
$$w_{\text{new}} = w_{\text{old}} - \alpha \Delta w$$

$$b_{\text{new}} = b_{\text{old}} - \alpha \Delta b$$



$$\text{M.S.E.} = \frac{1}{2} (t - y)^2$$

$$\text{Sigmoid} \Rightarrow \sigma(x) = \frac{1}{1 + e^{-x}}$$



$\Delta w < 0$ $[w_{\text{new}} - w_1]$
 \hookrightarrow to reach w_{new} from w_1 , we need to do subtraction

$$\begin{aligned} Z_{1in} &= x_1 v_{11} + x_2 v_{21} + v_{01} \times 1 \\ &= 1 \times 0.4 + 1 \times 0.3 + 0.2 \\ &= 0.3 \Rightarrow Z_1 = \sigma(0.3) = 0.5744 \end{aligned}$$

$$x_1 = 1$$

$$x_2 = 1$$

$$\begin{aligned} Z_{2in} &= x_1 v_{12} + x_2 v_{22} + v_{02} \times 1 \\ &= 1 \times 0.1 + 1 \times 0.6 + 0.2 \times 1 \\ &= 0.7 \Rightarrow Z_2 = \sigma(0.7) = 0.6684 \\ y_{in} &= 0.5744 \times 0.6 + 0.6684 \times 0.4 + 1 \times 0.2 = 0.8788 \\ y &= \sigma(0.8788) = 0.7065 \end{aligned}$$

Compute values to 4 decimal places

suppose $t = 1$

$$\begin{aligned} \text{Loss} &= \frac{1}{2} (t - y)^2 \\ \Delta w &= \frac{\delta L}{\delta w} \Rightarrow \frac{\delta (\frac{1}{2} (t - y)^2)}{\delta w} \end{aligned}$$

there is no 'w' term here, but in y , there is a y_{in} and in y_{in} there is a w term

$$\frac{\delta L}{\delta w} = \frac{\delta \text{Loss}}{\delta y} \times \frac{\delta y}{\delta y_{in}} \times \frac{\delta y_{in}}{\delta w}$$

$$\frac{\delta y}{\delta y_{in}} = \frac{\delta (\frac{1}{1+e^{-y_{in}}})}{\delta y_{in}} = \frac{\delta (e^{-y_{in}})}{(1+e^{-y_{in}})^2} = \frac{e^{-y_{in}}}{(1+e^{-y_{in}})^2} \Rightarrow \frac{1}{1+e^{+y_{in}}} \times \frac{-e^{+y_{in}}}{(1-e^{-y_{in}})^2} = \frac{y_{in}}{(1-y_{in})^2}$$

$$\frac{\delta (\frac{1}{2} (t - y)^2)}{\delta y} = 2 \times \frac{1}{2} \times -(t - y)(-1) = -t + y - t$$

$$= \frac{e^{-y_{in}}}{(1+e^{-y_{in}})^2} \Rightarrow \frac{1}{1+e^{+y_{in}}} \times \frac{-e^{+y_{in}}}{(1-e^{-y_{in}})^2} = \frac{y_{in}}{(1-y_{in})^2}$$

$$\frac{\delta y_{in}}{\delta w} = \frac{\delta (w_1 z_1 + w_2 z_2 + w_0)}{\delta w_i}$$

= ~~w_1, w_2, w_0~~ ~~Based on which weight we are updating, rest are considered as constants.~~

$$\Rightarrow \frac{\delta L}{\delta y} \times \frac{\delta y}{\delta y_{in}} \times \frac{\delta y_{in}}{\delta w} = (y-t)(y)(t_1-y)(z_1)$$

⇒ $(0.7065 - 1) (0.7065) (1 - 0.7065) (0.5744)$

in example, ⇒ $(0.7065 - 1)$

$$\Rightarrow \Delta w_1 = -0.0349$$

$$w_1 = 0.6 - 0.25(-0.0349)$$

$$= 0.6087$$

$$\Delta w_2 = (0.7065 - 1)(0.7065)(1 - 0.7065)(0.6684)$$

$$= -0.04067$$

$$w_2 = 0.5 - 0.25(-0.04067)$$

$$= 0.5102$$

$$\Delta w_0 = -0.0608$$

$$w_0 = 0.2 - 0.25(-0.0608)$$

updating $v_{11}, v_{12}, v_{21}, v_{22}$

$$\frac{\delta L}{\delta v_{11}} = \frac{\delta L}{\delta y} \times \frac{\delta y}{\delta y_{in}} \times \frac{\delta y_{in}}{\delta z_1} \times \frac{\delta z_1}{\delta v_{11}} \times \frac{\delta v_{11}}{\delta x_1}$$

w_1 $(z_1 \times 1-z_1)$ x_1

$$= (y-t)(1-y)(y)(w_1)(z_1)(1-z_1)(x_1)$$

$$\Rightarrow \Delta v_{11} = -0.0023$$

$$v_{11} = 0.4 - 0.25(-0.0023) =$$

$$\Delta v_{11} = (y-t)(1-y)(y)(w_1)(z_1)(1-z_1)(x_1) = -0.0023$$

03/02/2025

- There is no difference b/w ANN and DNN except that a DNN has many more hidden layers.
 - All YOLO models are trained on COCO datasets.
 - Take a pretrained model for any purpose, as they are guaranteed to have good feature extractors.
 - Remove its output layer since we want a different classification than what it was originally trained for. Then we replace it with our own neural network.
 - Fine tuning a model for a purpose ⇒ appending your own neural network to the output of the pretrained model. (Transfer Learning and fine Tuning)
- Add a good number of layers, not just 1 or 2.

Don't use updated values for updation anywhere until all values are updated.

(similar differentiation to)
 $\frac{\delta y}{\delta v_{11}}$

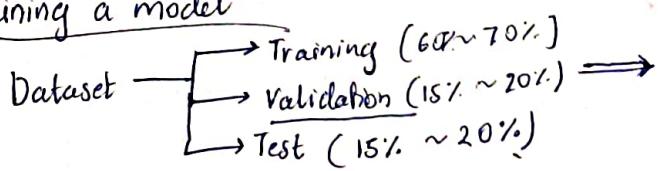
when appending your own neural network,

* & never keep sudden decrease in number of neurons.

i.e. - Pretrained model \rightarrow 1024 neurons \rightarrow 2048 neurons \rightarrow 1024 neurons \rightarrow 2 neurons

sudden convergence will cause difficulty in learning.

Training a model



Bias \Rightarrow difference b/w mean of expected value and calculated mean

$$= E(f(x)) - f(x)$$

Theoretically calculated mean.
i.e.

Variance \Rightarrow how much the calculated values/lines vary from the expected line.

$$= E((f(x) - E(f(x)))^2)$$

For each epoch,

1. Training on 1 epoch

2. validation on validation data

- continue till training loss increases but validation is decreasing \rightarrow overfitting.

Early stopping

Avoiding overfitting has many different methods.

A fundamental technique used to avoid overfitting.

Dropouts \Rightarrow Randomly dropping out neurons based on some probability.

It is to avoid problem of very few of the neurons learning the correct structure while rest remain useless. This method causes neurons to forcefully learn.

Handling Overfitting [based on number of epochs]

early stopping

Dropouts

Data augmentation

Ensemble learning

Adding noise

Regularization

Batch Normalization

When data is small and doesn't contain enough variations, the model is likely to overfit. So, in order to increase variations, properties of the image like contrast, brightness, scale, etc. are changed.

The chosen data augmentation needs to be appropriate to the dataset.

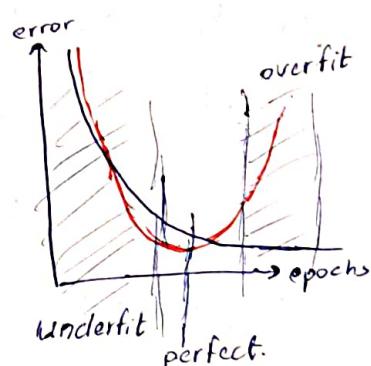
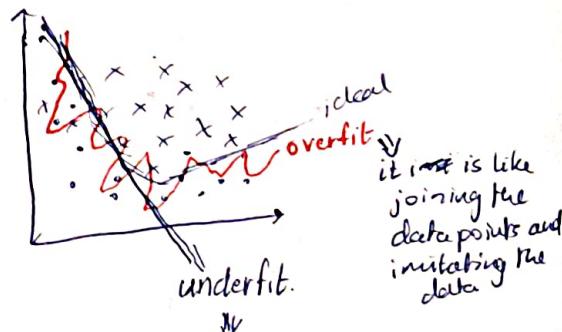
Using ~~optimisation techniques~~ different training sets of data for multiple models that are identical

- ideally, assuming that the error in prediction by each of the models is independent of each other, error = $\frac{1}{n} \times \text{error of each model}$

complexity of model depends on multiple factors

- no. of layers
- no. of neurons
- no. of epochs of training

too much complexity is as bad as too less complexity



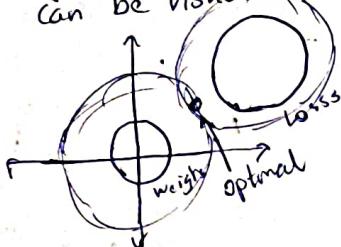
Having 2 consecutive dropouts before a layer simply means that the probability of dropout is increased for the next layer.

Geometric data augmentation causes issue in YOLO models as the coordinates of bounding boxes are definitely going to change

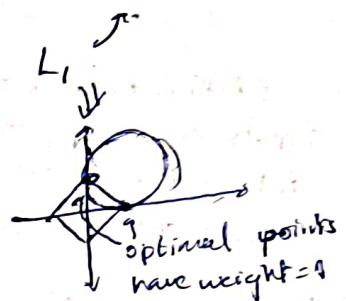
Regularization

- L_1 regularization
- L_2 regularization

the equation of L_2 regularization can be visualized as



- Instead of minimizing error, we minimize True error
- In order to make sure that the weights don't shoot up.
 - When weights shoot up, it can be said that the model is not able to find clear patterns within the data due to which the weights constantly change.
 - Penalizing the weights → keeping the weights in check.
- True error can also be written as $T.E = \text{Error} + \alpha \|w\|_1$
- L_1 is used for dimensionality reduction and not for overfitting handling. L_2 is for overfitting.



Adding Noise

we assume that the noise follows a Gaussian distribution

- E_i represents the noise random variable
- $\tilde{x}_i = x_i + E_i$ → inputs with added noise.
- $\tilde{y}_i = \sum w_i \tilde{x}_i$ → $E(\tilde{y}_i - y_i) = E((\tilde{y}_i - y) - \sum w_i E_i)^2$
expected $= E(y - y)^2 + (\sum w_i E_i)^2 / \sum w_i$
- same as L_2 regularization \Rightarrow Error + $\alpha \|w\|^2$

Training a YOLO model -

1. install YOLOv... using pip
2. Have a dataset that is compatible with YOLO
3. Use command line of YOLO to train model.

because covariance between E_i and $y - y$ is 0.

Batch normalization

→ Creating batches of data and applying normalisation of each.

suppose we take a batch of data, for each attribute within a row, we calculate the mean and standard standard deviation.

$$\Rightarrow \tilde{A}_i = \frac{A_i + \mu}{\sigma} \xrightarrow{\text{scale shift}} (A_i - \bar{A}_i + \beta) \xrightarrow{\text{within each batch}}$$

↓ moving avg. across batches

$$\mu_{\text{mov}} = \gamma \mu_{\text{mov}} + (1 - \gamma) \mu_i$$

$$\sigma_{\text{mov}} = \gamma \sigma_{\text{mov}} + (1 - \gamma) \sigma_i$$

Convolutional Neural networks are designed to work with grid-like structures.

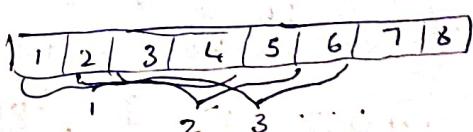
1-D CNNs can be used to work on time-series data.

What is convolution?

suppose we have a 1-D array of 8 elements, and we have a function that accepts 4 inputs and gives 1 output.

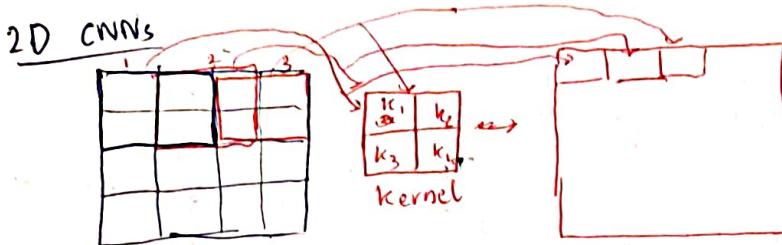
convolution is taking 4 inputs at a time

- uses concept of sliding window



in a convolutional neural network, convolution of input and filter is used to compute output.

x $w \rightarrow$ w as kernel (synonymous with weights in normal neural networks)

$$S_t = (x * w)_t$$


CNN is not much different from ANN, basically the inputs are aggregated by convoluting it with the kernel and summing it all up to predict a value for that given cell in the matrix.

Basic example of convolution \Rightarrow image editing (blurring, brightness, etc.)

$$\begin{bmatrix} \cdot & \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot & \dots \end{bmatrix}_{\text{2D matrix of multiple values}} * \begin{bmatrix} k_1 & k_2 & k_3 \\ k_4 & k_5 & k_6 \\ k_2 & k_8 & k_9 \end{bmatrix}_{\text{FxF Kernel}} = \begin{bmatrix} \cdot & \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot & \dots \end{bmatrix}_{\text{output 2D convoluted matrix}}$$

$m \times n$ (when there is padding at ends)

Types of convolution

1. No-zero padding / valid convolution \Rightarrow not padding with zeroes on edges of matrix to help fit the kernel function such that the cells on edges are placed aligned with the current kernel.

2. Just enough zero Padding / Same convolution

3. Full convolution

4. Stride (S) \Rightarrow when using the sliding window, instead of shifting sliding window by 1, we shift it by the value of S , a.k.a. stride. Doing this reduces the dimensions of the output as well.

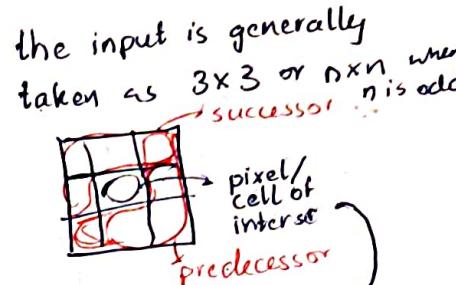
$$W_2 = \frac{W_1 - F + 2P}{S} + 1 \quad H_2 = \frac{H_1 - F + 2P}{S} + 1$$

Done with the assumption that images have high definition. By using a stride of the image has high definition, the next assumption is that the neighbouring pixels share same/similar information so we can ignore them. Resulting image can be called as a downsampled image. The process of reducing dimensions of the image is called downsampling.

5. Depth (K) - For an image, we can have multiple dimensions (R, G, B). There will be 1 2D matrix kernel for each dimension of the image. The number of such kernels is basically the depth of the kernel.

\Rightarrow Input $\Rightarrow 227 \times 227 \times 3$ Kernel $\Rightarrow 11 \times 11 \times 3$ Filters = 96 Stride = 4 Depth = 96

$$\frac{227}{4} = 56 \quad \frac{227 - 11 + 2}{4} + 1 = 55 \quad \Rightarrow 55 \times 55 \times 96 \text{ output.}$$



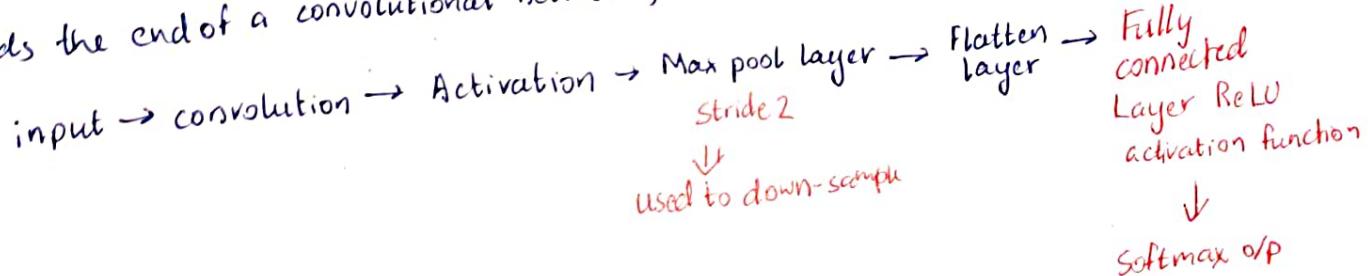
A value is set for this pixel based on surrounding cell values, as well as other current and the ~~current~~ kernel.

DL
An ANN has Dense connections (every neuron is connected to every other neuron)
But a CNN has sparse connections (not every neuron is connected to every other neuron)

ANNs will also work on Images however, the dense connections of the neurons because the ANN to expect similar input every time. However, the image that is passed into the ANN may be slightly shifted/scaled/rotated which confuses the ANN.
Therefore ANN will work on images of specified formats only, which make them unsuitable to work with images.

CNN does parameter sharing

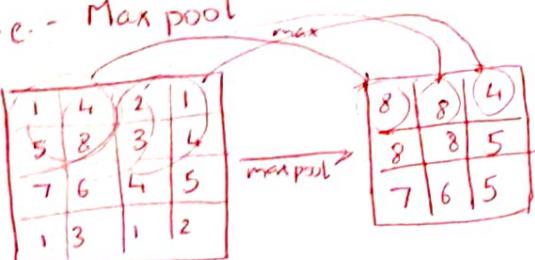
Towards the end of a convolutional network, the connections are dense.



Pooling

- helps to fetch statistics of data in the neighbourhood

i.e. - Max pool



↳ Helps to reduce computations while minimizing loss of data.

Max pooling helps to make the model scale-invariant.

Usually, pooling is used along with striding so that there is down sampling

LeNet (an example CNN)

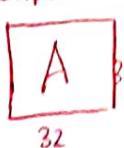
↳ Developed to extract postal codes

↳ Input $\Rightarrow 32 \times 32$, black and white only.

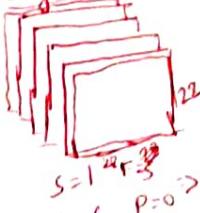
$k \Rightarrow$ depth of output, not Kernel.

no. of parameters \Rightarrow dimensions of kernel $\times k +$ bias

Input



convolution layer 1



pooling layer (14×14)

S=2 F=2
K=6 P=0
params=0
+
because no parameters are being fetched here

Convolution layer 2

S=1 F=5
K=16 P=0
params=2400
~~10x10~~ $\times 6 \times 16$

pooling layer 2 (5×5)

S=2 F=2
K=16 P=0
params=0

Fully connected layer (120 neurons)
params: 48120

Fully connected layer 84 neurons

Fully connected 10 neurons \rightarrow Output

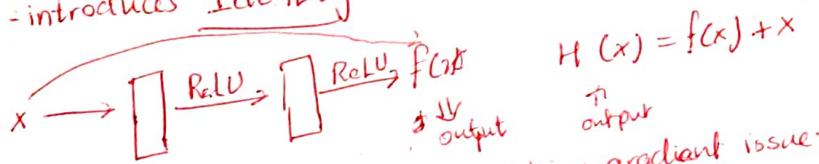
ILSVRC (ImageNet Large Scale Visual Recognition Challenge) (ImageNet) \rightarrow 1,281,167 images
 ↳ when a model is proposed for object classification, it needs to be trained on this dataset and prove its accuracy in classification.

VGG Net

- aims to minimize no of parameters
- same kernel size throughout the network — principle \Rightarrow Applying 2 3×3 kernels consecutively produces same size output as a 5×5 kernels
- ↳ By reducing parameters, the computation within each layer can be reduced.

Residual Networks (ResNet)

- introduces Identity connections.



→ Helps to overcome vanishing gradient issue.

[we are taking residue from previous layer]

- ↳ The input is passed along with output of previous layers into the upper layers of the network.

State Space models
Region Mambas

RCNN \Rightarrow Region based CNN \Rightarrow Detection + classification
YOLO \Rightarrow

- ↳ 2 stages
 - 1. Selecting bounding boxes
 - 2. Classification of object in bounding box. (Using CNN)

YOLO is a single stage detector; i.e. - it takes the object bounding box and search for class of object, only during training stages.



Neck, back bone of model.

↑
Feature extraction

Recurrent Neural Network, Many to many \Rightarrow language 1 to language 2 translation

Many to one \Rightarrow sentiment analysis

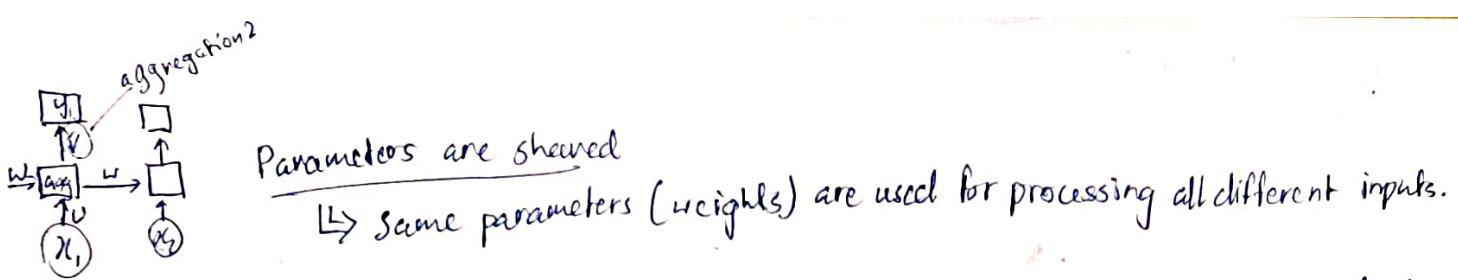
One to many \Rightarrow generation of data.

previous state
current state input

$$s^{(t)} = f(s^{(t-1)}, \theta)$$

current state
In RNN, the output is not passed in as input again, the state is passed again

to the next cycle of input



Aggregation (of inputs and history) → Activation → Aggregation again → Activation → output (softmax)

Training the network

$$o = \sigma(C + Vh_t) \text{ applied to each element within the aggregated value}$$

$$h_t = \tanh(w_{hi}x + b)$$

backpropagation

$$\nabla L = \frac{\partial L}{\partial v}$$

$$\frac{\partial L}{\partial v} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o} \cdot \frac{\partial o}{\partial v}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o} \cdot \frac{\partial o}{\partial h_t} \cdot \frac{\partial h_t}{\partial w}$$

$$h_t = \tanh(b + wh_{t-1} + ux_i)$$

$$\frac{\partial L}{\partial w} = \sum_{r=1}^t \frac{\partial h_r}{\partial h_r} \cdot \frac{\partial h_r}{\partial w}$$

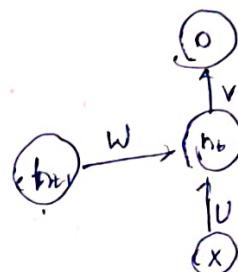
Types of RNN applications

1. One to one
2. One to many → one example input to give a sequence of outputs.
3. Many to one → a sequence of inputs is given and a single output is given.
4. Many to many → sequence of inputs → sequence of outputs.

make a vector/stack of

requires $h_{t-1}, h_{t-2}, \dots, h_0$
↳ back propagation through time

iterating using history and weights.
to update the weight w , we need to consider all previous history.



Disadvantages

1. Vanishing Gradients (in case the weights are too small &/negligible)
2. Exploding Gradient (in case the weights are too big)

Suppose x_t → input,
 \hat{y}_t → predicted output
is a matrix
another matrix

$$s_t = w \cdot a_{t-1} + u x_t + b$$

$$a_t = \tanh(s_t)$$

$$o_t = V \cdot a_t + c$$

$$\hat{y}_t = \text{softmax}(o_t)$$

prediction

$$s_t = [0.5] + [0] - [1]$$

$$= [-0.5]$$

$$a_t = [\tanh(-0.5)] = [-0.008]$$

$$V_1 = 0.3$$

$$o_t = [0.3 \cdot (-0.008)] + [0.1] = [0.0976]$$

example input and output	
in	out
1	0.5
2	0.1
3	0.5
4	0

take $a_0 = [0.5]$ (1×1 matrix)

$$w = [1]$$

$$x_t = [0]$$

$$b = -1$$

$$V = [0.3]$$

LSTM → extension of RNN where the memory that is important in long term and short term is used in current pred.

$$\hat{y} = \sigma(0.976) \\ = 0.9243$$

$$\text{loss} = 0.5 - 0.9243 \\ = -0.0243$$

LSTM

based on input, the memory of the previous state is retained.

input, output, forget gates are

taking example of NLP, (sentiment analysis)

we would

1. Forget stopwords
2. Selectively read information having sentiment words
3. Selectively write

$$S_{t-1} \xrightarrow{\substack{\text{Book} + \\ 60 \text{ slides}}} O_t = h_{t-1} \xrightarrow{\substack{\text{decides what} \\ \text{to remember} \\ (\text{a matrix})}} \rightarrow \text{process is called} \\ \text{selective write}$$

$$O_{t-1} = \sigma(w_0 h_{t-1} + u_0 x_{t-1} + b_0)$$

called output gate

$$h_{t-1} \xrightarrow{\substack{\text{Selective} \\ \text{Read} \\ \text{Written}}} W \xrightarrow{\text{f}_t} \overline{S_t} \times i_t + S_{t-1} \xrightarrow{\substack{\text{Temporary state} \\ \text{Selective} \\ \text{forget}}} S_t$$

selective read

selective forget

selectively read and forget about from previous state.

Suppose we want to caption a video, we wake captions into RNN.

RNN (CNN)

This model of an LSTM is very computationally expensive.

Hence introduction of Gated Recurrent Units (GRU)

↳ in GRU, the forget gate is basically taken as 1 - it

↳ in gated recurrent units, the forget gate is basically taken as 1 - it

↳ GRU is optimized on computation hence may perform worse than LSTM but not better than LSTM.

GRU may or may not train faster than LSTM

Applications of Encoder-decoder
Transliteration
Text summarization

Bi-LSTM \Rightarrow bidirectional LSTM

Encoders and Decoders

Encoders and decoders can be anything, like RNN, LSTM, etc.

For example, suppose an image is to be converted to text,

- we first have a CNN to be encoder \rightarrow Extract features of the image-
- then we have another model, i.e. RNN to generate text from the input feature.

Textual entailment \Rightarrow based on text 1, text 2 is constructed.

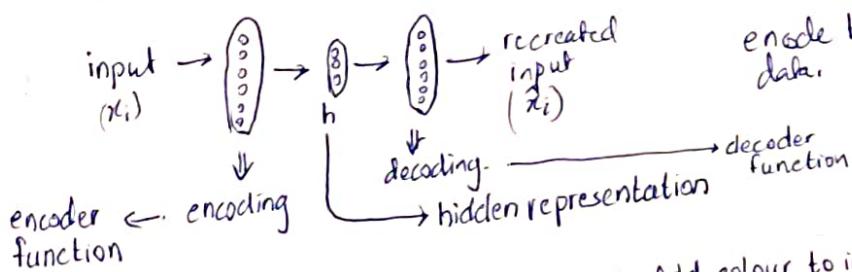
2 RNNs are used, output of 1 RNN is used as input to next RNN

Encoder, decoder is different from U-net and Transformer

Auto encoders

→ Used for PCA

→ A neural network that takes x as input, encodes, decodes and gives out x as output
encode to create a latent representation of the data.
→ compressed.



Auto encoders can be used to

- 1. Add colour to images
- 2. Add clarity to images. (add noisy image as input, loss is to be calculated as loss of data b/w originating and prod. img)

Undercomplete Autoencoder - hidden layer size is smaller than input layer size.

Overcomplete Autoencoder - hidden layer size is larger than input layer size

→ should have non-linear function in order to ensure that the neural network is learning instead of acting as identity functions.

Reconstruction error \Rightarrow error b/w input and recreated input (Output)
→ Used when input data is very complex and some extra data would help to identify dependancies in the data.

Regularized Autoencoder (L2 regularization)

Train any architecture of autoencoder
freely change the dimension of the hidden representation

Sparse Autoencoder (for overcomplete autoencoder)

similar to dropouts is done.
→ Neurons are dropped in hidden layer so that only some number of neurons are present in HL and input layer.

Contractive Autoencoders

Also uses regularization

A matrix is created such storing the change in output of each neuron with

respect to each feature

computing Jacobian coefficient for each feature

used to remove unimportant features.

$$f \Rightarrow \begin{bmatrix} \cdot & \cdot & \cdot & \dots & \cdot \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \cdot & \dots & \cdot \end{bmatrix}$$

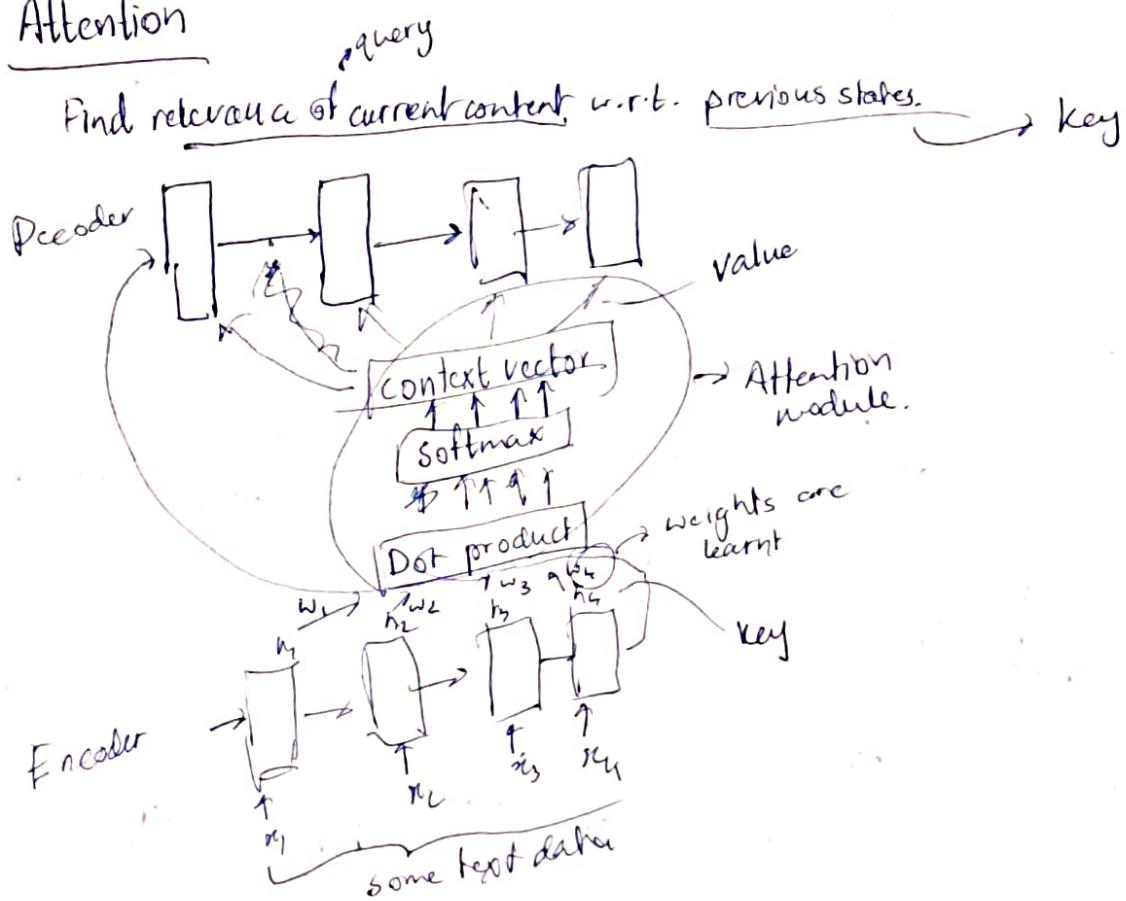
$$f_1 \Rightarrow \begin{bmatrix} \cdot & \cdot & \cdot & \dots & \cdot \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \cdot & \dots & \cdot \end{bmatrix}$$

$$f_2 \Rightarrow \begin{bmatrix} \cdot & \cdot & \cdot & \dots & \cdot \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \cdot & \dots & \cdot \end{bmatrix}$$

Denoising Autoencoders

Adding Gaussian noise to the input data while having predicted output to be image without noise. Hence using the ~~data~~ model to learn to remove the noise.

Attention



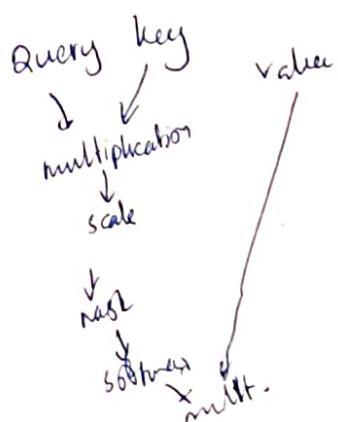
Attention and Transformers were originally made for NLP

$$\underline{(Q \cdot K)V}$$

↓
softmax

Self-attention

↳ find relevance of part of current context w.r.t. another part in the same sentence



Transformers are basically use of multiple self-attention module

The model may have to train for more epochs after attention is added