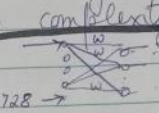


# CNN - CONVOLUTIONAL NEURAL NETWORK

## Introduction

→ We studied fully connected Multilayer Perceptron MLP Neural Network, using sigmoid function to activate a neuron, and back propagation to adjust weights. They are not suitable for images because:-

- ① One pixel going to one perception at input stage means for  $1024 \times 1024 \times 3$  (RGB) of them at the input and then fully connected at hidden layer.  
→ Excessive Computation complexity. Too many weights to be trained! 

- ② MLP architecture does not exploit spatial correlation which is typical in images.

- ③ It cannot make sense of shifts in positions - Translations of an object cannot be handled.

→ CNNs use filters - 2D slices of weights which are much fewer than number of pixels. These filters are applied all over the image. Their weights are "learned" by the CNN through supervised learning with back propagation. Thus, CNNs are most suitable for image analysis.

(i) CNNs are the only ML technique that automatically learn features in an image  $1024 \times 1024$ .

(ii) The filter weights are shared by different portions of the image. Thus, there are fewer weights. Of course there are several filters that are learned for any given image.  $(25 \times 25)$

(iii) CNNs are translation-invariant in learning features.

(iv) CNNs are deep rather than wide as MLP.

(v) CNNs typically use ReLU activation function instead of sigmoid to avoid vanishing gradient problem.

Topic: Convolutional Neural Network CNN ①  
Date: \_\_\_\_\_ Page No: \_\_\_\_\_

Convolutional Neural Networks (CNN) are deep neural networks with a series of NN layers:

- 1/P → convolutional, nonlinear, pooling... fully connected → 2/P

Thus, the input is fed into the first convolutional layer and after a series of transformations finally pass through a fully connected layer for classification at the output.

CNNs are most used for image classification like the "visual cortex" of human brain where different neuronal cells fire in the presence of different features of the visual field (e.g. some fire in the presence of edges), so also in a CNN different parts of an input image trigger different sets of neurons in the conv. layer. Thus, each convolutional neuron has its own "receptive field" in the image. Note that this is in contrast with a Multi Layer Perceptron Feed Forward NN, which is fully connected from input to hidden layers.

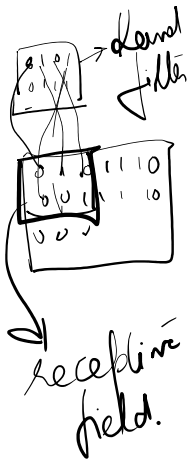
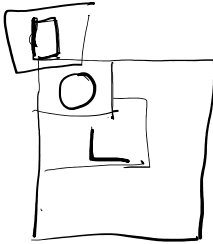
Discrete convolution is  $\sum_{m=0}^{i-1} f[m] g[n-m] \rightarrow$  Shift, multiply, add

Despite the name "Convolution", the mathematical function that is performed is a "Dot Product" between an (array of numbers) and the receptive field of an image.

- The array is called Filter or Kernel
- The numbers in this array are the weights

So, assume the image is of size 32 pixels  
 $\times 32 \text{ pixels} \times 3$  (for each component of RGB)  
 with each pixel an 8-bit integer.

NEW DAY

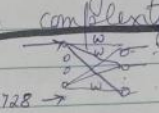




# CNN - CONVOLUTIONAL NEURAL NETWORK

## Introduction

→ We studied fully connected Multilayer Perceptron MLP Neural Network, using sigmoid function to activate a neuron, and back propagation to adjust weights. They are not suitable for images because:-

- ① One pixel going to one perception at input stage means for  $1024 \times 1024 \times 3$  (RGB) of them at the input and then fully connected at hidden layer.  
→ Excessive Computation complexity. Too many weights to be trained! 

- ② MLP architecture does not exploit spatial correlation which is typical in images.

- ③ It cannot make sense of shifts in positions - Translations of an object cannot be handled.

→ CNNs use filters - 2D slices of weights which are much fewer than number of pixels. These filters are applied all over the image. Their weights are "learned" by the CNN through supervised learning with back propagation. Thus, CNNs are most suitable for image analysis.

(i) CNNs are the only ML technique that automatically learn features in an image  $1024 \times 1024$ .

(ii) The filter weights are shared by different portions of the image. Thus, there are fewer weights. Of course there are several filters that are learned for any given image.  $(25 \times 25)$

(iii) CNNs are translation-invariant in learning features.

(iv) CNNs are deep rather than wide as in MLP.

(v) CNNs typically use ReLU activation function instead of sigmoid to avoid vanishing gradient problem.

Topic: .....

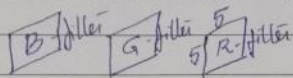
CNN

Date: .....

Page No: .....

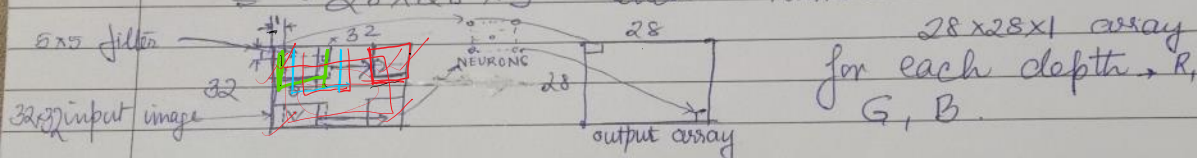
(2)

The kernel is an array of  $(5 \times 5 \times 3)$  weights.



The filter slides through the image from top left to bottom right, allowing each convolutional neuron to perform a dot product in its own receptive field of  $5 \times 5 \times 3$  pixels with the filter. Thus each convolutional neuron performs 75 multiplications & adds them up to produce a single no. A bias is also added. Assume the filter slides by one pixel, right & down.

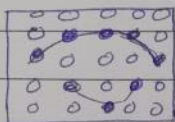
The kernel, like a flashlight moves over the image to produce  $(32 - 5 + 1) \times (32 - 5 + 1) \times 3$  numbers.



For a single filter, we get an output of  $28 \times 28 \times 1$ . With  $n$  filters, we can get  $28 \times 28 \times n$  such matrices, as many neurons as 784.

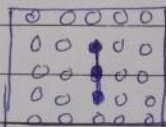
The filters represent different kinds of possible features in the image such as

→ curves



black dots have high value say 45  
white dots are zero

→ straight edges.



→ colour :- Say red region for cheeks, in R-FILTER

NEW DAY™



Topic: .....

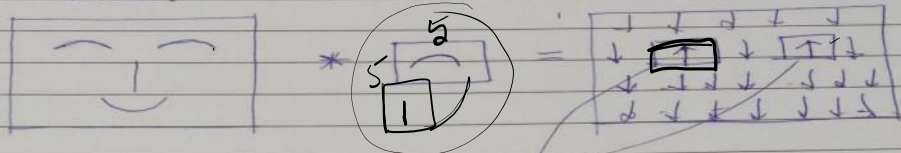
CNN

Date: .....

Page No.: .....

3

The imp. point is  $\rightarrow$  if the input image contains that feature in the receptive field for a conv. neuron, then it will output a large value.



These two receptive fields will give high value  $\uparrow$  for above filter, others will give low values  $\downarrow$ .

Thus, the dot product discovers the "correlation" between filter and image to learn certain aspects of the image.

When there are several filters, then more features can be learnt.

The significance of a CNN is that we do not have to pre-design or handcraft these filters, rather the CNN learns the weights of the filters by itself, through supervised training with backpropagation.

The CNN is a deep architecture with multiple convolution layers. As we go further in the series from input to output, each conv. layers gathers as input the features extracted from previous layer and increases its "receptive field" size. The aim is to

detect higher forms of features. For example

input conv layer detects → next conv layer detects → next conv layer → ... → final conv layer

curves & edges → combination of basic features - semi-circle, curve + edge, hor + vert edge, red areas, patterns → higher combination of features, facial parts → a proper face with say brown skin, black eyes, dark hair

Thus, a deep NN begins with low level features like curves and edges and gradually increases the receptive fields to extract high level features such as face, paws, wings, beak, etc. blue eyes, red portions, etc. patterns of colours & shades etc.

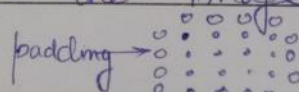
### PADDING

It may be noticed that the image pixels at the edges are not touched by the kernel as frequently as the middle. This means they are not searched thoroughly for features. To enhance their chances, "padding" is followed. This overcomes two problems

- 1) losing corner information
- 2) Shrinking output of convolution.

$n \times n$  image to  $(n - w + 1) \times (n - w + 1)$ , where  $n$ :- image dimension,  $w$ :- filter dimension.

When there is padding, additional zeros are added on the border of the image.





Topic: .....

CNN.

Date: .....

Page No. : .....

For a padding width of size  $p$ , the output image is  $(n + 2p - w + 1) \times (n + 2p - w + 1)$  ✓

Valid padding means no padding while 'Same' padding is when the pad restores size of original image:  $n + 2p - w + 1 = n$  Thus  $p = \frac{w-1}{2}$

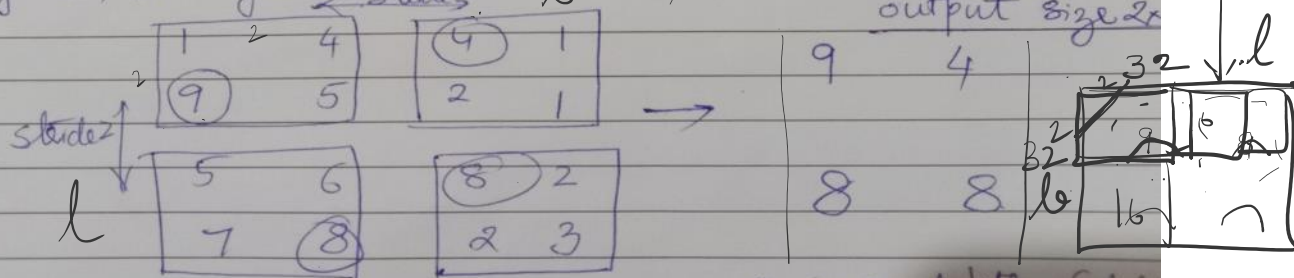
STRIDE:- So far filters were moving over image with a stride = 1. For a stride of  $s$ :-  
Output is  $\left[ \frac{n+2p-w}{s} + 1 \right] \times \left[ \frac{n+2p-w}{s} + 1 \right]$ , where

$s$  is both horizontal & vertical stride. Stride reduces the size of an image, taking advantage of spatial relatedness.

Conv

POOLING Layer.

Once a feature is captured by convolution layer, it is not necessary to use the entire picture but certain highlights only. Pooling layer uses special filters to capture the average or maximum of the relevant numbers. Thus for a  $4 \times 4$  block within an image, using max pooling with a  $2 \times 2$  filter, stride of 2, we get:-



$n_c$  is no. of filters (depth)

In general, if  $l \times b$  is image block,  $w$  is filter width & stride then output =  $\frac{l-w+1}{s} \times \frac{b-w+1}{s} \times n_c$

31379131

Topic: .....

CNN

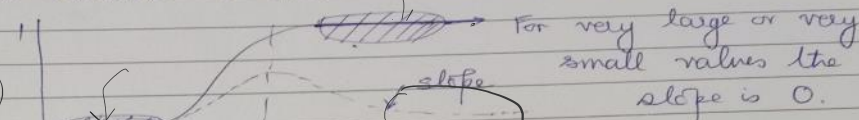
Date: .....

Page No. 6

6

ReLU { Rectified Linear Unit Activation Function }

The sigmoid activation function has a narrow transition between 0 and 1. Hence it is called a "soft switch".

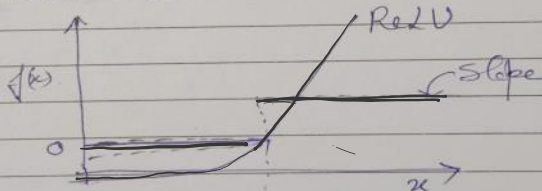
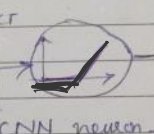


Since output error is propagated gradually from output to input by differentiating the cost function (squared error) w.r.t. weights and their activation function  $f'(z)$  at each layer  $l$ , gradually the terms become zero for highly activated neurons. Thus the weights are hardly adjusted in consecutive iterations.

To avoid this situation, one must use a non linear function with +ve and constant slope.

DOT-PRODUCT

kernel  
\*  
Receptive field  
= a CNN neuron



$$f(x) = \max(0, x)$$

ReLU  
 $g(x)$

Now the slope is large enough to remain significant for learning, so that weights are adjusted according to Eq.  $W_{l+1} = W_l - \eta \left( \frac{dC}{dW_l} \right)$ , at each layer  $l$ ,  $\eta$  being the user defined learning rate.

$\max(E_3)$   
 $E \ll 1$

FULLY CONNECTED LAYER:

$\max(E_3)$

The final layer of a CNN is fully connected M&P. It works on a greatly filtered and

$\alpha(e^2 - 1)$

$E \ll 1$

NEW DAY





Topic: CNN

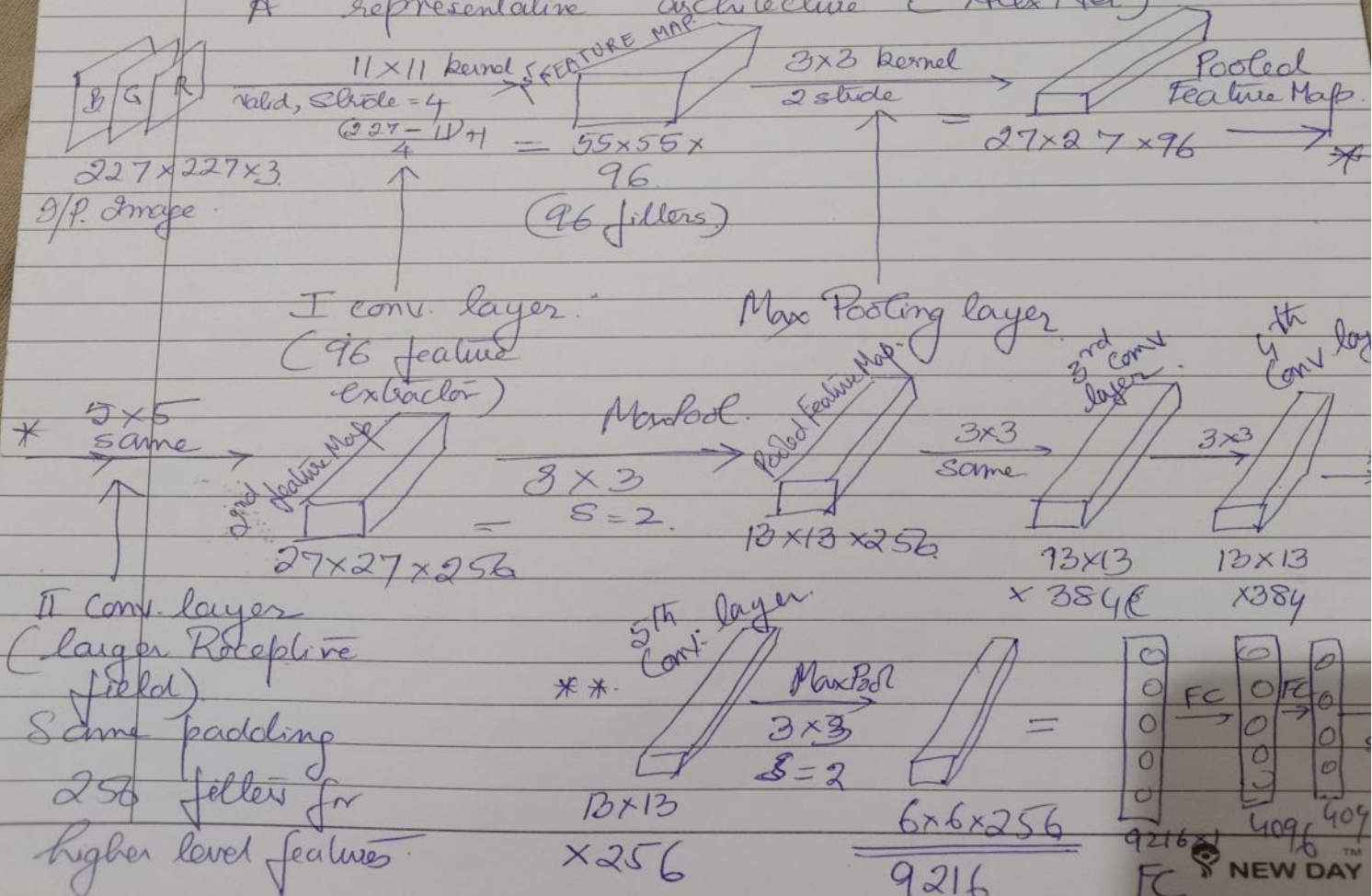
Date:

Page No.:

7

reduced form of the original image due to a series of convolutions & pooling. Therefore this is the "flat" part of the CNN, where results are classified just like any MLP. There can be one or more FC layers. The activation function used in the FC layer may be ReLU if aggregating image information or Softmax if performing classification.

A representative architecture (Alex Net)



Topic: .....

CNN

Date: .....

Page No.: .....

8

Such elaborate CNNs have millions of weight parameters to be trained besides a host of other <sup>hyper</sup>parameters like kernel width, height, no. of kernels, stride, padding, regularization to contain overfitting, learning rate.

→ Some Regularization methods :-

Too many weights, all trained thoroughly to some training data, becomes "rigid" → meaning they cannot generalize to even slightly different test data thus giving low accuracy.

To prevent this, typically following are done :-

1) Early stopping :- Training and validation are done simultaneously using different parts of the available data set. Initially both improve accuracy of prediction/classification. After a while, training accuracy improves but validation accuracy drops. That when further training is stopped. Next, testing follows, with much reduced overfitting.

2) Dropout :- Train the network with reduced nodes. Each node is either kept with prob.  $p$  or kept out with prob.  $(1-p)$ . Input nodes are typically kept in so that input info is not lost. The FC output layer has most of the weight, so they are dropped out with higher prob. Intermediate nodes are kept with  $\sim 0.5$  prob. Thus the network is trained on fewer nodes during each training iteration, with "out" nodes re-inserted before next iteration. This relaxes overfitting.

During testing, each node's output is simply weighted with prob.  $p$  (full network is tested). The network gives more <sub>node</sub> robust results, correctly predicting/classifying varied test inputs.

3) Data perturbations :- To make the training data more varied, it is perturbed by artificial & deliberate changes such as cropping, erasing, rotating objects etc, while retaining same label as the original.



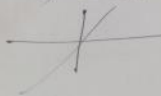
# Different variants for ReLU.

1) Linear

$$f(z) = z \times m \quad f'(z) = m$$

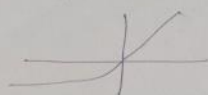
2) ELU

+ve :- large range of activations, not just binary.



-ve :- slope is constant and not dep. on change in  $x$ .

2) ELU :- Exponential Linear Unit



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

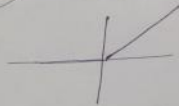
$$f'(z) = \begin{cases} 1 & z > 0 \\ \alpha e^z & z < 0 \end{cases}$$

+ve :- Unlike ReLU, can produce -ve outputs.

-ve :- ~~Exponential~~ Exploding activation for +ve  $z$ .

3) ReLU

$$f(x) = \begin{cases} x & x > 0 \\ 0 & x \leq 0 \end{cases} \quad f'(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$



+ve :- avoids vanishing gradients prob

-ve :- can be used only in hidden layers

dead neurons for -ve region.

4) Leaky ReLU.



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

$$f'(z) = \begin{cases} 1 & z > 0 \\ \alpha & z < 0 \end{cases}$$

+ve :- fixed ~~dying~~ ReLU

-ve :- only ~~hidden~~ layers Not

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad \text{K no. of classes.} \quad (10)$$

Softmax activation

→ A vector of output is converted into a vector of probabilities in a well-calibrated manner.

→ Outputs are normalized so that their sum is zero.

→ Each prob. is membership degree in that output class.

$$n \text{ output classes } \left\{ \begin{array}{l} \text{class} \\ \text{Raw values} \\ o_1 \\ o_2 \\ o_3 \\ \vdots \\ o_m \end{array} \right\} \rightarrow \begin{array}{l} \text{prob} \\ \text{values} \\ e^{o_1} / \sum e^{o_i} \end{array}$$

→ Softmax is a soft version of argmax.

→ It is used because ~~Linear~~ <sup>linear</sup> & Sigmoid are inappropriate for multi-class classification.

(i) Linear (output = input) → no activation generates any value.

(ii) Sigmoid is OK for binary classification {binomial distribution}, as well as for non mutually exclusion multi-class classification. But not for mutually exclusive multi-class classification (multinomial prob. distribution).