

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

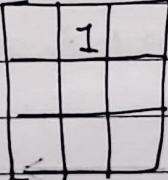
Page No.

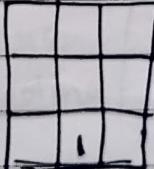
Date

Unit 2 :- CNN

- i) Alexnet
- ii) Resnet
- iii) LeNet
- iv) UGGNet
- v) RNN
- vi) Transfer learning

~~original matrix * filter → new grid
grid~~
feature map.

9 * loopy
pattern = 

6 * -11- = 

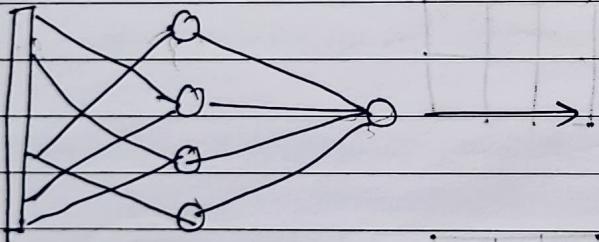
ReLU → makes model non-linear.

[S.D./19]
 Pooling layer → Reduce dimension of input image.

- ~~Advantages~~ Benefits → i) Reduce dimensions & comp.
 ii) Reduce overfitting
 iii) Model is tolerant towards variations & distortions

i/p image → Convolutional + ReLU → Pooling → Convolutional + ReLU → Pooling

Feature extraction



Classification

- Convolution →
- Connection sparsity reduces overfitting
 - conv + pooling gives location invariant feature detection
 - Parameter sharing

ReLU → • Introduces non-linearity.

- speeds up training; faster to compute.

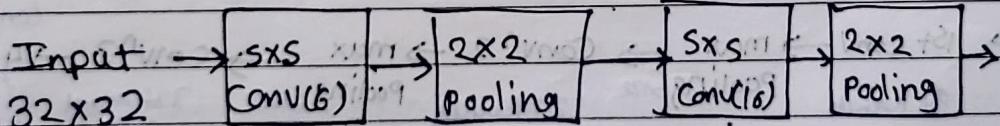
Pooling → • Reduces dimension & computation

- Reduces Overfitting

• Makes the model tolerant towards small distortion & variations

* AlexNet & LeNet

* LeNet



* Padding

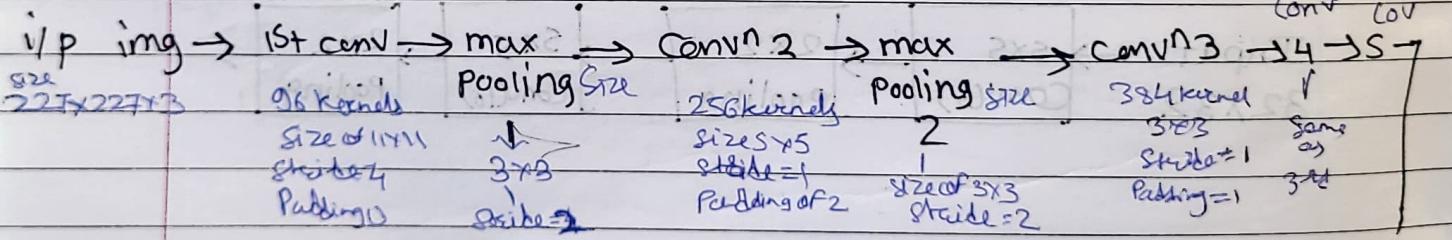
$$(n + 2P - f) + 1 \rightarrow \text{output size}$$

When we apply filter in info in corner is lost. To overcome this we use zero padding

$f \rightarrow$ filter size

* AlexNet

- winner of 2012 Imagenet competition
- consists of 8 layers
- 5 are convolution layer
- 3 fully connected layers
- I/p size $\rightarrow 227 \times 227 \times 3$
- 1st convolution layer \rightarrow 96 filters, 11×11
stride = 4



- max pooling $\rightarrow 3 \times 3$ with padding 2

101-299 X

~~form~~

$$\text{O/P size} = \left(\frac{n+2p-f}{s} \right) + 1$$

- ReLU is applied after each convolutional & fully connected layer except last
 - Dropout is applied before 1st & 2nd fully connected layers

Unit 4: * Transfer Learning

problem defining exactly what is and is not

11 - repeat loop example 1

2010-09-16 09:00:17 D01X8.P 201 - 45

Section 20(1) of the Act – moulding the horns of guidance
for the benefit of both the people and the environment
without regard to the existing and existing
right of landholders to continue to do so

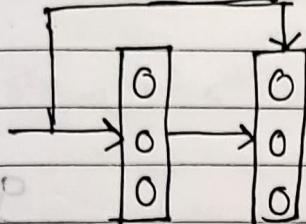
* Res-Net

- Deep networks are hard to train because of vanishing gradient problem
- On ImageNet dataset evaluate residual nets with depth of 152 layers - 8 deeper than VGGNet but still having lower complexity.
- An ensemble of these residual nets achieve 3.57% error on the ImageNet test set. This result won 1st place on ILSVRC 2015 classification task.
- Resnet 50 is a variant of Resnet model which has 48 conv. layers along with 1 maxpool & 1 average pool layer.
- It has 3.8×10^9 floating point operations

{ vanishing gradient problem - As more layers using certain activation functions are added to neural networks, the gradients of the loss function approaches zero, making the network hard to train }

- Residual networks use the concept of skip connection.

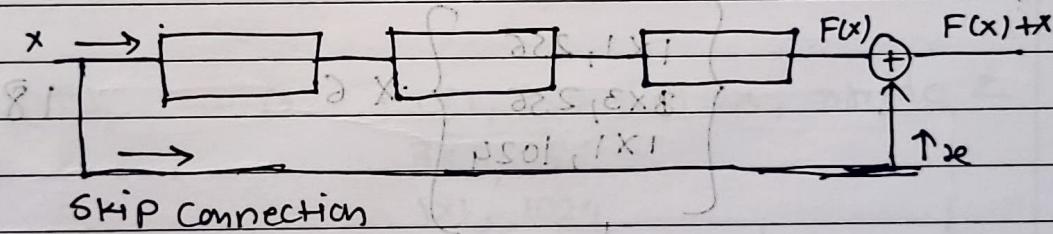
* Skip connection



original i/p is added to the o/p of convolutional block.

since, some of the connections are skipped, the value of the loss function gradients won't become zero.

* Residual block



$$y = F(x) + x$$

$F(x)$ → loss function.

{ make $F(x)=0$, so that we can make $y=x$ }

* Resnet - 50

$7 \times 7, 64, \text{stride} = 2$

$3 \times 3, \text{max-pool}, \text{stride} 2$

$$\left\{ \begin{array}{l} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right\} \times 3$$

$$\left\{ \begin{array}{l} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right\}$$

$$\left\{ \begin{array}{l} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right\} \times 6$$

$$\left\{ \begin{array}{l} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right\} \times 3$$

Average pool, $100 \times 1 \times 1$, softmax

50 layers

1st layer \rightarrow 7x7, 64 filters with stride 2

2nd layer \rightarrow maxpool layer 3x3 with stride 2

3rd, 4th & 5th \rightarrow 1x1, 64 filters

3x3, 64 filters

1x1, 256 filters

7th \rightarrow 1x1, 128 filters, stride 2

3x3, 128 filters

1x1, 512 filters

8th, 9th, 10th \rightarrow 1x1, 128

3x3, 128

1x1, 512

11th: \rightarrow 1x1, 256 filters, stride 2

3x3, 256

1x1, 1024

12th, 13th,

14th, 15th & 16th \rightarrow 1x1, 256

3x3, 256

1x1, 1024

17th: \rightarrow 1x1, 512 with stride 2

3x3, 512

1x1, 2048

18th & 19th \rightarrow 1x1, 512

3x3, 512

1x1, 2048

20th \rightarrow Average pool with strides.

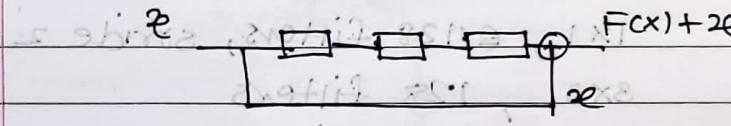
Subtraction error at, $x = \text{avg pool}$ for

Stride size 8x8 \rightarrow fully connected layer (fc), 1000

$28 \times 28 \times 32, 1 \times 1 \leftarrow 14 \times 14, 64$

$28 \times 28 \times 32, 1 \times 1$

* Identity block $28 \times 28, 1 \times 1$



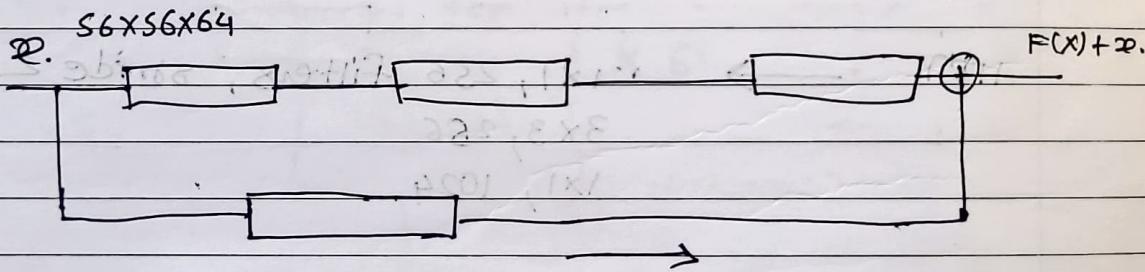
ip size = o/p size.

$28 \times 28, 1 \times 1 \leftarrow 14 \times 14, 64$

$28 \times 28, 1 \times 1$

* Convolutional block 1×1

$28 \times 28 \times 128$



ip size \neq o/p size

we can add x only if size of $F(x) = x$,
but when $x \neq F(x)$ we add a convolution layer in
shortcut path so that, $F(x) = x$

we use 1×1 convolution to make size of $F(x)$, is equal to x .

$$\left(\frac{n+2P-f}{s} \right) + 1$$

$$\left(\frac{86+2 \times 0 - 1}{2} \right) + 1 \rightarrow 28$$

* Hyperparameters

- Batch Normalization is used after each conv. layer.
- Xavier initialization with SGD+momentum is used. A basic term of 8×8 image is used.
- The learning rate is 0.1 & is divided by 10 . Validation error becomes constant after 20 epochs.
- more over, batch size is 256 & weight decay is $1e-5$.
- There is no dropout in ResNet.

* VGG 16

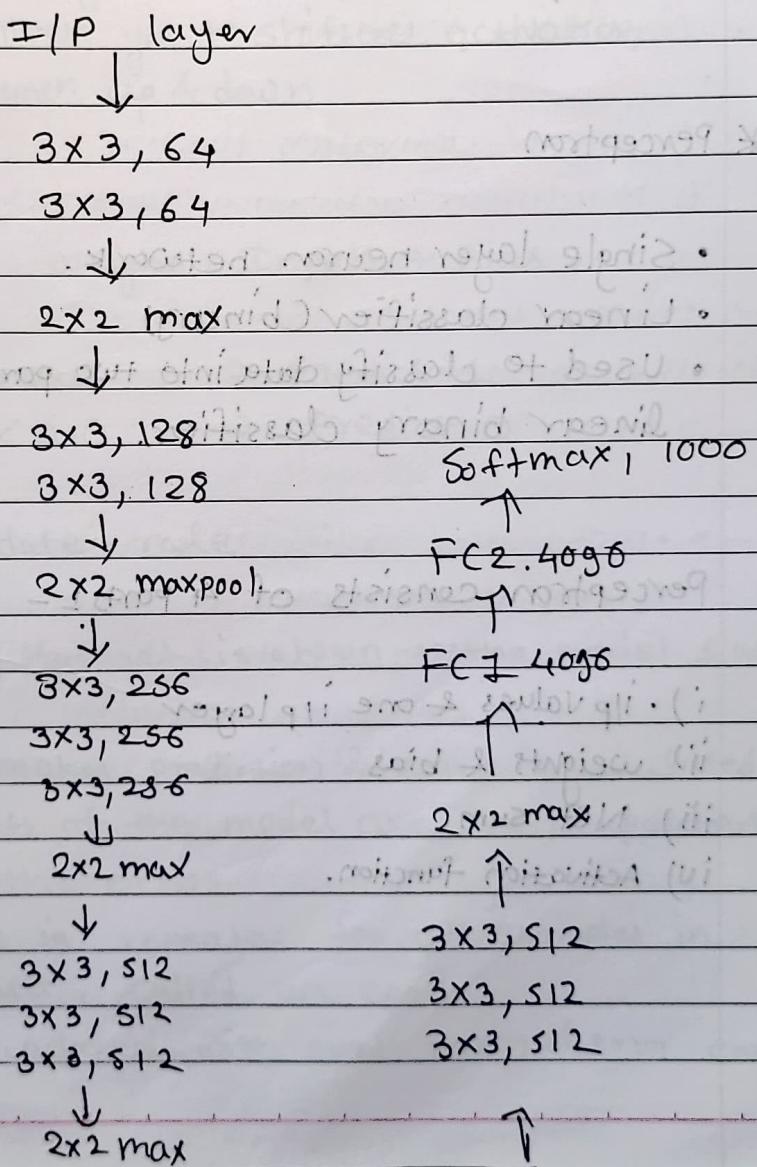
- Introduced by simonyan & zisserman in 2014.
- VGG Net was runner up of the Imagenet visual recognition (ILSVRC) classification in 2014.
- VGG stands for visual geometry group.
- VGG 16 has 16 conv. layers & VGG19 has 19 conv. layers.
- The VGG architecture is the basis of ground breaking object recognition models.
- VGGNet 16 achieves almost 92.7% top-5 test accuracy in ImageNet.
- It replaces large kernel size filters with several 3×3 kernel sized filters.
- VGGNet - 16 supports 1.16 layers & can classify images into 1000 object categories.

* Architecture

i) VGG 16

- 13 convolutional layers + 3 fully connected layers.
- Input size $\rightarrow 224 \times 224$
- VGG's conv. layers leverage the minimal receptive field. i.e. 3×3 , the smallest possible size that still captures up/down & left/right.

- There are also ~~not~~ 1×1 conv. layers acting as linear transformation i/p.
- This followed by ReLU. ReLU reduces training time.
- All hidden layers in VGG uses ReLU.
- VGG does not usually leverages Local Response Normalization (LRN) as it increases memory consumption & training time.
- Out of the 3 fully connected layers 2 have 4096 neurons & last one has 1000.



* VGGNet vs AlexNet

* VGGNet vs ResNet

* Unit 1 - Deep Learning basics

To generalize if a $m \times m$ image is convolved with $n \times n$ kernel. The o/p image size is $(m-n+1) \times (m-n+1)$

* Perceptron

- Single layer neuron network.
- Linear classifier (binary)
- Used to classify data into two parts, hence known as linear binary classifier.

Perceptron consists of 4 parts :-

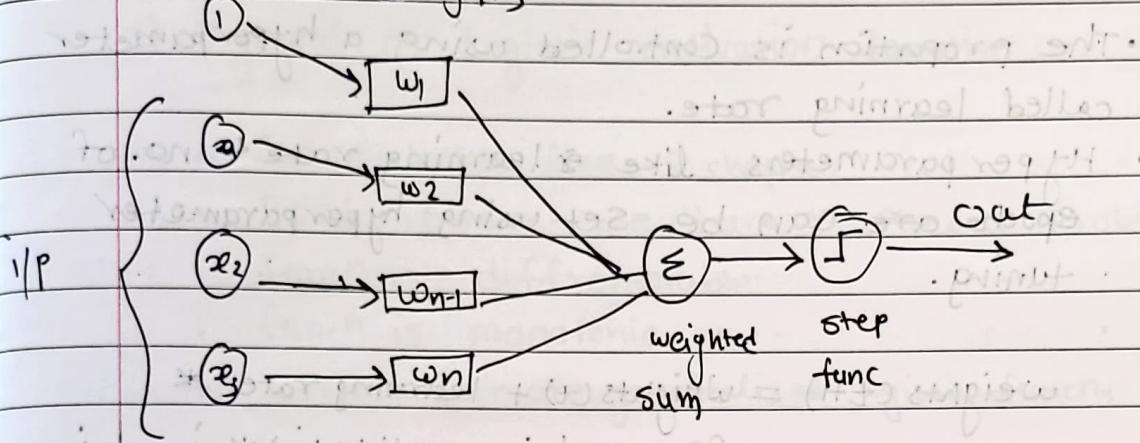
i) ip values & one ip layer

ii) weights & bias

iii) Net sum

iv) Activation function.

const. weights



weights \rightarrow strength of particular node

bias \rightarrow Allows you to shift activation function up & down

Activation \rightarrow weights \times Inputs + Bias

if activation $> 0.0 \rightarrow$ predict 1
if activation $< 0.0 \rightarrow$ predict 0

* Perceptron update rule

- The training dataset is shown to the model at a time.
- the model makes prediction & error is calculated.
- The weights of the model are then updated to reduce the errors.
- This process is repeated for all examples in training dataset, called an epoch.
- model is updated with small portion of error each batch.

- The proportion is controlled using a hyperparameter called learning rate.
- Hyperparameters like learning rate & no. of epochs can be set using hyperparameter tuning.

$$\text{weights}(t+1) = \text{weights}(t) + \text{learning rate} * (\text{expected}_i - \text{predicted}_i) * \text{input}_i$$

* Activation functions

- Linear activation function
- Non-linear activation function.

i) → linear function is a line or linear funcn. Therefore the opf of funcn. will not be confined between any ranges.

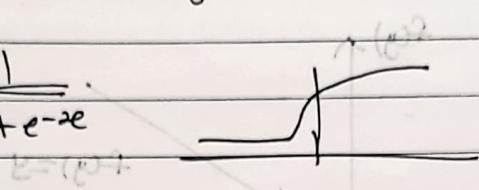
ii) → most used activation functions

- easy for the model to generalize or adapt with variety of data & differentiate between opf.
- The non-linear activation functions are mainly divided on the basis of range or curves

1) Sigmoid or Logistic activation function

- Curve looks like S shape
- Used where we have to predict probability
- func is differentiable
- func is monotonic
- softmax is more generalized logistic func.

$$\phi(z) = \frac{1}{1+e^{-z}}$$



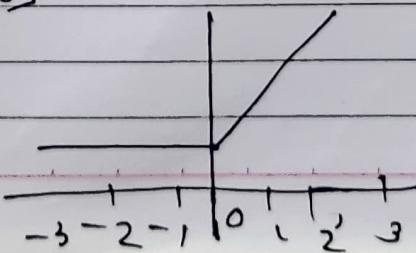
2) Tanh or hyperbolic tangent activation function.

- tanh is also like logistic sigmoid but better. The range of the tanh func is from (-1 to 1).
- The advantage is negative i/p will be mapped strongly to negative & zero i/p will be mapped near zero in tanh graph.
- The func is differentiable.
- The func

3) Relu

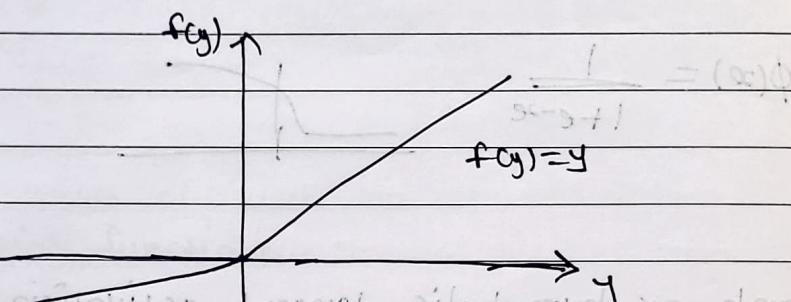
- Most used activation func
- Half rectified from bottom
- Range $\rightarrow [0 \text{ to } \infty]$

$$F(z) = \max(0, z)$$



4) Leaky ReLU

- The leak helps to increase the range of ReLU, usually, value of $\alpha \rightarrow 0.01$.
- when α is not 0.01 then it is randomized ReLU
- Range $\rightarrow (-\infty, \infty)$
- Monotonic in nature & their derivatives too



Leaky ReLU is a smooth version of ReLU. It is differentiable at all points except at zero. At zero, it is non-differentiable. It is used in neural networks because it is differentiable everywhere and has a constant negative slope for negative values of y.