# NPTEL Week 9 Live Sessions

## on Deep Learning (noc24_ee04)

A course offered by: Prof. Prabir Kumar Biswas, IIT Kharagpur

- **Quiz 7, Quiz 8 Solution**
- **VGG16, ResNet implementation**

By

**Arka Roy**
**NPTEL PMRF TA**

**Prime Minister's Research Fellow**
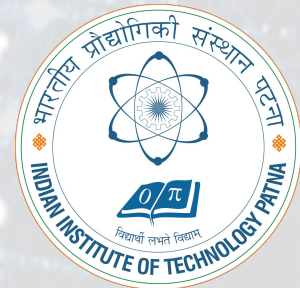**Department of Electrical Engineering, IIT Patna**
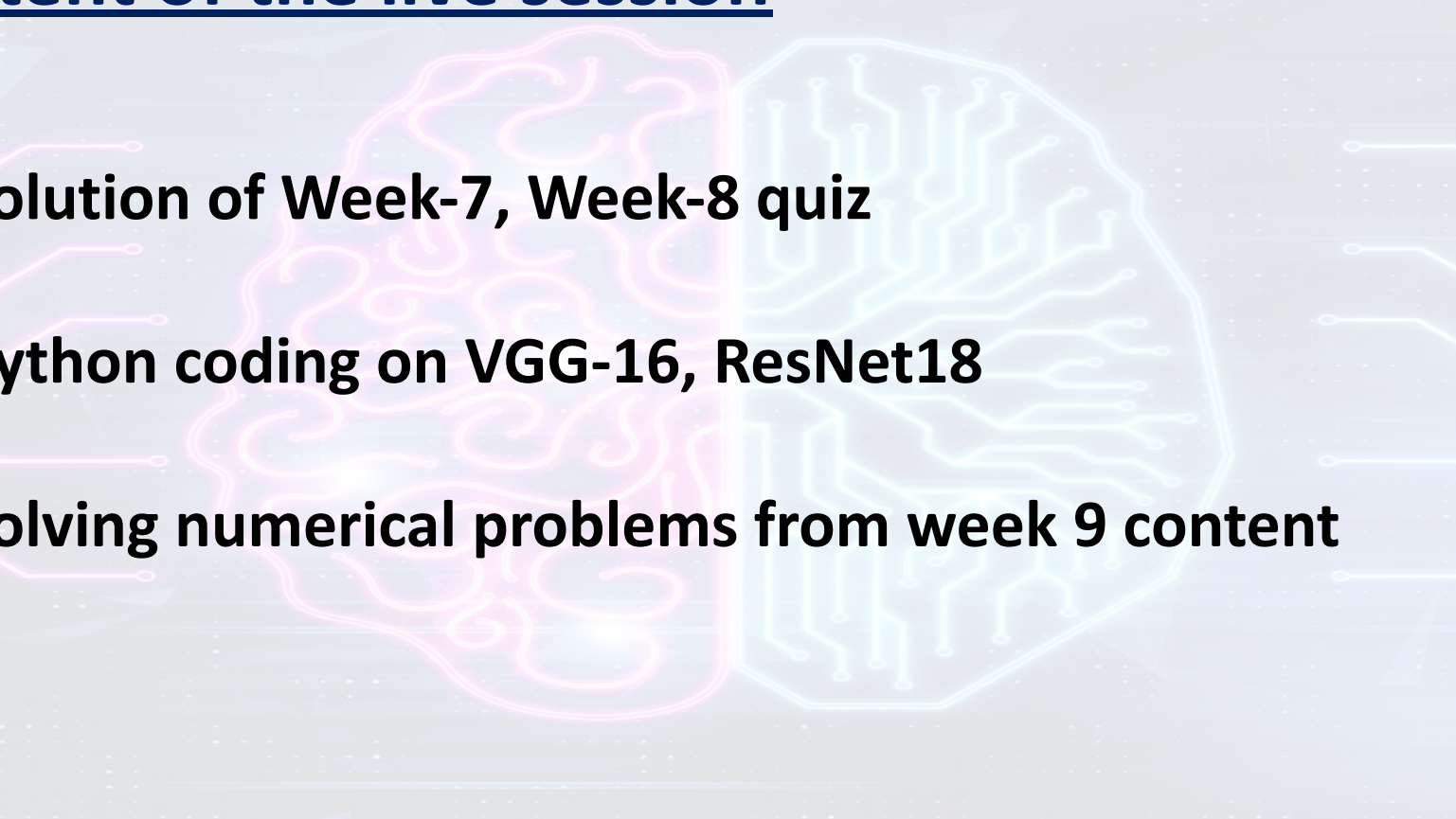**Web:** **https://sites.google.com/view/arka-roy/home**

PMRF
Prime Minister's Research Fellows
Ministry of Education
Government of India

# Content of the live session

1. Solution of Week-7, Week-8 quiz

2. Python coding on VGG-16, ResNet18

3. Solving numerical problems from week 9 content

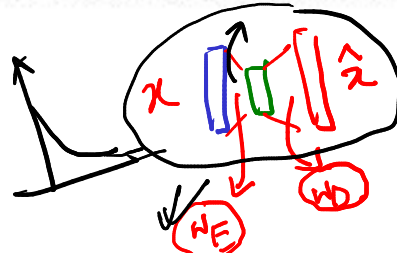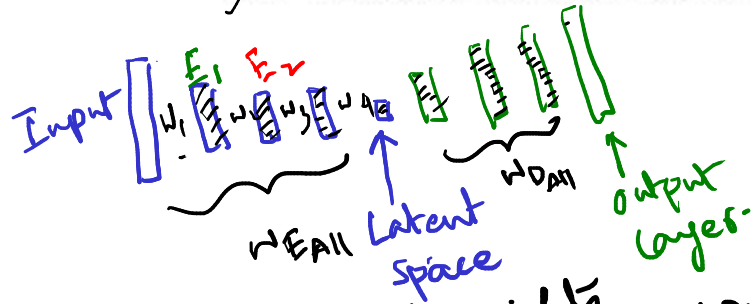① What is the main advantage of layer-by-layer pre-training for deep autoencoders?

$y = Wx + b = [W_1 W_2 \cdots W_n] \begin{bmatrix} x_1 \\ x_i \\ x_n \end{bmatrix}$

Stacked
Autoencoders

a) It reduces the total number of weights and simplifies the optimization process

b) It provides better initial weight values for the entire network

$= [W_1 W_2 \cdots W_n b] \begin{bmatrix} x_1 \\ x_i \\ x_n \\ 1 \end{bmatrix}$

c) It allows for parallel training of different hidden layers.

d) It guarantees perfect reconstruction of the input data with minimal error.

Input  $W_1$  $E_1$  $W_2$  $E_2$  $W_3$  $W_4$  $W_{D_{All}}$  output layer.

$W_{E_{All}}$  Latent Space

$x$  $\hat{x}$   $\ell = MSE(x, \hat{x})$

$W_E$  $W_D$   $\arg\min \ \ell(x, \hat{x})$
$W_E, W_D.$

These pretrained Weights you can use as the first initialized weight values rather than initializing them using a rand() function.

$W(0) \rightarrow rand$

$W(n+1) \leftarrow W(n) - \eta \dfrac{\partial E}{\partial W}.$

② Select the correct option about Denoising autoencoders?

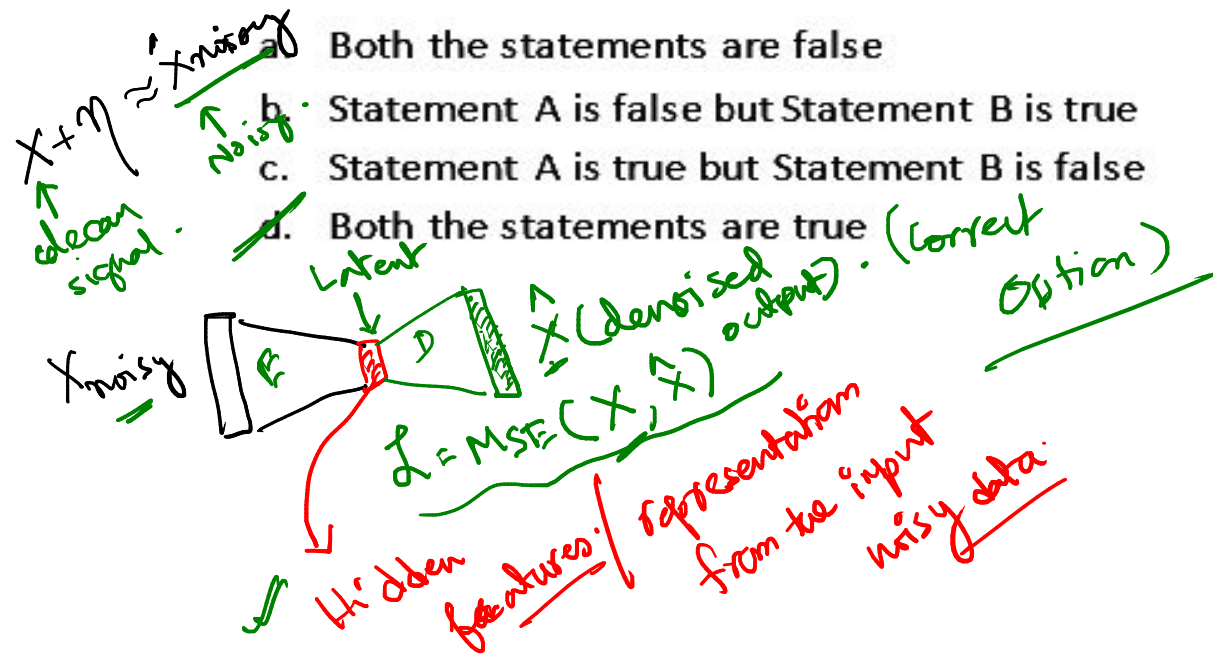Statement A: The loss is between the original input and the reconstruction from a noisy version of the input

Statement B: Denoising autoencoders can be used as a tool for feature extraction.

a. Both the statements are false
b. Statement A is false but Statement B is true
c. Statement A is true but Statement B is false
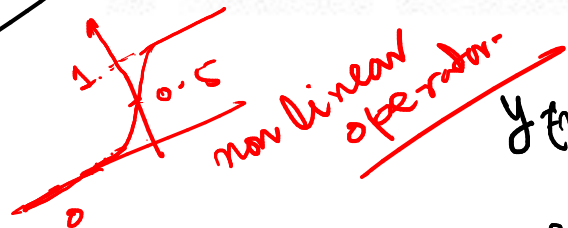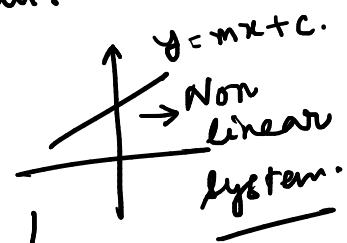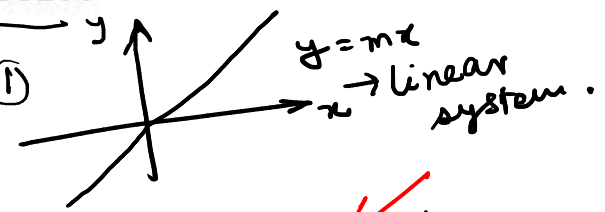d. Both the statements are true (Correct Option)

$X + \eta \simeq X_{noisy}$

Noisy

clean signal.

Latent

$X_{noisy}$   E   D   $\hat{X}$ (denoised output). (Correct)

$\mathcal{L} = MSE(X, \hat{X})$

Hidden features / representation from the input noisy data.

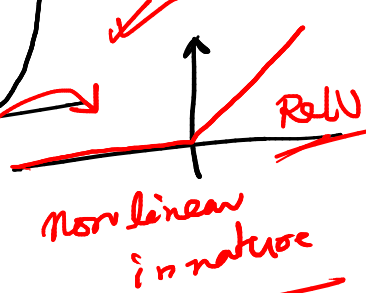③ **Which of the following is a linear operator?** → Superposition Theorem.

$V = IR.$
$\phi = LI.$
$Q = CV.$

a. Sigmoid
b. Rectified Linear Unit
c. Convolution function
d. None of the above

①

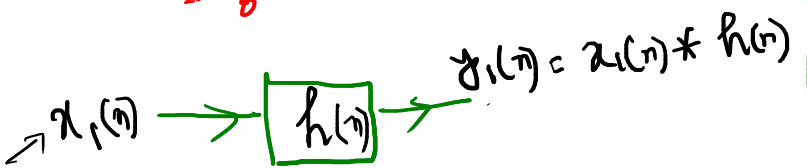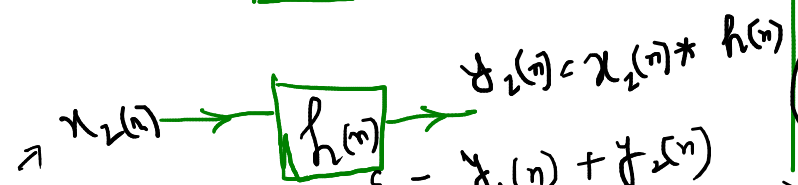$y = mx$ → linear system.

$y = mx + c.$ → Non linear system.

Piecewise linear model of diode.

non linear in nature

$Relu(x) = max(0, x).$

Non linear. ReLU

non linear operator.

$x(n) \rightarrow \boxed{h(n)} \rightarrow y(n)$

$y(n) = h(n) * x(n)$
$= x(n) * h(n).$

$\rightarrow x_1(n) \rightarrow \boxed{h(n)} \rightarrow y_1(n) = x_1(n) * h(n)$

$\rightarrow x_2(n) \rightarrow \boxed{h(n)} \rightarrow y_2(n) = x_2(n) * h(n)$

$y_{total}^s = y_1(n) + y_2(n)$
$= x_1(n) * h(n) + x_2(n) * h(n).$

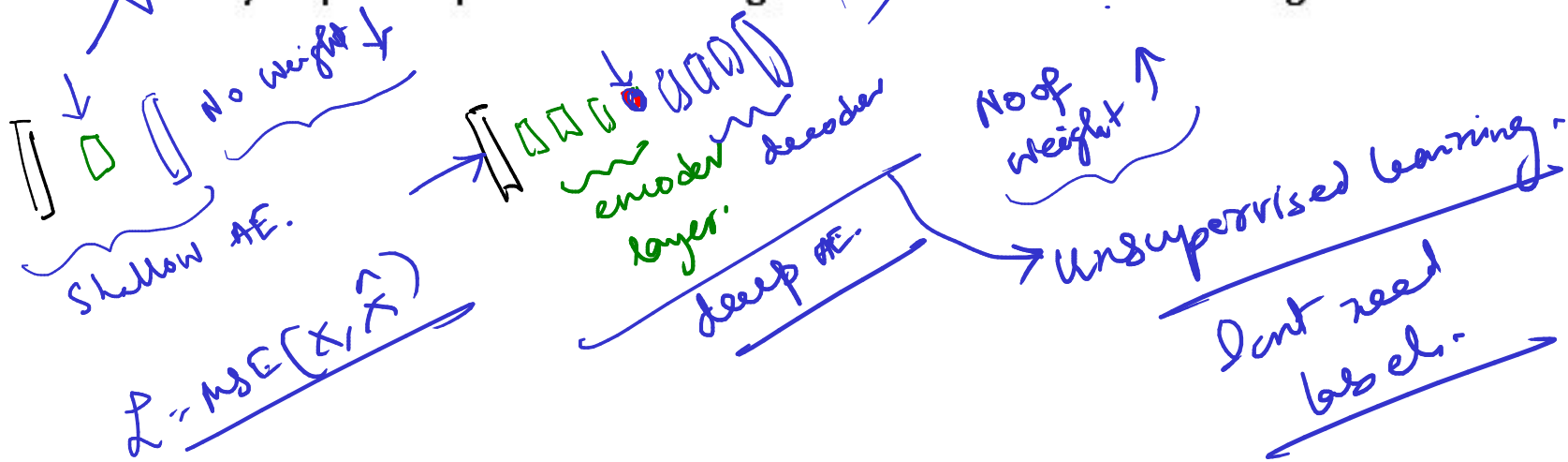$x_1(n) + x_2(n) \rightarrow \boxed{h(n)} \rightarrow y_0(n)$

Conv → linear op.

$y_0(n) = (x_1(n) + x_2(n)) * h(n).$

$y_0(n) = x_1(n) * h(n) + x_2(n) * h(n)$

Super position theorem obeyed $y_0(n) = y_{total}^s(n)$

Which statement is TRUE about deep autoencoders?

a. They have a single hidden layer for representing the latent space.

b. They have few parameters compared to shallow autoencoders.

c. They excel at capturing complex relationships and features in high-dimensional data.

d. They require supervised learning with labeled data for training.

No weight

Shallow AE.

$\mathcal{L} = MSE(X, \hat{X})$

encoder decoder
layer.

deep AE.

No of weight

Unsupervised learning.

Don't need labels.

**Which of the following is false about autoencoder?**

True.

a. Autoencoders possesses generalization capabilities
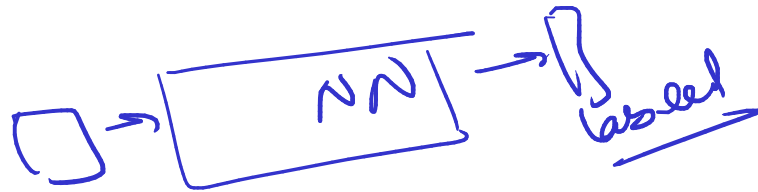
b. Autoencoders are best suited for image captioning task ✗ False

c. Its objective is to minimize the reconstruction loss so that output is similar to input $MSE(x, \hat{x})$ · True

d. It compresses the input into a latent space representation and then reconstruct the output from it True

✓ Anomaly detection
↗ Image segmentation
✓ feature extraction.

NN → labeled

The Dirac delta function is $\infty$ when t=0. Fill in the blanks.
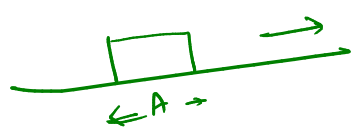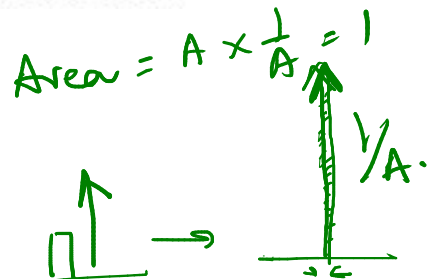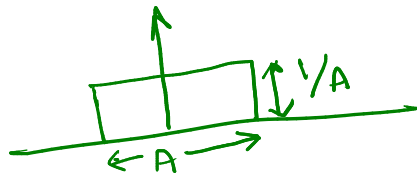
a. 1
b. 0
c. Infinity
d. None of the above

$\rightarrow$ Continuous time signal.

$t \in$ Real.

$\delta(t)$

unit impulse function.

$\delta(n)$.

Discrete time system.

$t \rightarrow n$
$\uparrow$
sample.

$n \in$ integer.

$\delta(n) = 1 \; ; n=0$
$= 0 \; ; n \neq 0.$
elsewhere

Area $= A \times \frac{1}{A} = 1$

$\frac{1}{A}$

$1/A$

$A \rightarrow 0.$

$A \rightarrow 0.$
$\frac{1}{A} \rightarrow \infty.$

① $\delta(t)$ exists at $t=0$ only

② $\delta(t) = 0 \; ; t \neq 0.$

$A \downarrow \quad \frac{1}{A} \uparrow$

$\rightarrow$ Area $= \lim_{A \to 0} \left(A \times \frac{1}{A}\right) = 1$

③ Amplitude of $\delta(t)$ at

$t=0 \; ; \lim_{A \to 0} \frac{1}{A} \rightarrow \infty$

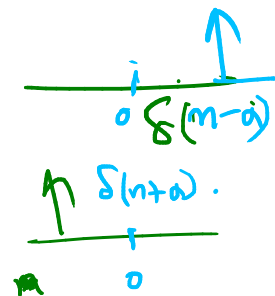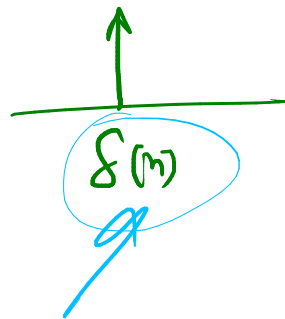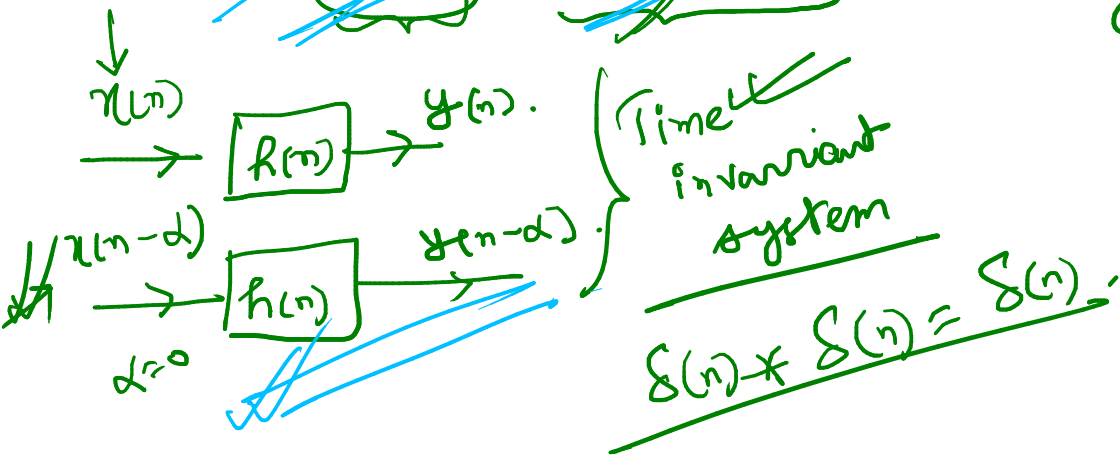$\delta(t) \rightarrow \infty.$

④ $\int_{-\infty}^{\infty} \delta(t) \, dt = 1.$

Impulse response is the output of _____ system due to impulse input applied at time=0. Fill in the blanks from the options below.

a. Linear
b. Time Varying
c. Time Invariant
d. Linear And Time Invariant

Sample $(n)=0$

$n(n)$

$x(n)$ → $h(n)$ → $y(n)$.

$x(n-\alpha)$ → $h(n)$ → $y(n-\alpha)$. } Time invariant system

$\alpha=0$

$\delta(n) * \delta(n) = \delta(n)$.

$\delta(n)$ → $h(n)$ → $y(n)=$ $\delta(n)$. $n=0$

$n=0$

$y(n) = \qquad \delta(n).$

$\delta(n)$

$\delta(m-\alpha)$

$\delta(n+\alpha)$.

$0$

Given the image below where, Row 1: Original Input, Row 2: Noisy input, Row 3: Reconstructed output. Choose one of the following variants of autoencoder that is most suited to get Row 3 from Row 2.

Original Images



Noisy Input



Autoencoder Output
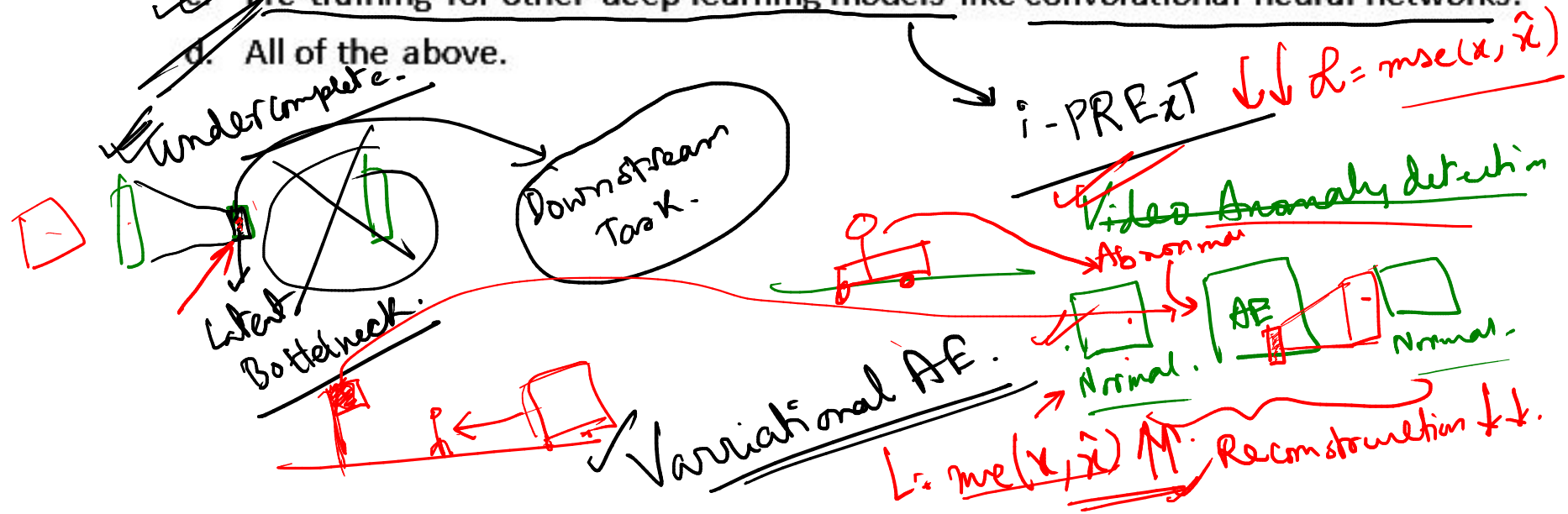


→ Denoising AE.

a. Stacked autoencoder
b. Sparse autoencoder
c. Denoising autoencoder

Which application of deep autoencoders utilizes the extracted latent space representation?

a. <u>Anomaly detection</u> by comparing data points to the known latent space distribution.

b. Image segmentation by classifying individual pixels based on their latent features.

c. Pre-training for other deep learning models like convolutional neural networks.

d. All of the above.

Undercomplete.

Latent

Bottelneck.

Downstream Task.

Variational AE.

i-PREzT

$\mathcal{L} = mse(x, \hat{x})$

Video Anomaly detection

Abnormal

AE

Normal.

Normal.

$L = mse(x, \hat{x})$

Reconstruction f.f.

**(1)** The input image has been converted into a matrix of size 256x256x3 and 4 kernel/filter of size 7x7 with a stride of 2 and no padding. What will be the size of the convoluted matrix?
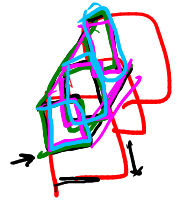
a. 127x127x3

b. 128x128x4

c. 124x124x3

d. 125x125x4

$$I/P = 256 \times 256 \times 3.$$

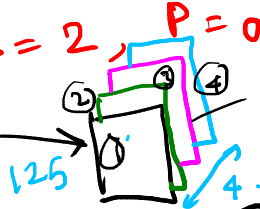$$f = 4, \qquad K = 7 \times 7, \quad S = 2, \quad P = 0.$$

$$O/P = \frac{I/P - K + 2P}{8} + 1$$

$$= \frac{256 - 7}{2} + 1$$

$$O/P = \frac{249}{2} + 1 = \lceil 125.5 \rceil$$

$$= 126.$$

$$256 \times 256 \times 3$$

$$C_{in}$$

$$k = 7 \times 7, \quad S = 2$$

$$P = 0$$

$$f = 4.$$

no of filter $= C_{out}$

$$125$$

$$4.$$

$$H \times W \times C_{out}.$$

$$125 \times 125 \times 4.$$

The figure below shows image of a face which is input to a convolutional neural net and the other three images shows different levels of features extracted from the network. Can you identify from the following options which one is correct?



a. Label 3: Low-level features, Label 2: High-level features, Label 1: Mid-level features

b. Label 1: Low-level features, Label 3: High-level features, Label 2: Mid-level features

c. Label 2: Low-level features, Label 1: High-level features, Label 3: Mid-level features

d. Label 3: Low-level features, Label 1: High-level features, Label 2: Mid-level features

**Which of the following statement is False about ReLU layer?**

$y = ReLU(x) = max(0,x)$

a. ReLU has expression $f(x) = max(0,x)$ ✓ True.

$y = 0$     $y = x$     $\frac{dy}{dx} = 1$

b. The derivative of ReLU is 1 if $x>0$; o otherwise True. $\frac{dy}{dx} = 0$   $dy/dx$

c. Implementation of ReLU has more computational cost than tanh or sigmoid

False

d. ReLU activation function was introduced in AlexNet architecture.

sigmoid (x)

$f(x) = \sigma(x)(1-\sigma(x))$

tanh(x)??

$y =$

$\frac{dy}{dx} = ?$
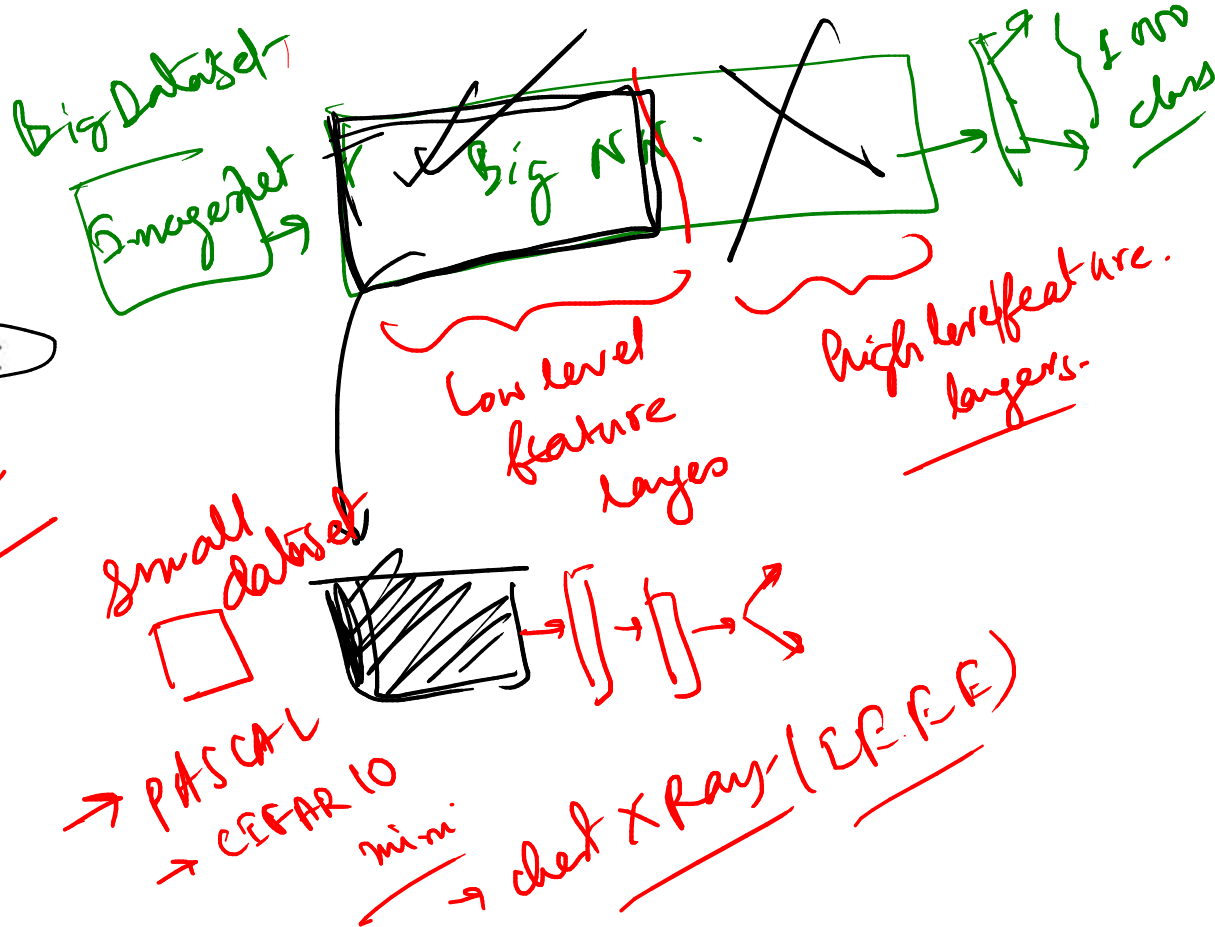
For a transfer learning task, which layers according to you can be more generally transferred to another task?

a. Higher layers
b. Lower layers
c. Task specific
d. Cannot comment

Big Dataset

Imagenet

Big N.N.

1000 class

Low level feature layers

High level feature layers.

Audio pretrained model → Finetuning.

Jannet

Small dataset

PASCAL
→ CEFAR 10
mini

→ chest X Ray (C.R.R.F)

Suppose your input is a 256 by 256 color (RGB) image, and you use a convolutional layer with 100 filters that are each 5x5. How many parameters does this hidden layer have **(with bias)**

a. 2501
b. 2600
c. 7500
d. 7600

$y = Wx + b$

3 much weight

$y = W_1 x + b_1$
$= W_2 x + b_2$
$= W_3 n + b_3$

1st filter
$5 \times 5 \times 3$

I/p.

$256 \times 256 \times 3$.

$f = 100$    $k = 5 \times 5$.

Bias.

$256 \times 256 \times 3$

Total no of Parameters.

Parameter.
$5 \times 5 \times 3 \rightarrow$ Single Kernel.
$\Rightarrow$ 100 such Kernel $\rightarrow (5 \times 5 \times 3) \times 100$.

Weight value.

for $1$ Kernel $\rightarrow 1$ bias
for $100$ Kernel $\rightarrow 100$ bias
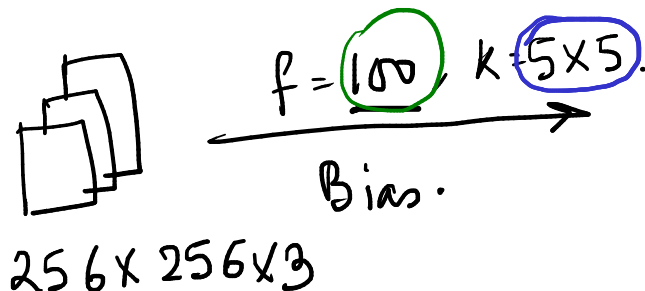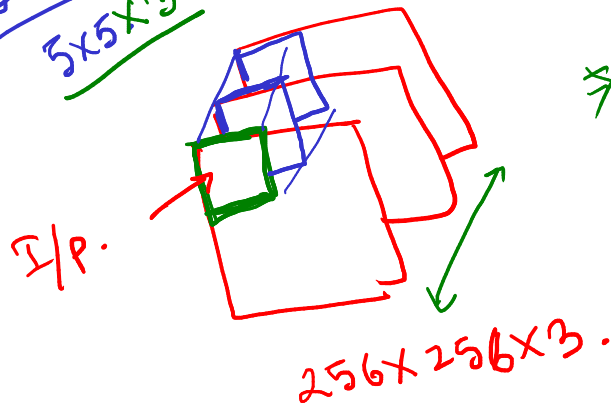
Total parameter = Total weight + Total Bias.
=

Suppose your input is a 256 by 256 color (RGB) image, and you use a convolutional layer with 100 filters that are each 5x5. How many parameters does this hidden layer have (**with bias**)
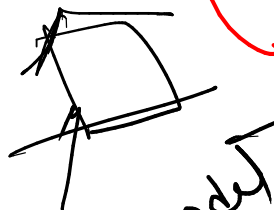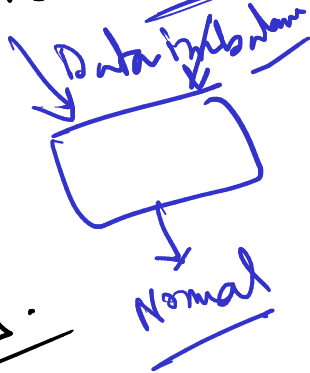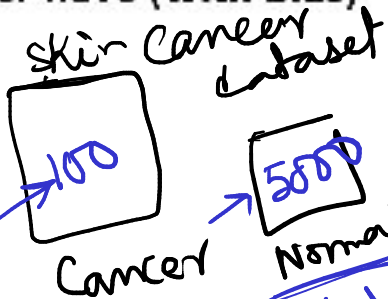
a. 2501
b. 2600
c. 7500
d. 7600

Total weight + Total Bias

= Total Parameter

$(5 \times 5 \times 3) \times 100 + 100$

$(5 \times 5 \times 3)a + a.$
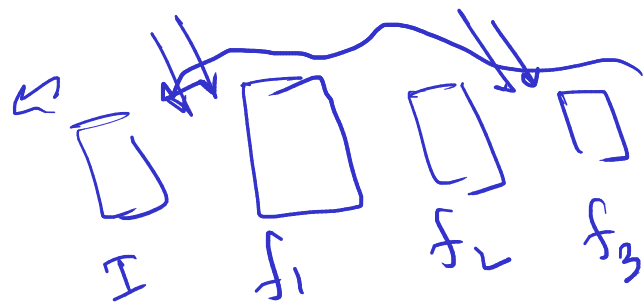
$a[(5 \times 5 \times 3) + 1].$

$= 7600$

$100.$

CNN model
is biased
→ on Tajmahal
→ If it has more
inclination towards
Tajmahal

Skin Cancer dataset

100

Cancer    5000    Normal.

Data Validation

Normal

Statement 1: Adding more hidden layers will solve the vanishing gradient problem for a 2-layer neural network

Statement 2: Making the network deeper will increase the chance of vanishing gradients.

    a. Statement 1 is correct

    b. Statement 2 is correct

    c. Neither Statement 1 nor Statement 2 is correct

    d. Vanishing gradient problem is independent of number of hidden layers of the neural network.

$$O = f_1 f_2 f_3 f_4 f_5 f_6 (I).$$

$$I \quad f_1 \quad f_2 \quad f_3 \quad f_4 \quad f_5 \quad f_6 . O \qquad \frac{\partial O}{\partial I} = \text{chain } I \text{ chain}$$

$$\frac{1}{4}\left(1 - \frac{1}{4}\right) \quad \sigma(n)(1 - \sigma(n))$$

$$W(n+1) \Leftarrow W(n) - \eta \frac{\partial E}{\partial w}.$$

Which of the following is false about CNN?

*True*

a. Output should be flattened before feeding it to a fully connected lyer

*False* b. There can be only 1 fully connected layer in CNN *Wrong*

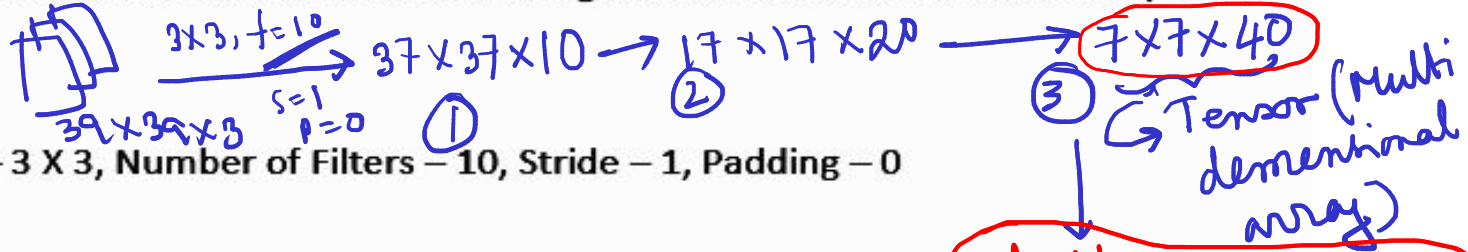c. We can use as many convolutional layers in CNN

*True*

d. None of the above

Vgg-16 . (13)

Resnet 18 (18)
50
101

DenseNet101

Let us consider a Convolutional Neural Network having three different convolutional layers in its architecture as:

$$39 \times 39 \times 3 \xrightarrow[\substack{s=1 \\ p=0}]{3\times3,\ f=10} 37 \times 37 \times 10 \xrightarrow{} 17 \times 17 \times 20 \xrightarrow{} 7 \times 7 \times 40$$

① ② ③ → Tensor (Multi dimensional array)

**Layer-1**: Filter Size − 3 X 3, Number of Filters − 10, Stride − 1, Padding − 0

**Layer-2**: Filter Size − 5 X 5, Number of Filters − 20, Stride − 2, Padding − 0

**Layer-3**: Filter Size − 5 X5 , Number of Filters − 40, Stride − 2, Padding − 0

Flattening operation.

Fully connected layer

Vector: M×1 = 1960×1

Dense

Layer 3 of the above network is followed by a fully connected layer. If we give a 3-D image input of dimension 39 X 39 to the network, then which of the following is the input dimension of the fully connected layer.

$$\frac{39-3+0}{1}+1 = 37$$

$$\frac{37-5+0}{2}+1$$
$$17$$

$$\frac{17-5+0}{2}+1$$

a. 1960
b. 2200
c. 4563
d. 13690

$$M = 7 \times 7 \times 40$$
$$= 49 \times 40$$
$$= 1960$$

Consider a CNN model which aims at classifying an image as either a rose, a marigold, a lily or orchid (consider the test image can have only 1 of the images at a time). The last (fully-connected) layer of the CNN outputs a vector of logits, L, that is passed through a _____ activation that transforms the logits into probabilities, P. These probabilities are the model predictions for each of the 4 classes.

Fill in the blanks with the appropriate option.

a. Leaky ReLU
b. Tanh
c. ReLU
d. Softmax

Multi class classification.

Softmax → Probability of occurance for each class.

$y = \max(0, x)$

$y = x$

$y = 0$

$y = x$

$y \sim x;\ x > 0$

$= \alpha x;\ x < 0$

$y = \alpha x.$

$\alpha \in [0, 1].$

$0.3.$

(10) Imagine you're training a CNN for Autonomous driving vehicle to distinguish between pedestrian, bicycle, bike and cars in images. You have two options:
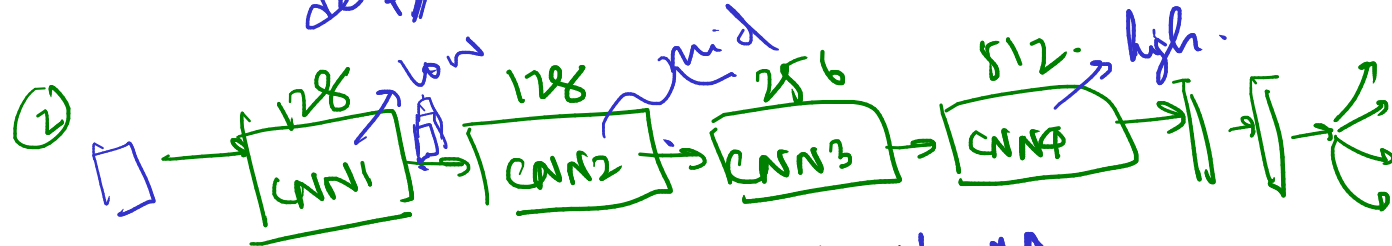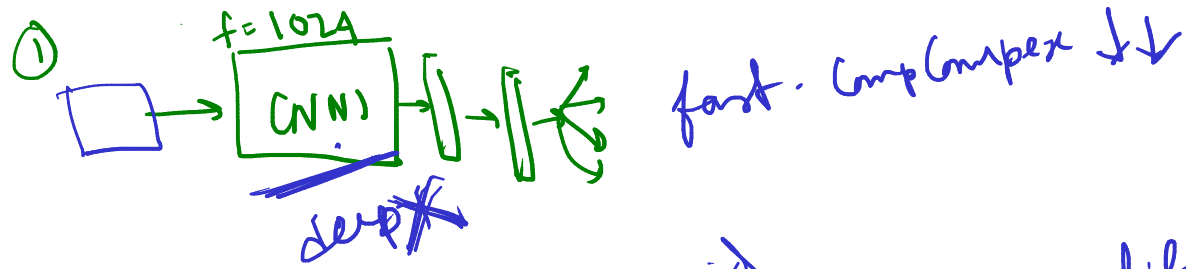
Option 1: A shallow CNN with just one convolutional layer having 1024 kernels and a few fully connected layers.

Option 2: A deeper CNN with 4 convolutional layers having 128 kernels in first layer,128 kernels in second layer, 256 kernels in third layer and 512 kernels in fourth layer and then fully connected layers.

Both options use the same total number of kernels convolutional layer (=1024)

Choose the incorrect statement:
a.  Option 1 will give higher inference speed since it can operate convolutions in parallel whereas option 2 can't be parallelized as results of next convolutional layers are dependent on past layers — } True.
b.  Option 2 , Deeper CNN, with multiple layers, can perform hierarchical feature extraction, thereby has higher representational power and accuracy — } True.
c.  Option 1 , Shallow CNN, with one convolutional, can extract 1024 features from image, thereby has higher representational power and accuracy — } False.
d.  Option 2 , Deeper CNN, with multiple layers, can extract more abstract features that depend on features of shallower layers and therefore has higher representational power and accuracy — } True

① 

f=1024

CNN1  → fast. Comp Complex ↑↓

depth ~~x~~

② 

128 → low  128 → mid  256  512 → high.

CNN1 → CNN2 → CNN3 → CNN4 → high.

→ high computational complexity ↑↑

→ wide variety of feature

low
mid
high.

Dataset :→ 216. (sample/Images).

Batch Size = 6

Six epochs

36 images.

36 buckets.

216
no of
images.

# Thank You

Epoch 1 :—
[ = = = — ...... =]

36tu.
↓
= ]

total Loffs :

total acc
val los

Epoch 1
iteration 1 :→

Epoch 2
[ = — ...... —]

iteration 2 :—  iteration 3 ---- ---- iteration: 36 :—

Batch of 6
CNN
$L_1$
$L_2$
$L$