

Unit-5Pretrained models & other architectures

I. 1. LeNet

2. AlexNet

3. VGGNet

4. DenseNet & Google Net (though not in syllabus)

5. ResNet

6. Transfer LearningII. Improving Deep Neural Networks

1. Hyperparameter tuning -

2. Regularization & optimization.

3. Data Augmentation techniques.

IV. GANs

1. GAN

2. Reinforcement Learning.

Questions from previous year SEE

- (Q) Sketch the core components of a Reinforcement Learning paradigm (10m)
- (Q) Evaluate the following pre-trained deep learning models: (10m)
- i) Google Net
 - ii) VGG-16.
- (or)
- (Q) List different hyperparameters used in Deep learning models. Justify the effect of tuning atleast 4 hyperparameters to improvise the model performance. (10m)
- (Q) Explain the pros & cons of the transfer learning paradigm (10m)

5)

LeNet

i) LeNet (LeNet-5) is a CNN proposed by Yann LeCun in 1989 originally designed to demonstrate the task of handwritten digit classification by a CNN.

Architecture

i) The LeNet CNN architecture has 7 layers : 3 convolutional layers, 2 subsampling layers & 2 fully connected layers.
(Input: $32 \times 32 \times 1$)

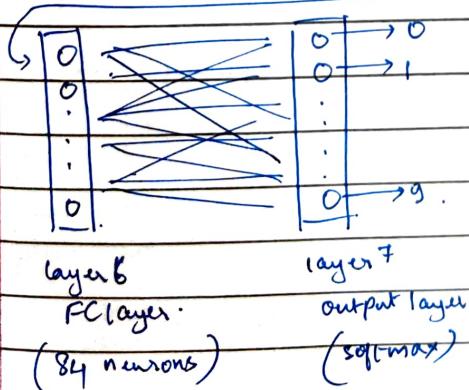
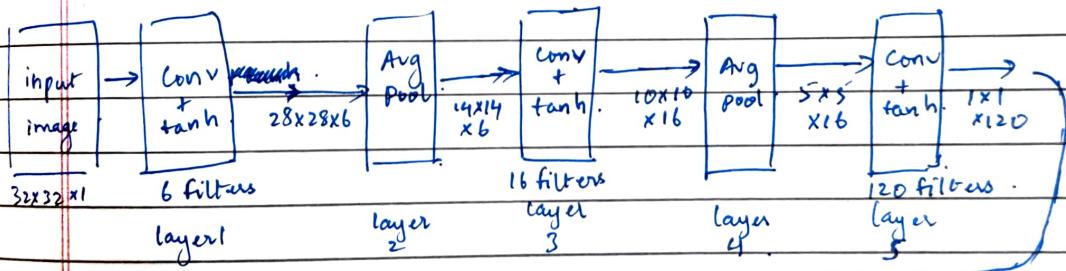
ii) The convolutional layers are used to extract spatial features.

iii) Average pooling is used for subsampling.

iv) tanh function is used as the activation function.

v) The fully connected layers are nothing but Fully connected FFNN used for the final classification.

vi) The 7 layers are arranged as follows



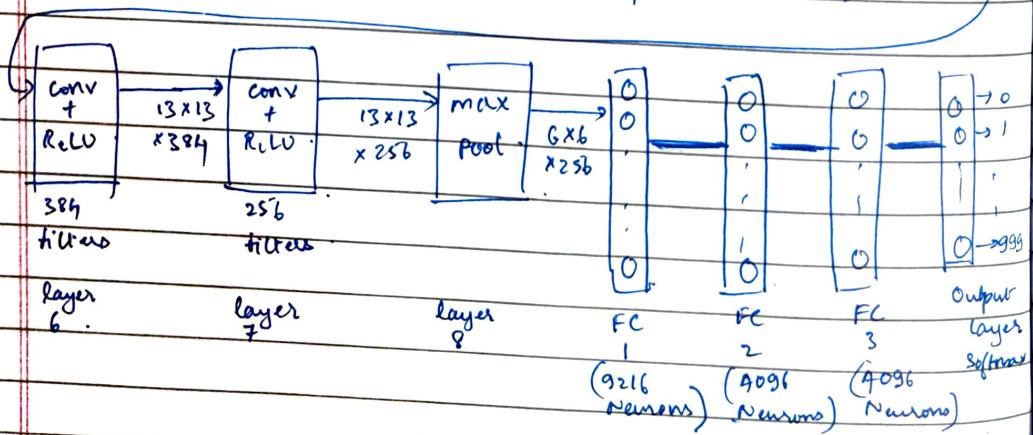
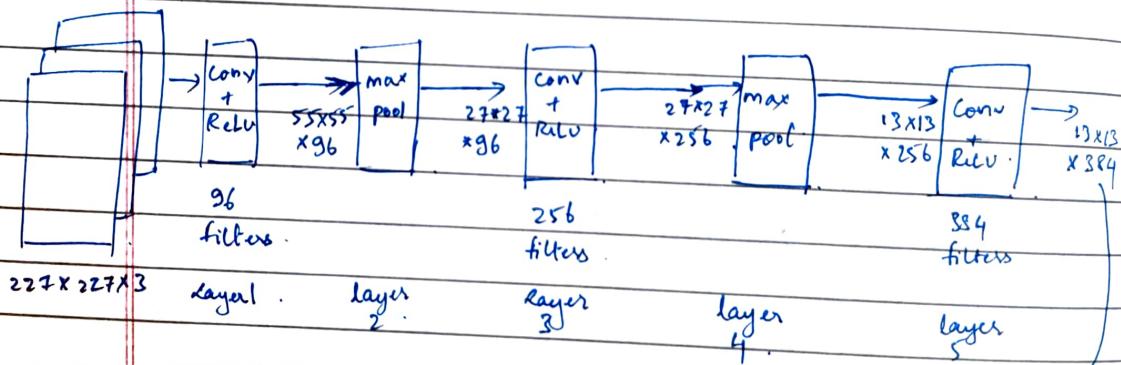
2) AlexNet

i) AlexNet is a CNN architecture proposed by Alex Krizhevsky in 2012 to classify the images of the ImageNet dataset (more than 1 million images with 1000 class labels).

.) This was the first ^(new) architecture that used GPU to boost training performance

.) Architecture:

- i) the AlexNet architecture consists of 12 layers: 5 convolutional layers (ReLU activated), 3 max-pooling layers, ~~2 fully connected layers~~, 3 Fully connected layers, & 1 soft max layer.
- ii) Max pooling is used for subsampling.
- iii) ReLU function is used as the activation function.
- iv) Input dimension: $227 \times 227 \times 3$ (after padding).
- v) Stochastic Gradient Descent (SGD) is used as the optimizer.
- vi) Data augmentation like flipping, jittering, cropping, color normalization etc are carried out.



if Dropout \rightarrow Forward pass \rightarrow Backward pass \rightarrow Forward pass \rightarrow MP \rightarrow C \rightarrow MP \rightarrow C \rightarrow C \rightarrow MP \rightarrow FC \rightarrow FC \rightarrow FC

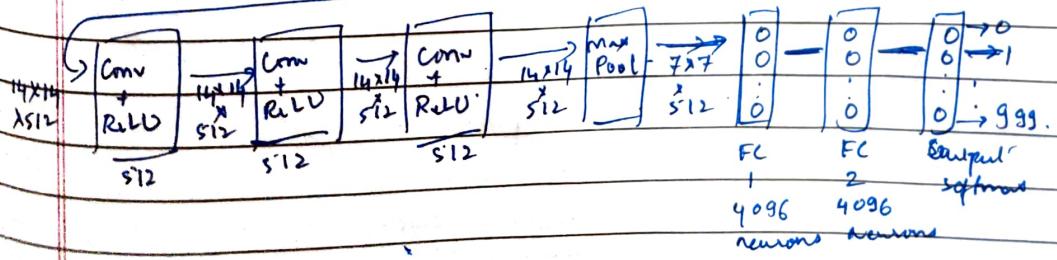
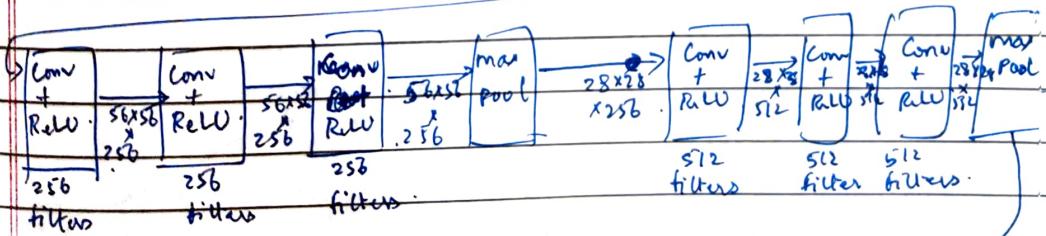
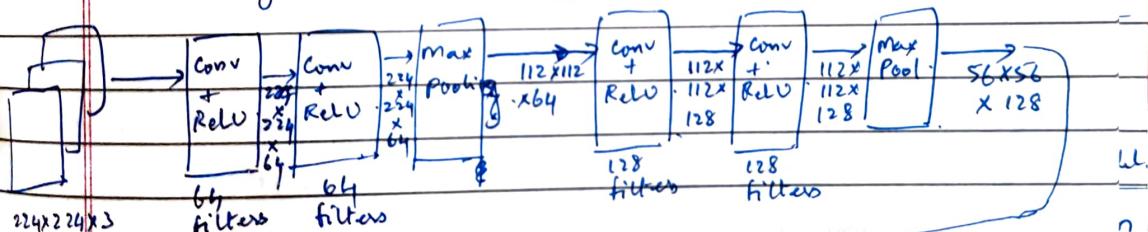
$S \leftarrow FC$

3) VGGNet

- The VGGnet is a deep CNN developed by the Visual Geometry Group in 2014 & one of the most liked models submitted at ILSVRC 2014 (in ImageNet competition).
- The 2 types VGG architecture
 - i) VGG-16
 - ii) VGG-19
- The numbers 16 & 19 refers to the no. of weight layers of the model (pooling layer is not considered a weight-layer).
- Both VGG 16 & 19 have a uniform architecture.

3.i) VGG-16 architecture

- 16 weight layers: 13 Conv layers (ReLU activation) + 3 FC layers (2 hidden layers \rightarrow ReLU activation)
1 output layer \rightarrow softmax activation
- 5 pooling layers (max pooling)



CCM \rightarrow cCM \rightarrow CCCM \rightarrow CCCM \rightarrow CCCM \rightarrow FC1 \rightarrow FC2 \rightarrow S

FC 1
4096 neurons
FC 2
4096 neurons
Output softmax

3.2) VGG-19 architecture

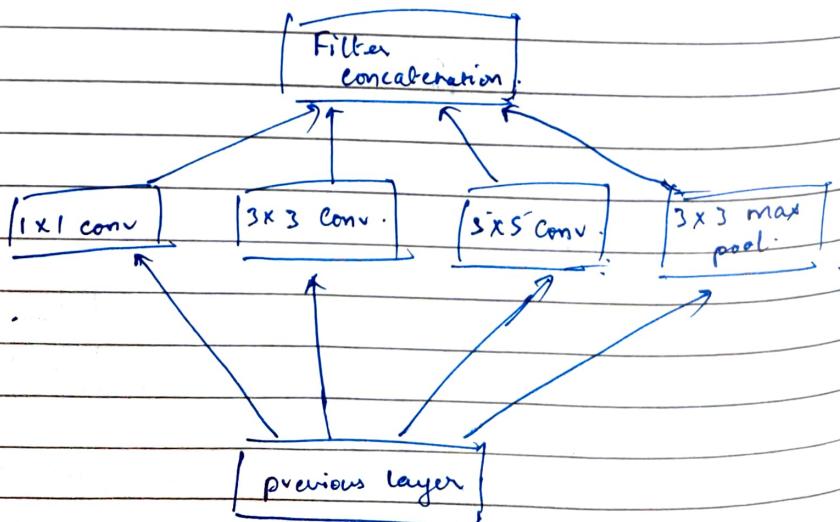
- Has the same basic idea as that of VGG-16, but has 3 additional conv layers.

$CCM \rightarrow CCM \rightarrow CCCM \rightarrow CCCM \rightarrow CCCM \rightarrow FC1 \rightarrow FC2 \rightarrow S$.

Q) ~~What is Google Net?~~

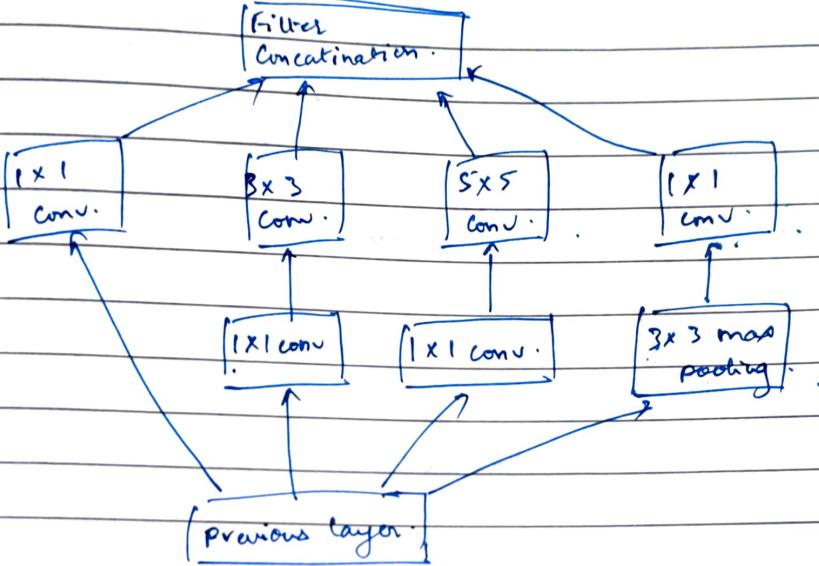
- Google Net (Inception v1) was proposed by Google in 2014 in the research paper titled "Going Deeper with Convolutions".
- It was the winner of the ILSVRC 2014 image classification challenge.
- Architectural features of Google Net:

- The overall architecture is 22 layers deep. The ~~overall~~ architecture was designed to improve computational efficiency (by ~~so~~ decreasing the number of parameters). It ~~only~~ allows parallel convolutions of different dimensions in the same block.
- To achieve parallel convolutions of different dimensions in the same block, Google Net uses an Inception module (naive version).



Inception Module, naive version.

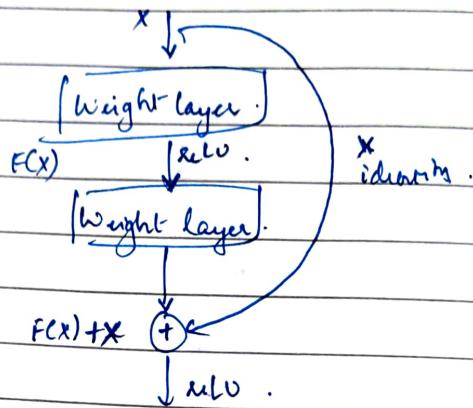
- iii) To decrease the number of parameters (dimension reductions)
 GoogLeNet uses 3 additional 1×1 convolutions in the naive Inception module (in each block) & a Global Average Pooling at the end of the architecture.



Inception module with dimension reductions

-) GoogLeNet was the winner at the 2014 ILSVRC competition & had the top 5 error rate of 6.67% in classification tasks.
- 5) Resnet
 Resnet also called Residual Network was developed by Microsoft research in 2015.
 In deep neural networks, due to the large no. of layers, a problem called as vanishing/exploding gradients occur due to which the training & test errors increase.
 To solve the problem of vanishing/exploding gradients, the Resnet architecture introduced the concept of Residual Blocks.
 A residual block uses the concept of skip connection, wherein the output of a layer is connected to further layers by skipping some layers in between.

-) @ The ResNet architecture was built by stacking these Residual blocks together



4) Transfer Learning (TL)

-) Transfer learning is ML technique wherein a model developed for a particular task is reused as the starting point ~~for a new model~~ for another model for a similar / related task.
-) Instead of training a new model from scratch, TL allows you to leverage the learned parameters from a pretrained-model. (i.e. use the knowledge gained say e.g.: the knowledge of extracting certain features etc).
-) Key concepts related to TL.

i) Pretrained models: there are deep CNNs previously trained on large datasets that require significant computational resources, such as image classification on large image dataset, video/Audio segmentation etc.

ii) Transfer:

iii) Feature Extraction: In TL, the early layers of a pre-trained model are kept (which capture the general features like edges, textures, shapes etc) & only the final layers (usually the FC layers) are retrained as required by the new task.

iii) Fine-tuning: Sometimes, the weights of the early layers from a pre-trained model are fine-tuned (adjusted) helping the model to adjust to the specifics of the new task.

i) Pros of TL.

i) Reduced Training Time: Since the model uses the pre-trained weights from a related task, the training time on the new task is significantly reduced.

ii) Improved Performance with Limited Data: TL is beneficial in tasks with a small dataset.

This is because the pre-trained model already passes knowledge from a larger related dataset.

iii) Efficient Use of Computational Resources: Since we don't have to train our model from scratch, TL facilitates efficient use of computational resources particularly.

iv) Better Generalization: Models trained using TL usually often generalize better than models trained from scratch.

ii) Cons of TL.

i) Negative Transfer: If the source & the target tasks are not sufficiently related, then TL can harm the performance. This phenomenon known as -ve transfer, occurs when the pre-trained model's features are not related to the required application to the new task.

ii) Overfitting Risk: Fine-tuning a pre-trained model on a small dataset can lead to overfitting.

iii) Complexity in Model Selection: Determining which pre-trained model to choose, which layers to freeze. Fine-tune can

iv) Data privacy concerns: Using a pre-trained model trained on sensitive \oplus proprietary data raises data privacy concerns.

II) Empowering Deep Neural Networks

i) Hyperparameter Tuning

- Most deep learning ~~most~~ algorithms have several hyperparameters which control many aspects of the ~~most~~ algorithm's behavior.
- Some hyperparameters affect the time & space cost of running the algorithm.

Some hyperparameters ~~can~~ affect the quality of the model obtained by the training process

- The hyperparameters are set before the training process & unlike the model parameters (weights & biases), hyperparameters are ~~not~~ learned during the training process

- 2 ways to set hyperparameters -

i) Manual Hyperparameter tuning: Manually set the hyperparameters (trial & error) based on some best practice/heuristic.

ii) Automatic Hyperparameter Tuning: Based on some algorithm like Grid Search, Random Search.

- Few common hyperparameters of a deep learning algorithm (DNN)
 - i) Learning rate (η).
 - ii) Depth of the Neural Network: No. of layers of the net (No. of hidden layers)
 - iii) Width of a layer: No. of neurons in a layer
 - iv) Convolution Kernel width/dimension
 - v) Implicit zero padding -

vii) Dropout rate

2) Effect of these hyperparameters on the model's capacity.

Hyperparameter	Increases model capacity when...	Reason
i) No. of hidden units (in no. of hidden layers)	increased	Increasing the no. of hidden units increases the representational capacity of the model.
ii) Learning rate (η)	tuned optimally.	Too high @ too low η results in a model with low effective capacity due to optimization failure.
iii) No. of neurons in a layer	increased	directly proportional to the representational capacity of the model.
iv) Convolution kernel width	increased	Increasing the kernel width increases the no. of parameters in the model.
v) Implicit zero padding	increased	Padding before convolution keeps the representation size large.
vi) Dropout rate	decreased	
vii) Regularization & Optimization.		

2) Regularization

Regularization techniques are used to prevent the model from overfitting to the training data by adding a penalty term to the loss function.

Overfitting occurs when a model performs well on training data but poorly on unseen data.

1.) Commonly used Regularization techniques :

- i) L1 & L2 regularization.
- ii) Dropout.
- iii) Batch normalization.

2.) Optimization .

Optimizers or Optimization algorithms are techniques to minimize the value of the loss function as effectively as possible & thus return the model's parameters corresponding to that are global/local minima.

$$\text{i.e. } w^* = \underset{w}{\operatorname{argmin}}(\text{Loss})$$

3.) Commonly used Optimization techniques .

i) Gradient Descent

↳ Stochastic Gradient Descent (SGD) .

↳ Mini-Batch Gradient Descent .

ii) Adam .

3.) Data Augmentation Techniques .

These are techniques used to increase the diversity & quantity of the training data for deep learning models .

b.) Advantages of data augmentation

- i) Improve model's performance .
- ii) Reduce overfitting .
- iii) make model more robust to variations in i/p data .

c.) Few Basic image data augmentation techniques

- i) Flipping .
- ii) Rotating .
- iii) Cropping .
- iv) Scaling .

- .) Few advanced image data augmentation techniques .
i) Color Jitter
ii) Random Erasing.
iii) Mixup.

IV.

- 2) Reinforcement learning .

- 3) Sketch the core components of a Reinforcement learning paradigm.

Ans.) Reinforcement Learning (RL) is an ML paradigm wherein a learner/agent learns to make decisions by interacting with an environment.

.) Every action of the agent is associated with a reward (either positive or negative). The agent's goal is to maximize cumulative reward overtime.

- .) Core components of a RL learning paradigm .

i) Agent : The learner / decision maker in the RL environment.

The agent interacts with the environment, takes actions & learns from the outcomes of the actions.

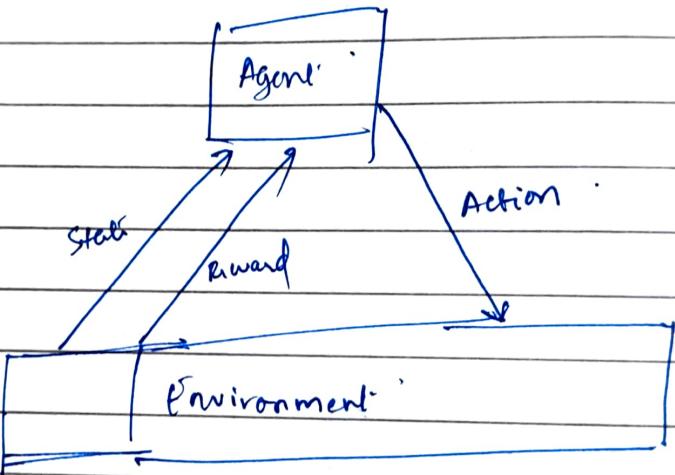
ii) Environment : The external system with which the agent interacts

iii) States(S) : The set of all possible states of the environment.
A current state provides the representation of the current situation of the environment.

iv) Actions(A) : Set of all possible moves an agent can make.
The agent selects an action based on the current state & its policy.

v) Policy $\pi(t)$: A strategy that the agent follows to divide actions based on the current state.

vi) Reward (R): A scalar value received after each action indicating the immediate benefit of the action.
 can either be a +ve reward (good action) or -ve reward (bad action).



$$(s_0) \xrightarrow[a_0]{r_0} s_1 \xrightarrow[a_1]{r_1} s_2 \xrightarrow[a_2]{r_2} \dots$$