

Mid Term Assessments

Name waqar kaleem khan

Enrollment No 01-249191-013

Course Deep Learning

Teacher Dr Imran Siddqi

Submission Date 28-05-2020

Class MSDS(3-A)

Q NO 01 (i) Difference between using pre-trained models as feature extractor and fine-tuning the pre-trained model.

Ans:- The pre-trained models can be used as a separate feature extraction program in which what we do is that input can be pre-processed by the model or portion of the model to a given an output for example a vector of numbers for each input image. So that we can use that as ~~an~~ an input when we are training new model.

Alternatively the pre-trained model or desired portion of a model can be integrated directly to a new model. In this case the weights of the pre-trained model can be chosen to be not updated when training a new model otherwise the weights will be updated when training new model.

The fine tuning of a pre-trained model is that to tweak ~~to~~ parameters of already trained model. So the model can adopt the new task. The starting layer of the model learn very general features and as we go further to more layers the layers tend to learn more specifically about the task it is being trained on. So for fine-tuning we want to freeze or intact the initial layers and retained the later layers for our tasks.

2

(Q No 1 (ii))  
No of employees = 500

No of picture for each employee = 10

If we multiply 500 to 10 then overall

train dataset =  $500 \times 10 = 5000$

(ii) test data per for each employee is 2

Overall dataset =  $500 \times 2 = 1000$

So we have 5000 images for training  
and 1000 images for testing

number of all images = 6000

input size of image =  $100 \times 100 \times 3$

Steps we will follow to create the system  
are below

- (i) we will use Anaconda (jupyter notebook)
- (ii) Install the essential Libraries using "conda"  
or pip3 commands
- (iii) After installing import that libraries to  
jupyter notebook
- (iv) Next step is to load the train and  
test data set.
- (v) use numpy arrays to convert the pixel  
into arrays because in python we have  
one dimensional array so we use numpy.
- (vi) giving shape to the data according to  
the model.
- (vii) Then split data for the model in test  
train if the data is not splitted  
before best practices are 80% training  
and 20% testing means the data  
is prepared for model
- (viii) Load the model without the classifying part

(ix) Then we will have convolutional layer in model the data will be passed through the convolutional layer. To get the important feature in the image pixels layer which is closer to input layer will begin to detect simple feature like edges etc

(x) In the last layers add SVM classifier for classifying non-linear objects because we are using images data

(xi) Train the classifier

(xii) And in the last part the output layer have 500 neuron and using softmax activation function.

1

Q NO #2 (i) Maxpooling with filter size 3x3  
padding = 1 and stride = 1

2	2	2
5	5	5
1	1	1

adding padding = 1 to above image  
padding = 1 so we will add padding  
to the four sides of the image

0	0	0	0	0
0	2	2	2	0
0	5	5	5	0
0	1	1	1	0
0	0	0	0	0

Now applying filter for max pooling with  
stride = 1  
Note black color is used for filter

Step = 1

0	0	0	0	0
0	2	2	2	0
0	5	5	5	0
0	1	1	1	0
0	0	0	0	0

Max value<sub>1</sub> = 5

Step = 2

0	0	6	0	0
0	2	2	2	0
0	5	5	5	0
0	1	1	1	0
0	0	0	0	0

value<sub>2</sub> = 5

Step = 3

0	0	0	0	0
0	2	2	2	0
0	5	5	5	0
0	1	1	1	0
0	0	0	0	0

value<sub>3</sub> = 5

Step = 4

0	0	0	0	0
0	2	2	2	0
0	5	5	5	0
0	1	1	1	0
0	0	0	0	0

value<sub>4</sub> = 5

Step = 5

0	0	0	0	0
0	2	2	2	0
0	5	5	5	0
0	1	1	1	0
0	0	0	0	0

value<sub>5</sub> = 5

Step = 6

0	0	0	0	6
0	2	2	2	0
0	5	5	5	0
0	1	1	1	0
0	0	0	0	0

value<sub>6</sub> = 5

Step = 7

0	0	0	0	0
0	2	2	2	0
0	5	5	5	0
0	1	1	1	0
0	0	0	0	0

value<sub>7</sub> = 5

3

Step 8

0	0	0	0	0
0	2	2	2	0
0	5	5	5	0
0	1	1	1	0
0	0	0	0	0

value 8 = 5

Step 9

0	0	0	0	0
0	2	2	2	0
0	5	5	5	0
0	1	1	1	0
0	0	0	0	0

value 9 = 5

Maxed pool values = 
$$\begin{array}{|c|c|c|} \hline 5 & 5 & 5 \\ \hline 5 & 5 & 5 \\ \hline 5 & 5 & 5 \\ \hline \end{array}$$

B) In maxpooling the parameters did not increase it is used to preserve high values in the feature map

Q2(ii) Briefly discuss whether or not a CNN preserves relative spatial relationships between components

Ans CNN do not preserves relative spatial relationship between components for example that we have two images of face one have every object like eye and nose etc in right place The other image objects are not in same place so to a CNN both pictures are similar because the both contain similar elements.

## Q#3 Comprehensive analytical review on advancement in convolutional Neural Network

**Abstract:** In this review we will discuss about that how the convolutional Neural Networks have evolved and what are the popular architectures and what are the recent trends in the CNN's. If we see the history of deep learning then we can see that in last years on the variety of problem's the deep learning led to a very good performance which include the areas of computer vision natural language processing (NLP) speech recognition. The CNN's among different neural network have been studied most extensively. While improvements in the graphic processor units (GPU) and the rapid growth in the annotated data the research on convolutional neural network got focused and achieved state of art of results on the various tasks. In this report the CNN's improvements have been detailed using different aspects like loss function, activation function, different layers etc. and also in

Let's well discuss some of the applications of the CNN's

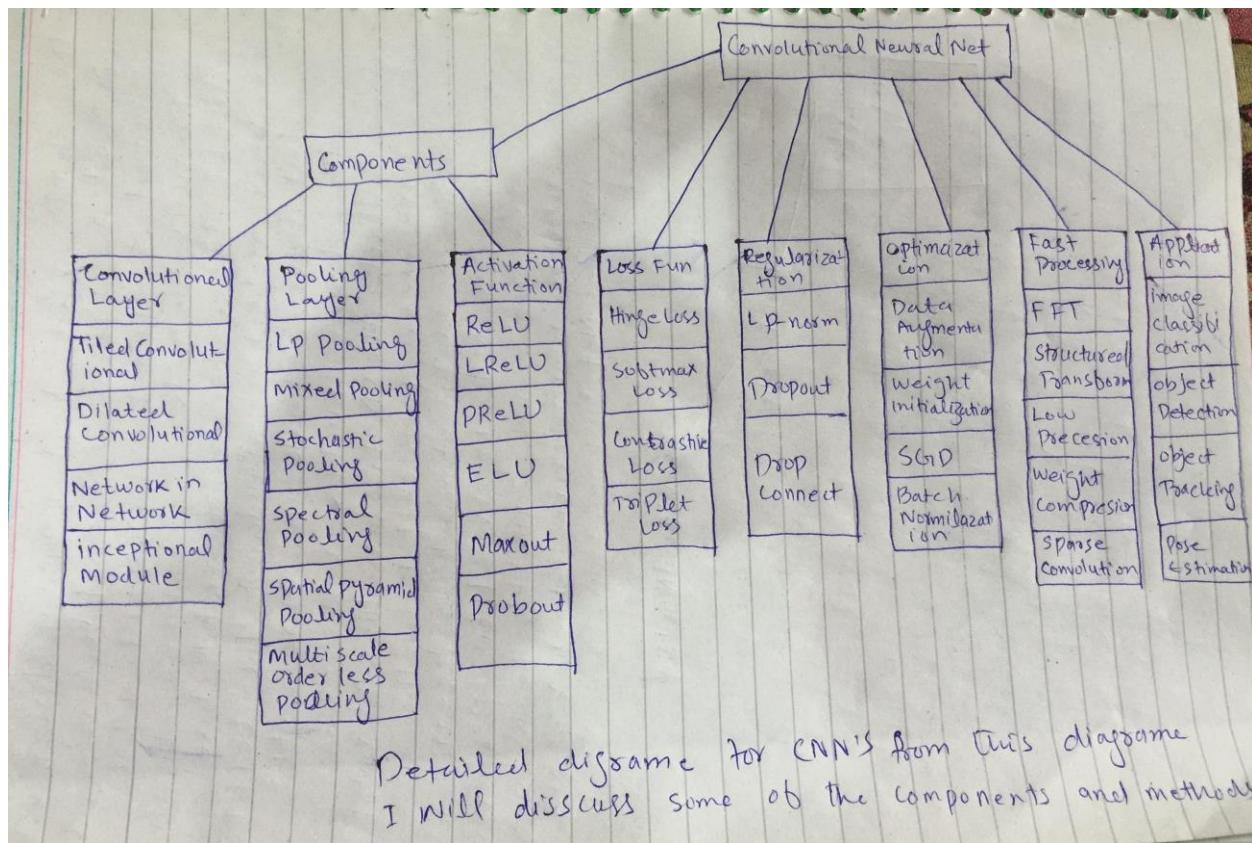
**Introduction:** In the introduction part we will discuss how CNN's work's what is the CNN's architecture and history.

CNN's are well known deep learning architecture which is a memory of human's or living visual mechanism. In 1959 two scientist named Hubel and Wiesel that cells in animal visual cortex are responsible for detecting light in receptive fields.

Two other computer scientist inspired by the above theory and proposed the neocognitron in 1980 which could be regard as the predecessor of CNN's. In 1990 LeCun in his research establish the modern framework of CNN. They developed a multi-layer Artificial neural network called LeNet5, which could be used to classify the handwritten digits but due to lack of large training data and not have a good computational power their model didn't perform good on the complex problems since 2006, many methods have been introduced to overcome the difficulties encountered in

3

The deep CNN's. Most notably Krizhevsky proposed a classical CNN architecture and showed significant improvements upon previous methods on the image classification task. The overall architecture of their method is similar to the LeNet's but their network was more with the deeper structure which is called AlexNet. When the AlexNet have been proposed then the other architecture have been proposed to improve their performance. From which some of the representative works are ZFnet, VGGNet and ResNet.



Some Basic Components of CNN's:- The basic components of the CNN's are very similar. Taking one of the most famous architecture LeNet-5 as an example. the architecture consist of three type layers which are convolutional, pooling, and fully-connected layers. The convolutional layer aims to learn feature extraction of the inputs. Convolutional layer is composed of several convolutional kernels which are used to compute different feature maps. Specifically each neuron of a feature map is connected to a region of a neuron to of a previous layer. The new feature map can be obtained by first convolving the input with a learned kernel and then applying an element wise non-linear activation function  $\sigma$ .

Mathematically  $(i, j)$  in  $k$ th feature map of  $l$ -th layer. is calculated.

$$z_{i,j,k} = w_k^T x_{i,j} + b_k$$

Explaining the above formula where the  $w_k^l$  and  $b_k^l$  are the weight vector and bias term of the  $k$ -th filter of the  $l$ -th layer and  $x_{i,j}^l$  is the input patch centered at location  $(i, j)$  of the  $l$ -th layer. The Kernel  $w_k^l$  that generates the feature

map  $Z_{ij,k}$  is shared. Such a weight sharing mechanism has several advantages like it can reduce the model complexity and made the network easier to train.

The Pooling layer aims to achieve shift-invariance by reducing the resolution of the feature map. It is usually placed between two convolutional layers. Each and every feature map of a pooling layer is connected to its corresponding feature map of preceding convolutional layer.

Mathematically we can denote the pooling as

$$y_{i,j,k} = \text{Pool}(a_{m,n}, k) \forall (m,n) \in R_{ij}$$

Explaining the above mathematical form where  $R_{ij}$  is a local neighbourhood around location. The typical average operation are average pooling and max pooling. In the first convolutional layer the kernel are designed to detect low level features such as edges and curves, while the high layer kernel are designed to learn how to encode more abstract features.

If we look to some of the CNN's architecture after several convolutional and pooling layer there may be one or more connected layers which aim to perform high level reasoning.

What in the fully connected layer happen is that all the neurons of the previous layer are connected to the current layer for generating a global semantic information.

After the fully connected layer the last layer of the CNN's architecture are the output layer. If your model have developed for some classification then mostly softmax activation function is used. Also sum can be used to combined with CNN's and can be used to solve different classification problems. Let us assume

that  $\theta$  denote all the parameters of a CNN for example (weight vector, bias) by minimizing the appropriate loss function the optimum parameter for a specific task can be obtain

$$L = \frac{1}{N} \sum_{n=1}^N l(\theta; y^{(n)}, o^{(n)})$$

$N$  = desired input-out relation  $(x^{(n)}, y^{(n)})$

$x^{(n)}$  = is the  $n$ th input data

$y^{(n)}$  = label corresponding target

$o^{(n)}$  = output of the CNN

Note:- In the first part of this review we explain how CNN's have been developed and how its working and some history of CNN

In the second part we will discuss the improvements or Advancements made to CNN's

Improvements or Advancements on CNN's:-

Since the success of AlexNet in 2012

after that there have been a lot

of improvements on CNN's we will

discuss some of the major improvements

but there are a lot of areas which

have been improved in the CNN's

Some of them are Convolution layer,

pooling layer, activation function,

loss function, regularization, optimization

etc.

Convolutional layer:-

In basic CNN's the

convolution filter is the generalized

linear model (GLM) for the patch under-

lying of the local image. This works

well for abstraction when instances of

latent concepts are linearly separable

For enhancing its representation ability

there is some improvements which will

be discuss below

(\*) Tiled Convolution:- The number of parameters

can be drastically

decrease by weight sharing mechanism

in the CNN's and it can also restrict the

model to from learning other kind of

invariance. Tiled CNN is a variation

for learning rotational and scale invariant features the tiled CNN's tiles and multiply feature maps. Within the same layer separate kernels can be learned and by square-root pooling the complex invariances can be learned implicitly. The convolution operation are applied every  $k$  units where  $k$  is the tile size to control the distance over which weights are shared. The units within each maps will have same weights when the tiles  $k=1$  and to the traditional CNN the tiled CNN becomes identical.

The experiment done on the NORB and CIFAR-10 dataset which results that  $k=2$  show the best results the tiled CNNs can perform better than traditional CNNs in small time series datasets.

\* Transposed Convolution:- Transposed convolution can be seen the backward pass of a traditional convolution. It's also known as deconvolution and fractionally strided convolution. In the deconvolution we associate a single activation with multi output activation while on the other hand the traditional convolution connects multiple input activation to a single activation. The dilation factor for the input feature map are given by the stride of deconvolution.

What deconvolution do is That first the input is upsample by the factor of stride then perform the convolution on the upsample input. Recently The deconvolution has been used for many task visualization, recognition etc.

### Pooling layer

Pooling layer is an important concept of CNN. It lowers the computational burden by reducing the number of connection between convolutional layers. In the pooling section we will discuss about some recent pooling methods

#### L<sup>p</sup> pooling:-

L<sup>p</sup> pooling is inspired by biologically complex cells. It has been theoretically analyzed which suggest that L<sup>p</sup> pooling provides better generalization than max pooling

mathematically L<sup>p</sup> pooling can be represent

$$y_{i,j,k} = \left[ \sum_{(m,n) \in R_j} (a_{m,n,k})^p \right]^{1/p}$$

$y_{i,j,k}$  is the output of the pooling operator at location  $(i,j)$  in the  $k$ -th feature map

$a_{m,n,k}$  is the feature value at location  $(m,n)$

If the  $p=1$  then L<sup>p</sup> corresponds to average pooling and when  $p=\infty$  L<sup>p</sup> reduce to max pooling.

\*) Mixed pooling :- mixed pooling is inspired by random Dropout and Dropconnect. The main purpose of the max pooling method is that combination of Max pooling and average pooling.

Mathematically representation of max pooling fun

$$y_{i,j,k} = \lambda \max_{(m,n) \in R_{ij}} a_{m,n,k} + (1-\lambda) \frac{1}{|R_{ij}|}$$

$$\sum_{(m,n) \in R_{ij}} a_{m,n,k}$$

As we can see  $\lambda$  in above formula which is random value which can either be 0 or 1 which means either using average pooling or max pooling. When the model is in process of forward pass then the  $\lambda$  value is recorded and during back propagation it can be used.

Activation function:- A proper activation function significantly improve the performance of the CNN for a certain task. In this section we will discuss one of all the most and recent used activation function.

ReLU: Rectified linear unit (ReLU) is one on the most notable, useable and non-saturated function

$$a_{i,j,k} = \max(z_{i,j,k}, 0)$$

Where  $z_{i,j,k}$  is the input of the activation function at location  $(i, j)$  on the  $k$ -th channel. ReLU is a piecewise linear function which retain

the positive part and prune the negative part to zero. The max operation of ReLU allows to compute much faster the tanh() and softmax(). With ReLU the deep networks can be trained without pre-training. At value 0 ReLU performance can be lost at backpropagation.

Loss Function:- For a specific task it is important to choose a specific loss function. Advancement done in the loss function is  
 (1) Hinge loss (2) softmax loss, (3) Contrastive loss  
 (4) Triplet loss We will explain one because of 2000 word policy by examiner.  
 (\*) Hinge Loss :- For the training of large margin classification Hinge loss is used

Like SVM the hinge loss function of multi-class SVM is defined in ~~eq~~ below

$$L_{\text{hinge}} = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \max(0, 1 - \delta(y^{(i)}, j) \bar{w}^T x_i)^p$$

$\delta(y^{(i)}, j) = 1$  if  $y^{(i)} = j$  otherwise = -1

if  $p=1$  Hinge loss ( $L_1$ -loss)

if  $p=2$  Hinge loss ( $L_2$ -loss) which is squared Hinge

**Regularization:-** In the deep CNNs regularization used to reduce the overfitting. Overfitting is the unnegetable in deep CNN's. For overcoming this problem the solution's we have is  $\ell_p$ -norm regularization, Dropout and Dropconnect. We will explain one of the above advancement in the Regularization.

- ④  **$\ell_p$ -norm Regularization:-** The term which have been penalize the model complexity the  $\ell_p$ -norm regularization add additional term to modifies the objective of function formula.

$$\epsilon(\theta, xy) = L(\theta, xy) + \lambda R(\theta)$$

$R(\theta) \Rightarrow$  is regularization term

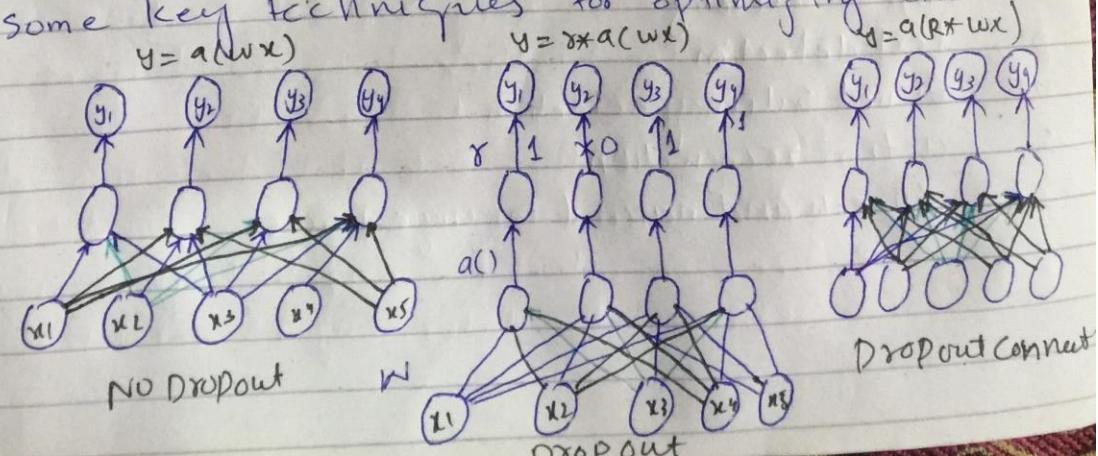
$\lambda \Rightarrow$  is the regularization strength usually the  $\ell_p$ -norm function employed as

$$R(\theta) = \sum_j \|\theta_j\|_p^p$$

if  $p \geq 1$  the  $\ell_p$ -norm is convex which means the optimization easier and renders this function

if  $p=2$  then  $\ell_p$ -norm regularization is commonly referred to as weight decay.

**Optimization:-** In this section we will discuss some key techniques for optimizing CNNs



**Data Augmentation:-** Deep CNNs are particularly dependent on the large quantities of training data availability. So the data augmentation is the elegant solution to alleviate the relative scarcity of data compared to the number of parameters. What data augmentation do is the transform the available data in to the new data without altering their nature. Some of popular augmentation methods are simple geometric transformation ~~and~~ such as mirroring, rotating, shifting.

Paulin propose a greedy strategy that selects the best transformation from a set of candidate transformations. but Their strategy involves a large number of model-training which be expensive computationally.

**Applications of CNNs:-** In the application section we will discuss or highlight some recent work done using CNN's listed below are areas where CNN's are applied are achieve state of art Performance.

- 1) Image classification
- 2) object tracking
- 3) text detection
- 4) visual saliency detection
- 5) action recognition
- 6) speech and natural language process