CISC 867 : Deep Learning.
Assignment #2
Name :- Hanan Fared Mohamed Omara.
ID :- 20 39 8559

1️⃣ input = 500 * 500 * 3 , 100 hidden units is used
→ The shape of the weight matrix of this layer ( without the bias)
= input * hidden units
= 500 * 500 * 3 * 100 = 75 million. shape → $(500*500*3, 100)$
→ The shape of the bias = 1 * 100 = (100,) where 100 is
the number of hidden units in the layer.

2️⃣ 10 filters , size of kernel = 5*5
The number of Parameters = $(F_h * F_w * L + 1)$ * number of filter
= (5 * 5 * 3 + 1) * 10
= 760

3️⃣ left (vertical edge detector) = $\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$

Right ( Horizontal edge detector) = $\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$

4️⃣ The bias will disappears when average large numbers

$$V_r = \frac{V_r}{1 - \beta_r} \quad \text{when } r \to \infty$$

and when it divide by small number will get larger.

1️⃣

**5** The Size batch $m$, input $z = (z^1, \ldots, z^m)$

output $M = \frac{1}{m} \sum\limits_{i=1}^{m} z_i$

$$\sigma^2 = \frac{1}{m} \sum\limits_{i=1}^{m} (z_i - M)^2$$

↳ normalization ⟹ $z_i \mathbb{P} = \dfrac{z_i - M}{\sqrt{\sigma^2 + \xi}}$

applying scale shifitting $z_i^u = \{ z_{i_{norm}} \} + \beta$

⟹ The Two reasons for using the batch norimalization layer.

① Batch Norimalization is a general Technique That can be used to norimalize the input to alayer.

② it can be used with most network Types, Such as Multi layers Perceptrons, Convolutional and neural Network it Makes anetwork more stable during Training, it allows higher Learning rate and increase The speed at which network Train, it solves internal Convarient Shift.

**6**

image size   256 * 256

first layer 32 feature   filter size 3*3   stride 1

the first layer has the same width and hight as the original image   5 next layer 3*3   stride 2

→ because first layer has the same width and height of original image

image

∴ we will use padding.

[8] → Just like Traditional dropout, inverted dropout
Keeps some Using neurons. This is known as
The "Keep Probability". The one difference is that's own
The Training of a neural network, inverted dropout scales The
activations by the inverse of The Keep Probability P = 1.

→ Inverted Dropout is implemented in Practice
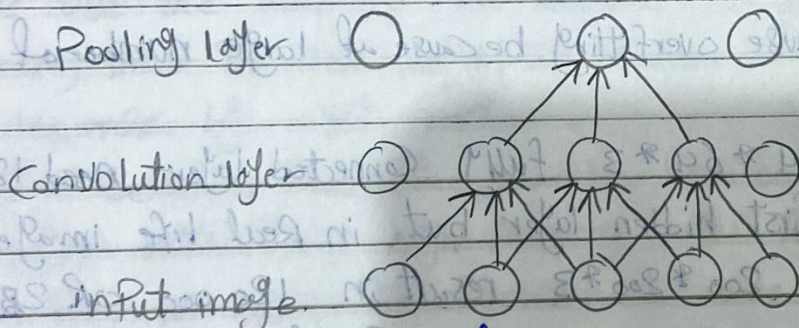in The various deep learning Frameworks because it helps.

Pooling layer

Convolution

256 * 256

when we applying the filter (3*3) we
get The Points in The Convolution output (5*5)
[receptive] and when applying The Pooling
window (3*3), we get The Points
That in The Pooling layers.

Pooling layer



Convolution layer

input image.

The size of The receptive field is **5 X 5**     or we can say 5*5*3

[7] input data block in a Convolutional network has dimension C*H*W
= 96 * 128 * 128
↳ channels    ↳ spatial dimension

applying Conv, filter = D*C * HF * WF = 128 * 26 *7 *7
stride 2, Pad 3

→ dimension of The output = $\dfrac{Size - filter + 2P}{S} + 1$

= $\dfrac{128 - 7 + (2*3)}{2} + 1 = 64.5$     i. output = 128 * 64 * 64

[8] → Just like Traditional dropout, inverted dropout randomly keeps some weights and sets others to zero. This is known as the " keep probability " P. The one difference is that, during the training of a neural network, inverted dropout scales the activations by the inverse of the keep probability $q = 1 - P$.

→ inverted Dropout is how Dropout is implemented in practice in the various deep learning frameworks because it helps to define the model once and just change a parameter to run train and test on the same model.

→ This prevents network's activations from getting too large, and does not require any changes to the network during evaluation.

[9] because it cause overfitting because of large number of parameters
for example :- $64 * 64 * 3$ fully connected layer need 12288 weights in the first hidden layer but in real life image have at least $200 * 200 * 3$ result in 120 000 or $225 * 225 * 3$ which result in 151875 weights in the first hidden layer also this layer require
• Large storage to store the weights
• Long time to train
• long time to classify an input image.
• Also when the place of class or object is changed it can't identify it.

[4]

10 • Given Two arrays A[] and B[] consisting of N and M integers respectively, The task is To construct a Convolution array C[] of size (N+M-1).

• The Convolution of 2 arrays is defined as $C[i+j] = \sum (a[i] * b[j])$ for every i and j

⇒ Our Example:
input A[] = $\{4, 1, -1, 3\}$ & B[] = $\{-2, 1\}$
Size of array , C[] = N+M-1 = 4+2-1=5 ⇒ (Length)
$C[0] = A[0] * B[0] = 4 * -2 = -8$
$C[1] = A[0] * B[1] + A[1] * B[0] = 4*1 + 1 * -2 = 2$
$C[2] = 3$ , $C[3] = -6$ & $C[4] = 3$
∴ The output = $\{-8, 2, 3, -6, 3\}$

11 
we use the learning rate adaptor to handle this Problem, every time the loss begins to Platea, The learning rate decreases by a set fraction.
⇒ The belief is that the model has become caught in region with the "high learning rate" Reducing the learning rate will allow the optimizer to more efficiently Find the minimum in the loss Surface. At this time, one might be concerned about Converging to local minimum.

12
because Conv. layer is more flexible than fully Connected because it isn't densely Connected the input doesn't affected To all output nodes. So, The number of weights Per layer is small which helps alot with high dimensional input. Such as image Processing . and they assume that also the Convolution layers are Explicit hierarchical representation

5
Nile

of features. The best thing in CNN architecture is no need for feature extraction.
- Reduces overfitting :- if the model is massively overfitting you can start adding dropout in small pieces.
- Translation invariant.

13
- The term "dropout" refers to dropping out the nodes (input and hidden layer) in a neural network. All the forward and backwards connections with a dropped node are temporarily removed, Thus creating new network architecture out of parent network. The nodes are dropped by probability of P

- in the original implementation of dropout layer, during Training a unit in layer is selected with a keep probability (1 - drop probability). This creates a Thinner architecture.?

- During the inference (test), we do not use a dropout layer, This means that all the units are considered during the prediction step. but, because of Taking all units from a layer, The final weights will be larger than expected and to deal with this is problem, weights are first scaled by the chosen dropout rate. With this, The network would be able to make accurate predictions.
⇒ To be more precise, if a unit is retained with probability P during training, the outgoing weights of that unit are multiplied by P during The Prediction Stage.

6