

## Review Questions (Assignment -4)

Name: Akshay Jain

Student ID: A20502846

Computer Vision AS4 Writeup (Review)

### Neural Networks

Sdn (a) Template matching refers to the image processing where we find similar templates in a source image by giving a base template to compare on. The process of template matching is done by comparing each of the pixel values of the source image one at a time to the template image.

Let's take a Neural Network such that:

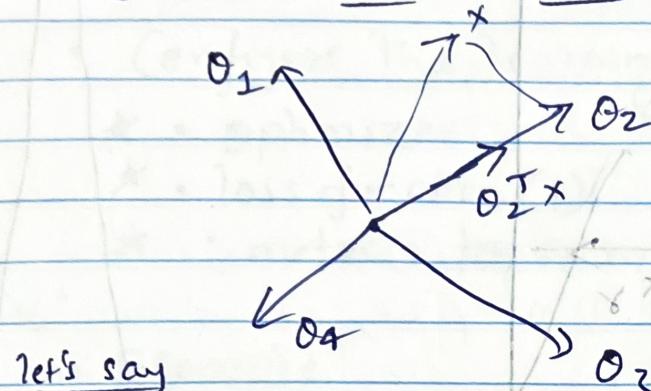


$x_1$    
input  $x_0, x_1, \dots, x_n$

& output  $\hat{y}$

$$\hat{y} = f(x_j, \theta) = \theta^T x + \theta_0 \quad \text{--- eqn ①}$$

Template Matching Interpretation can be understand using below fig. & eqn 2207 Computer vision



↳ the rows of the weight matrix  $\Omega^T$  represent templates  
↳ with  $K$  rows we have  $K$  templates (one template per class)

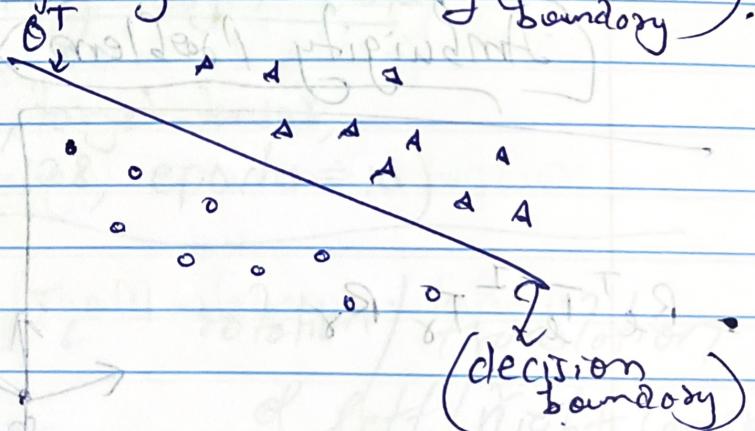
↳ the product  ~~$\Omega^T \times$~~   $\Omega^T x$  measure how well the input vector matches the template (dot product between vectors)

↳ High similarity to the template of a particular class indicates membership in this class.

### Decision Boundary Interpretation

Similarly, we can define  $\Omega^T$  for decision boundary interpretation that separates different examples from different classes by defining a boundary (or decision boundary).

Here;  $\Omega^T$  separates class 'o' & 'd' by defining a boundary.



$$\frac{S_{j,i}^{(i)}(S)}{(1)(S)}$$

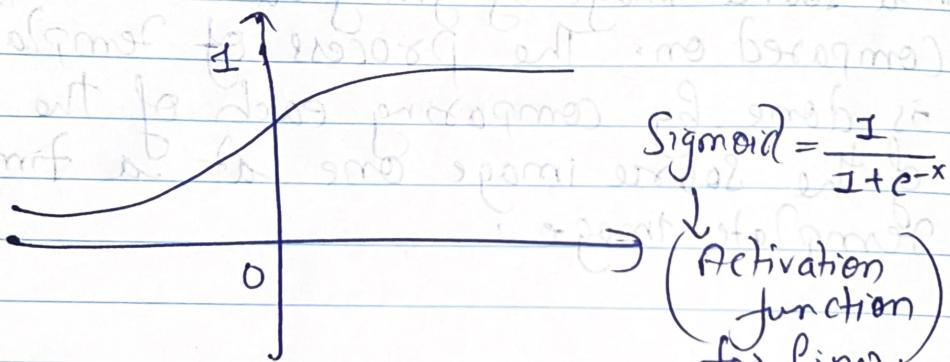
## Converting Similarity score to Probability

Given similarity scores (or decision boundary distance  $S_j^{(i)}$ ) from a linear classifier (higher better), Compute:

$$\hat{y}_j^{(i)} = P(y=j/x^{(i)})$$

↳ In a 2-class classification case:

$$\hat{y}_j^{(i)} = P(y=j/x^{(i)}) = \text{Sigmoid}(S_j^{(i)})$$



↳ In a K-class classification:

$$\hat{y}_j^{(i)} = P(y=j/x^{(i)}) = \frac{\exp(S_j^{(i)})}{\sum_{k=1}^K \exp(S_k^{(i)})}$$

$$\Theta^T \phi \rightarrow \sum_i \theta_i + x^T \phi = (\phi, \theta)^T \downarrow \text{Softmax}$$

So, Softmax is used for binary classification & Sigmoid function is used for multiclass classification.

SOL<sup>n</sup>(I) (C)

Expression of L<sub>1</sub> & L<sub>2</sub> loss :

L<sub>1</sub> & L<sub>2</sub> loss are used to minimize error b/w known & predicted values in regression problems.

$$L_1 : L_i(\theta) = \sum_{j=1}^k |\hat{y}_j^{(i)} - y_j^{(i)}|$$

$\hat{y}_j^{(i)}$  known / true       $y_j^{(i)}$  predicted

whereas;

$$L_2 : L_i(\theta) = \sum_{j=1}^k (\hat{y}_j^{(i)} - y_j^{(i)})^2$$

$y_j^{(i)}$  is predicted value &  $\underline{y_j^{(i)} \in R}$

## Huber Loss

As per wikipedia, In statistics the huber loss is a loss function used in robust regression that is less sensitive to outliers in data than the squared error loss & is given by :

$$S_\sigma(d) = \begin{cases} \frac{1}{2}d^2 & \text{if } |d| \leq \sigma \\ \sigma(|d| - \frac{1}{2}\sigma), & \text{otherwise} \end{cases}$$

$S_\sigma$  is "quadratic for small  $d$ "

"linear for large  $d$ "

with equal values & slopes of the different section at the two points where

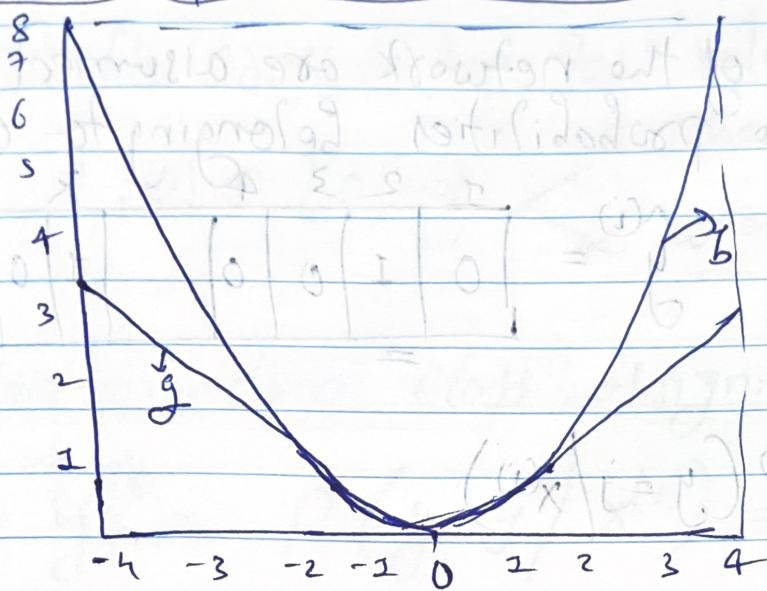
$$|d| = \sigma \quad (\text{equal values})$$

Also, variable ' $d$ ' often refers as residuals, that is to the difference b/w observed & predicted values.

$$\boxed{d = y - f(x)}$$

$\therefore$  we can also express  $\rightarrow S_\sigma(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \\ \sigma(|y - f(x)| - \frac{1}{2}\sigma) & \end{cases}$

## Graphical representation Huber loss



Huber loss ( $g, \sigma=1$ ) & squared error loss ( $b$ )  
as a function of  $(y - f(x))$

## Cross entropy loss

It is a metric used to measure how well a classification model in machine learning (deep learning) performs.

The loss (or error) is measured as a number between 0 & 1, with 0 being a perfect model.

→ In classification problems assumes that the true class label  $y^{(i)}$  is one-hot encoded.

From Class Notes

$$x^{(i)} \in C_k \Rightarrow y^{(i)} = \begin{bmatrix} 1 & 2 & 3 & \dots & k & \dots \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & \dots & 1 & \dots \end{bmatrix}$$

## Solution(1)(2)

Regularization refers to techniques that are used to calibrate machine learning / deep learning models in order to minimize the adjusted loss function & prevent overfitting or underfitting.

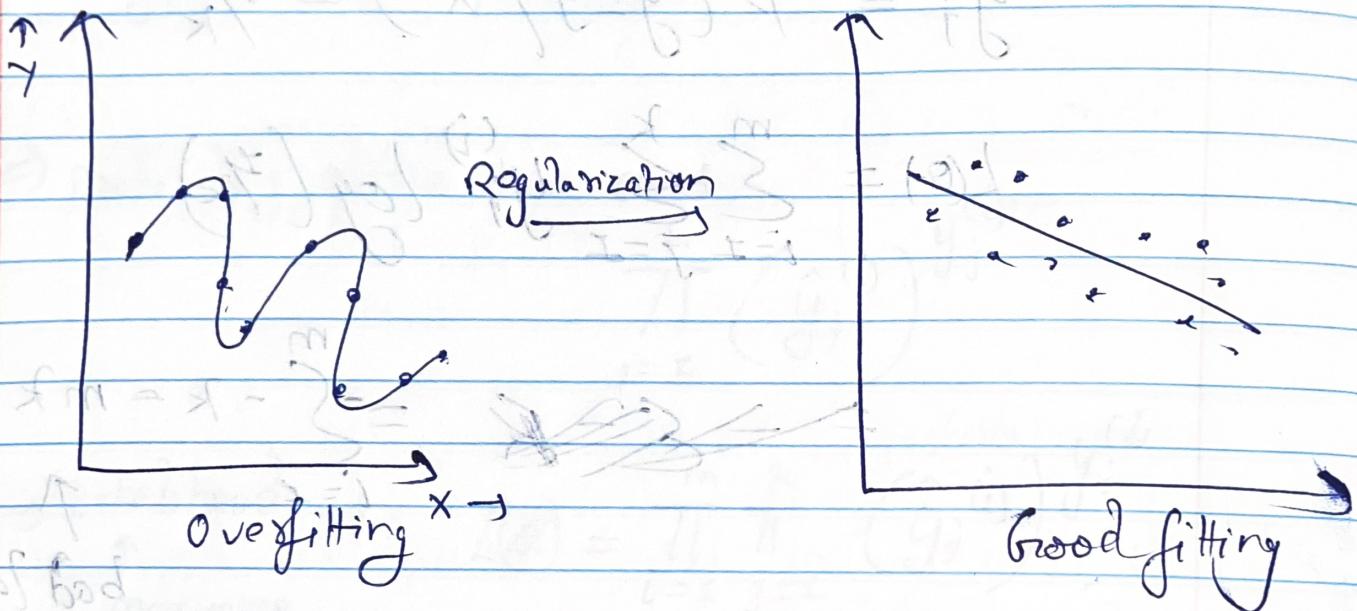
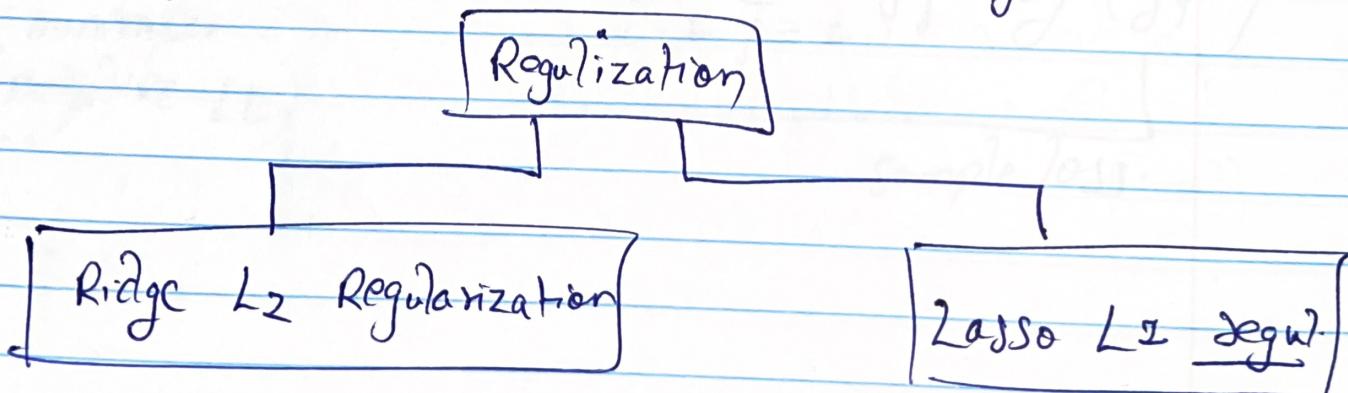


fig. Regularization on an overfitted model

These are 2 main types of regularization techniques:

- 1) Ridge Regularization
- 2) Lasso Regularization



The L<sub>2</sub> regularization is the most common type of regularization & also known as weight decay or ridge regularization.

During L<sub>2</sub> regularization, the loss function of the neural network is extended by so-called regularization term, which is called ( $J$ )

$$J(\text{Regularization term}) = \frac{1}{2} \|w\|_2^2 = \sum_{i,j} w_{ij}^2$$

$J$  is defined as Euclidean Norm of the weight matrices, which is sum over all squared weight values of a weight matrix.

The regularization term is weighted by the scalar alpha divided by two & added to the regular loss function. That is chosen for the current task.

New expression for loss function:

$$\hat{L}(w) = \frac{\alpha}{2} \|w\|_2^2 + L(w)$$

$$\hat{L}_w = \frac{\alpha}{2} \sum_{i,j} w_{ij}^2 + L(w)$$

where:  
 $\alpha$  = regularization rate

L1 regularization: In this case we simply, (Lasso) we  $\Omega$  (regularization term) as the sum of absolute values of the weight parameters in a weight matrix:

$$\Omega(w) = \|w\|_1 = \sum_i \sum_j |w_{ij}|$$

Like in L2 we can multiply regularization term by alpha & add the entire thing to the loss function.

$$L(w) = \alpha \|w\|_1 + L(w)$$

Loss function during L1 regularization

Solution (e)

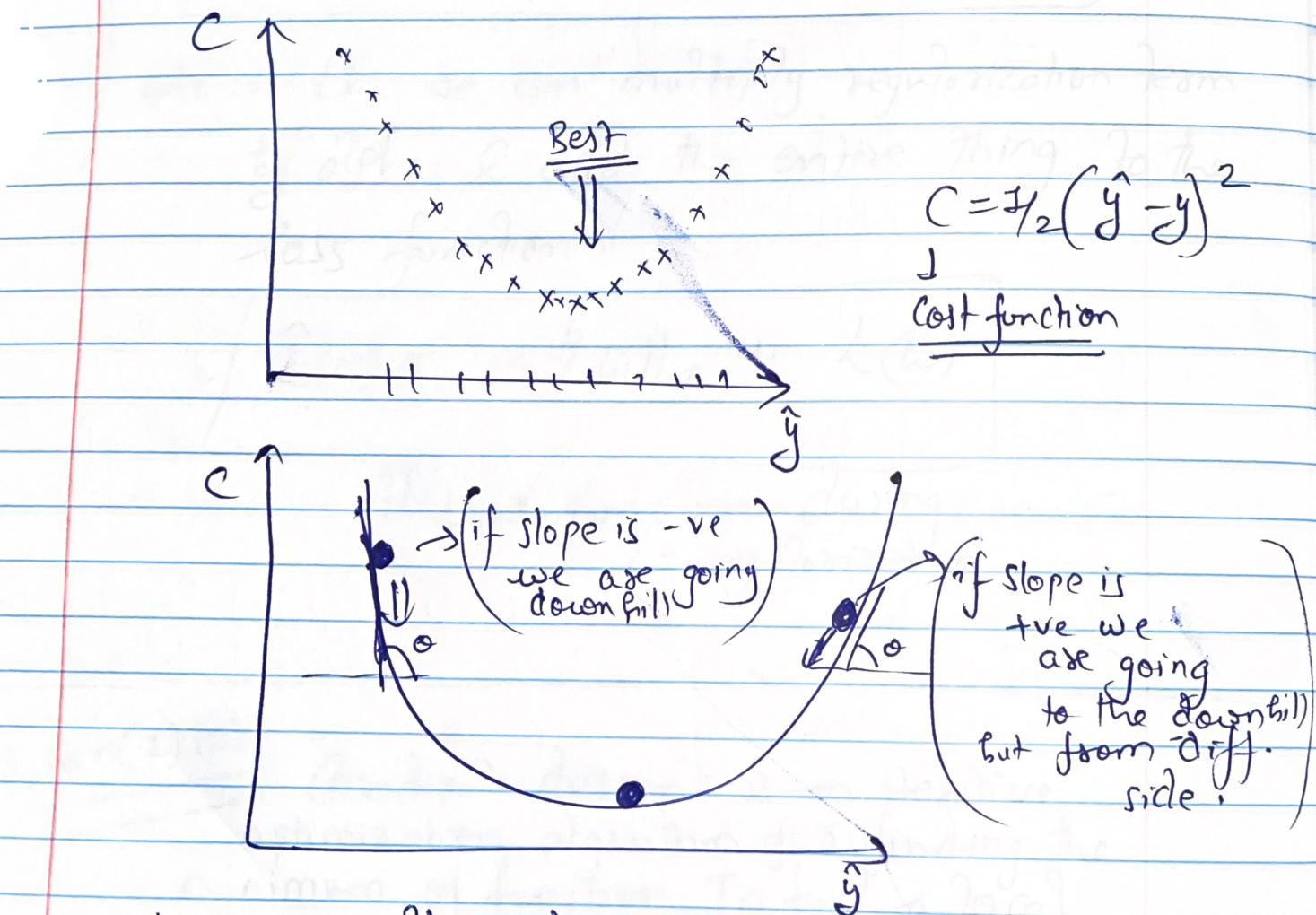
Gradient descent is an iterative optimization algorithm for finding the minimum of function. To find a local minimum of a function using gradient

descent, we take steps proportional to the negative of the gradient of the function at the current point.

In other words, we move in the direction opposite to our gradients.

## Solution (i) (f)

Gradient descent is basically the optimization done on the neural networks for minimizing the cost function. So, basically if we say there are thousands of weights that need to be adjusted & out of them which to choose so this 'Gradient Descent' approach helps to do optimization.



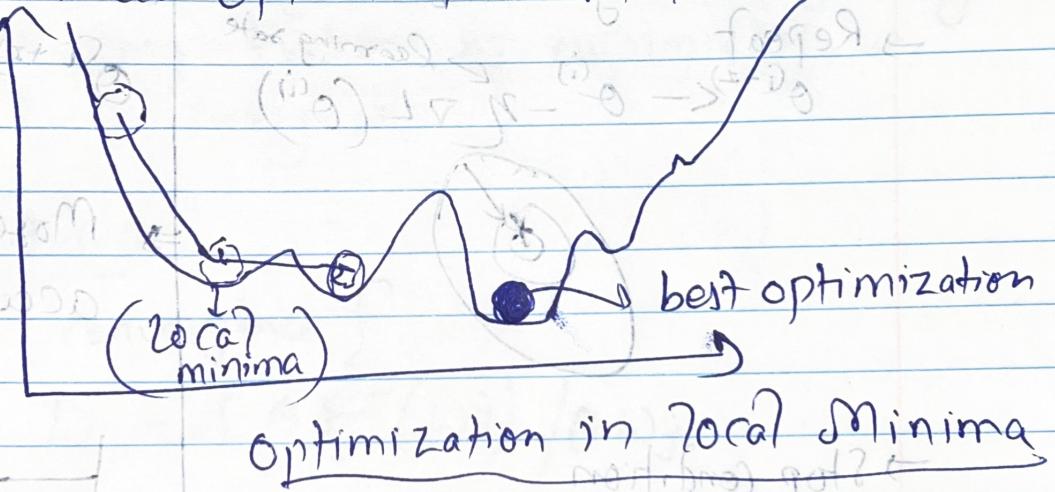
Now, we will see the slope of our cost function, which will help to do optimization.

→ Gradient Descent method)

So, using slope method we can easily find the best optimization.

### Limitations of Gradient Descent

It is good for problems having convex shape function but when the function is not convex rather hybrid it will do optimizations in local minima.



### Stochastic Gradient Descent (SGD)

It is defined by a modified gradient descent technique for doing optimization globally.

(bottom three rows handwritten)

In and all brief summary of differences b/w  
gradient descent & stochastic gradient descent  
as per class notes:

### Gradient Descent

→ Iterative optimization algorithm.

→ Start with guess  $\theta_0$ .

→ Repeat

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla L(\theta^{(i)})$$

Learning rate



→ Stop condition

$$|L(\theta^{(i+1)}) - L(\theta^{(i)})| < \epsilon$$

Local minimum

### SGD

→ Randomly order examples.

$$\rightarrow \text{For } j = 1 \dots m \quad (\text{less of } j^{\text{th}} \text{ example})$$
$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla L_j(\theta^{(i)})$$

→ More frequent but less accurate updates.

### SGD

update after every example

### mini batch

update using a subset of examples.

### GD

update based on all examples.

## Solution (I) (g)

Learning rate ( $\delta$  or  $\eta$ ) is one such hyper-parameter that defines the adjustment in the weights of our network w.r.t the loss gradient descent. It determines how fast or slow we will move towards the optimal weights.

The Gradient Descent algo. estimates the weights of the model in many iterations by minimizing a cost function at every step.

As per algorithm:

Repeat until convergence

$$w_j = w_j - \delta \frac{\partial F(w)}{\partial w_j}$$

}

where;  $\rightarrow w_j$  is the weight

$\rightarrow \theta$  is the theta

$\rightarrow F(w_j)$  is the cost function.

2 Cases

1) if  $\delta$  is too small we will need too many iterations to converge to best values!

2) if  $\delta$  is very large we will skip the optimal solution

Hence, using a good  $\delta$  is crucial.

Following approaches can be used for finding good  $\delta$ :

1) Choose a fixed learning rate.

The standard gradient descent procedure uses a fixed learning rate (eg: 0.01) that is determined by trial & error.

2) Use learning rate annealing

Learning rate annealing entails starting with a high learning rate & then gradually reducing the learning rate linearly during training. The learning rate can decrease to a value close to 0.

3) Use Cyclical learning rate.

Proposed by Leslie (2017)

Instead of monotonically decreasing the learning rate, this method lets the learning rate cyclically vary b/w reasonable boundaries.

4) Use an adaptive learning rate.

Another option is to use  $\delta$  that adopts based on the error output of the model.

Solution (1) (h)  $\Rightarrow$  Momentum is an extension to the gradient descent optimization algorithm, often referred to as Gradient Descent with Momentum.

$\Rightarrow$  It is designed to accelerate the optimization process, e.g. decrease the number of function evaluations required to reach the optima, or to improve the capability of the optimization algorithm e.g. result in a better final result.

$\Rightarrow$  A problem with GD algorithms is that the progression of the search can bounce around the search space based on the gradient. For example, the search may progress downhill towards the minima, but during this progression, it may move in another direction, even uphill, depending on the gradients of specific points (sets of parameters) encountered during the search.

This can slow down the progress of the search.

Momentum involves adding an additional hyperparameter that controls the amount of history to include in the update equation. i.e. the step to a new point in the search space.

The value of hyperparameter is in the range 0.0 to 1.0 & often close to 1.0. such as 0.8, 0.9 or 0.99.

We can break gradient descent update equation into 2 parts:

- 1) the calculation of the change to the position.
- 2) the update of the old position to the new position

→ The change in the parameters is calculated of the gradient for the point scaled by the step size.

$$\boxed{\text{change\_x} = \text{step\_size} * f'(x)}$$

New position is calculated by simply subtracting:

$$\boxed{x = x - \text{change\_x}}$$

If we think of updates over time, then the update at the current iteration or time( $t$ ) will add the change used at the previous time ( $t-1$ ) weighted by the momentum hyperparameter as follows:

$$\boxed{\text{change\_x}(t) = \text{step\_size} * f'(x(t-1)) + \text{momentum} * \text{change\_x}(t-1)}$$

The update to the position:

$$\boxed{x(t) = x(t-1) - \text{change\_x}(t)}$$

Note: The change in the position accumulates magnitude & dir<sup>n</sup> of changes over the iterations of the search, proportional to the size of momentum hyperparameter.

## Neural Network

Solution 1) e) i)

Formulae :-

Forward Propagation: As the name suggests, the input data is fed in the forward direction through the network. Each hidden layer accepts the input data, processes it as per the activation function & passes to the successive layer.

In order to generate some output, the input data should be fed in forward direction only.

Thus

At each neuron in hidden or output layers, the processing happens in 2 steps:

i) Pre-activation

ii) Activation.

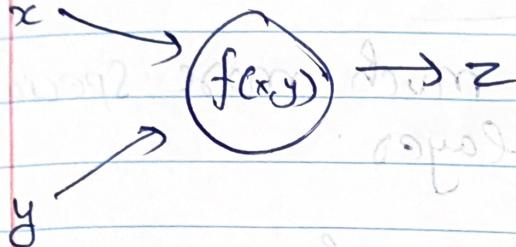
## Backward Propagation:

It is the practice of fine ~~training~~ tuning the weights of a neural net based on the error rate (i.e.) loss obtained in the previous epoch (i.e. iteration).

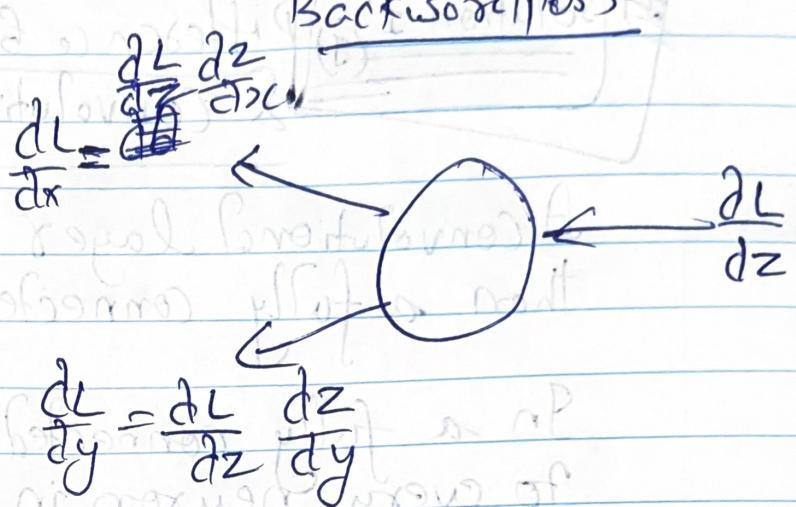
Proper training of weights ensures lower error rates, making the model reliable.

## Solution 1)

### Forward pass



### Backward pass



### Overall steps : (of propagation)

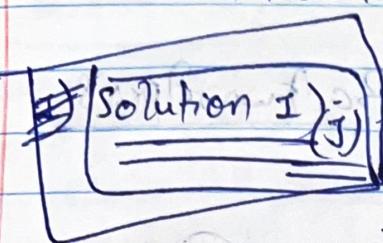
- In forward propagate step, the data flows through the network. To get the outputs.
- the loss function is used to calculate total error
- \* → then we use backward propagation algo. to calculate The gradient of the loss function with respect to each weight & bias.
- \* → Finally, we use gradient descent to update weights & biases at each layer.
- We repeat each step to minimize total error of the neural network.

88

## Review Question

1. Neural Network.

(Continued)



Pifference between fully connected & convolution layers:

A convolutional layer is much more specialized than a fully connected layer.

In a fully connected layer each neuron is connected to every neuron in the previous layer, & each connection has its own weight.

It is totally general purpose connection pattern & makes no assumptions about the features in the data.  
Also, very expensive in memory & computation.

In contrast, in convolutional layers each neuron is only connected to a few nearby neurons.

In the previous layer, the same set of weights is used for every neuron.

The fewer number of connections & weights make convolution layers relatively cheap in terms of memory & compute power needed.

Solution 1) (K)

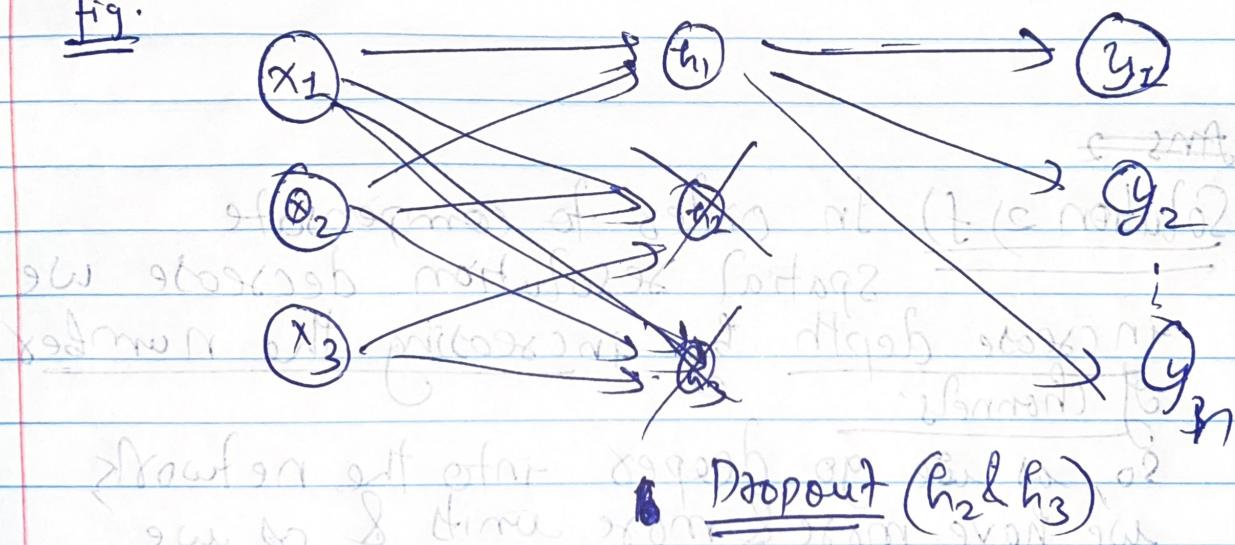
Dropout layer is used to

regularize deep neural networks,

→ It is a technique which is used to prevent overfitting.

→ Dropout works by randomly setting the outgoing edges of hidden unit to 0 at each update of the training phase.

fig.



In above fig: we have dropout h<sub>2</sub> & h<sub>3</sub> hence the hidden layer will only consist of h<sub>1</sub>.  
for evaluation of output ( $y_i$ ).

## Review Question 2

### Convolution Neural Network.

Solution 2) a)

RGB Image 4x4

Channel R

1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

3x3 filter

1	1	1
1	1	1
1	1	1

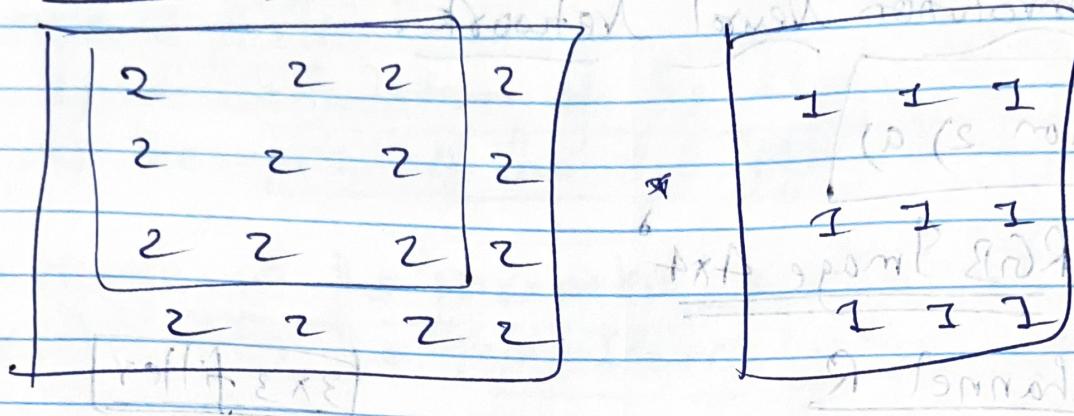
$\Rightarrow$  first element of  $4 \times 4$  matrix on convolution  
 $= (1 \times 1 + 1 \times 2 + 1 \times 2) \times 3$

$$= 9$$

$\Rightarrow$  Similarly, all the elements as all the elements are ~~the~~ same so output will also remain same = 9

9	9
9	9

## Channel G



$$= (2+2+2) \cdot 3 \\ = 18$$

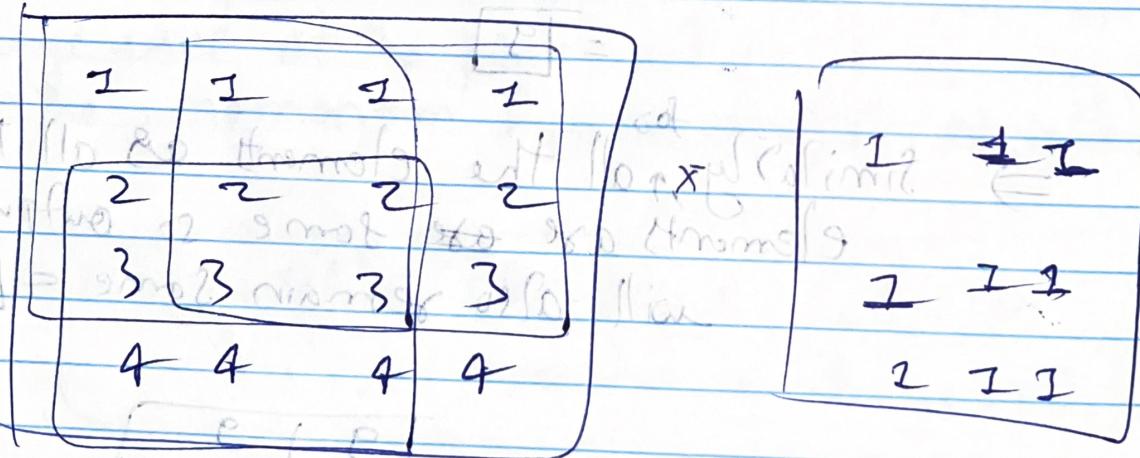
Similarly, other elements on convolution will be = 18

$\therefore$  Matrix on Convolution  $\Rightarrow$ 

18	18
18	18

 $\frac{2 \times 2}{=}$

## Channel B



First element =  $\begin{pmatrix} 1+1+1+ \\ 2+2+2+ \\ 3+3+3 \end{pmatrix} = 3+6+9 \\ = 18$

Second element = 18

$$\text{Third element} = \begin{pmatrix} 2+2+2+ \\ 3+3+3+ \\ 4+4+4 \end{pmatrix} = (6+9+12) \\ = 27$$

$$\text{Fourth element} = 27$$

So, matrix will be :

$$\begin{array}{|c|c|} \hline 18 & 18 \\ \hline 27 & 27 \\ \hline \end{array}$$

$2 \times 2$

So, we have

$$\begin{array}{|c|c|} \hline g & g \\ \hline g & g \\ \hline \end{array}$$

$$\begin{array}{|c|c|} \hline 18 & 18 \\ \hline 18 & 18 \\ \hline \end{array}$$

$$\begin{array}{|c|c|} \hline 18 & 18 \\ \hline 27 & 27 \\ \hline \end{array}$$

RGB channel  
convolution  
matrix

Element wise addition of all 3 channel :

$$\begin{bmatrix} g+18+18 & g+18+18 \\ g+18+27 & g+18+27 \end{bmatrix} = \begin{bmatrix} 45 & 45 \\ 54 & 54 \end{bmatrix}$$

On averaging each element with  $1/g$  while convolution elements will be :

$$\begin{bmatrix} 1+2+2 & 1+2+2 \\ 1+2+3 & 1+2+3 \end{bmatrix} = \begin{bmatrix} 5 & 5 \\ 6 & 6 \end{bmatrix}$$

Solution 2) b) with zero padding

<u>R</u>	→	<table border="1"> <tr><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td><td>0</td></tr> <tr><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td><td>0</td></tr> <tr><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td><td>0</td></tr> <tr><td>0</td><td>1</td><td>1</td><td>1</td><td>1</td><td>0</td></tr> <tr><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> </table>	0	0	0	0	0	0	0	1	1	1	1	0	0	1	1	1	1	0	0	1	1	1	1	0	0	1	1	1	1	0	0	0	0	0	0	0
0	0	0	0	0	0																																	
0	1	1	1	1	0																																	
0	1	1	1	1	0																																	
0	1	1	1	1	0																																	
0	1	1	1	1	0																																	
0	0	0	0	0	0																																	

3x3	<table border="1"> <tr><td>1</td><td>1</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>1</td></tr> </table>	1	1	1	1	1	1	1	1	1
1	1	1								
1	1	1								
1	1	1								

Berechnung  
mit Formeln  
X ist dann

$$= \begin{array}{|c|c|c|} \hline 4 & 0+6 & 0+6 \\ \hline 6 & 9 & 9 \\ \hline 6 & 9 & 9 \\ \hline 4 & 6 & 4 \\ \hline \end{array}$$

$$= \begin{array}{|c|c|c|} \hline 4 & 0+6 & 0+6 \\ \hline 6 & 9 & 9 \\ \hline 6 & 9 & 9 \\ \hline 4 & 6 & 4 \\ \hline \end{array}$$

: Formeln & No für mittlere Werte freihalten

$$\underline{\text{2nd element}} = \begin{array}{|c|c|c|} \hline 0 & 0+0 & 8+8+8 \\ \hline 1 & 1+1 & 8+8+8 \\ \hline 1 & 1 & 1 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

slide 10: die Formeln für die Mittleren Werte  
= (1. Element mit Formeln)

4th row or 13th element

$$= \begin{array}{|c|c|c|} \hline 0 & 1 & 1 \\ \hline 0 & 1 & 1 \\ \hline 0 & 0 & 0 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$$= 4$$

14<sup>th</sup> element

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\times \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} = 6$$

4<sup>th</sup> element

$$\begin{bmatrix} 0 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}$$

$$\times \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = 4$$

~~1x3 = 3x3 = 9 elements~~

~~81 - 2x2 = 6x6 = 36 elements~~

~~(8) = 4 elements~~

5<sup>th</sup> element

$$\begin{bmatrix} 0 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix}$$

$$\times \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} = 6$$

∴ final Convolved matrix with zero padding

$$\begin{bmatrix} 4 & 6 & 6 & 4 \\ 6 & 9 & 9 & 6 \\ 6 & 9 & 9 & 6 \\ 4 & 6 & 6 & 4 \end{bmatrix}$$

Channel  $\otimes G$

0	0	0	0	0	0
0	2	2	2	2	0
0	2	2	2	2	0
0	2	2	2	2	0
0	2	2	2	2	0
0	0	0	0	0	0

0	0	1	1	1
*	1	1	1	1
1	1	1	1	1

$$\begin{aligned} \text{First element} &= 2 \times 4 = 8 \\ \text{2nd \& 3rd} &= 2 \times 6 = 12 \\ \text{4th element} &= 8 \end{aligned}$$

$$\begin{aligned} s_m &= 2 \times 6 \\ g_m &= 2 \times 9 = 18 \end{aligned}$$

$$\text{Convolved } G \text{ matrix} = \begin{bmatrix} 8 & 12 & -12 & 8 \\ 12 & 18 & 18 & 12 \\ -12 & 18 & 18 & 12 \\ 8 & 12 & 12 & 8 \end{bmatrix}$$

channel B

0	0	0	0	0	0
0	1	2	1	1	0
0	2	2	2	2	0
0	3	3	3	3	0
0	4	4	4	4	0
0	0	0	0	0	0

1	2	1
3	3	1
1	2	2

$$1^{\text{st}} \text{ element} = 1 \times 2 + 2 \times 2 = 6$$

$$2^{\text{nd}} / 3^{\text{rd}} = 1 \times 3 + 2 \times 3 = 9$$

$$4^{\text{th}} = 1 \times 2 + 2 \times 2 + 3 \times 2 = 2 + 4 + 6 = 12$$

$$5^{\text{th}} = 1 \times 3 + 2 \times 3 + 3 \times 3 = 3 + 6 + 9 = 18$$

$$6^{\text{th}} = 1 \times 2 + 2 \times 2 + 3 \times 2$$

$$7^{\text{th}} = 2 \times 2 + 3 \times 2 + 4 \times 2$$

$$= 4 + 6 + 8 = 18$$

$$8^{\text{th}} = 2 \times 3 + 3 \times 3 + 4 \times 3 = 27$$

$$13^{\text{th}} = 6 + 8$$

$$= 14$$

$$14^{\text{th}} = 9 + 12$$

$$= 21$$

$$= \begin{bmatrix} 6 & 9 & 9 & 6 \\ 12 & 18 & 18 & 12 \\ 18 & 27 & 27 & 18 \\ 14 & 21 & 21 & 14 \end{bmatrix}$$

Element wise add<sup>n</sup> of all channel

$$= \begin{pmatrix} 4 & 6 & 6 & 4 \\ 6 & 9 & 9 & 6 \\ 6 & 9 & 9 & 6 \\ 4 & 6 & 6 & 4 \end{pmatrix} + \begin{pmatrix} 8 & 12 & 12 & 8 \\ 12 & 18 & 18 & 12 \\ 12 & 18 & 18 & 12 \\ 8 & 12 & 12 & 8 \end{pmatrix} +$$

$$+ \begin{pmatrix} 6 & 9 & 9 & 6 \\ 12 & 18 & 18 & 12 \\ 18 & 27 & 27 & 18 \\ 14 & 21 & 21 & 14 \end{pmatrix} = \begin{pmatrix} 18 & 27 & 27 & 18 \\ 30 & 45 & 45 & 30 \\ 30 & 54 & 54 & 30 \\ 26 & 39 & 39 & 26 \end{pmatrix}$$

Matrix after  
Convolution =  $\begin{bmatrix} 18 & 27 & 27 & 18 \\ 30 & 45 & 45 & 30 \\ 30 & 54 & 54 & 30 \\ 26 & 39 & 39 & 26 \end{bmatrix}$

Solution 2) c) Dilated Kernel (with  $dak = 2$ )

$$K = \begin{bmatrix} 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \end{bmatrix} \quad 5 \times 5$$

R. channel (Convolution)

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \quad 6 \times 6$$

$$\begin{bmatrix} 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \end{bmatrix} \quad 5 \times 5$$

$$I^{th} = 0 + 0 + (1+1) + 0 + (1+1) = 4$$

similarly,  
 $2^{nd}/3^{rd}/4^{th}$  elements.

~~Iteration 2)  $\hat{I}(i)$~~  = 
$$\begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix}_{2 \times 2}$$
 (feature map from R channel)

G channel

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 2 & 2 & 2 & 0 \\ 0 & 2 & 2 & 2 & 2 & 0 \\ 0 & 2 & 2 & 2 & 2 & 0 \\ 0 & 2 & 2 & 2 & 2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} * K =$$

~~If modified p[1,1] = (2+2) + (2+2) = 8~~

Similarly other elements.

if fixing feature so, convolved matrix (feature map) for G

$$\begin{bmatrix} 8 & 8 \\ 8 & 8 \end{bmatrix}_{2 \times 2}$$

B channel

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 2 & 2 & 2 & 2 & 0 \\ 0 & 3 & 3 & 3 & 3 & 0 \\ 0 & 4 & 4 & 4 & 4 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} * K = \begin{bmatrix} 12 & 12 \\ 8 & 8 \end{bmatrix}_{2 \times 2}$$

Element wise addn of all channel

$$= \begin{bmatrix} 4+8+12 & 4+8+12 \\ 4+8+8 & 4+8+8 \end{bmatrix}$$

$$= \begin{bmatrix} 24 & 24 \\ 20 & 20 \end{bmatrix}$$

Solution 2) d) Template matching is the process of moving the template over the entire image & calculating the similarity between the template & the covered window on the image.

It is implemented through 2D convolution.  
In convolution, the value of output pixel is computed by multiplying elements of two matrices & summing the results.

One of these matrices represents the image itself, while the other matrix is the template, which is known as convolution kernel.

$$\begin{bmatrix} s1 & s1 \\ 8 & 8 \end{bmatrix}$$

$$= 1 \times$$

$$\begin{bmatrix} 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Solution 2) e) MultiScale Analysis with a  
fixed window size

We can achieve multiple scale analysis by increasing the depth of the image. So, by this way same region in an image with a fixed window size can store more information at different depths. Hence, multiscale achieved.

Ans 2

Solution 2) f) In order to compensate spatial resolution decrease we increase depth by increasing the number of channels.

So, as we go deeper into the network, we have more & more units & as we know deeper layers have larger receptive fields & corresponds to more specific features & hence we can detect complex spatial shapes. That's is the major purpose for doing so.

~~Solution 2) g)~~

Formula to calculate feature map shape

$$f_m = \left[ \frac{i - k + 2p}{s} \right] + 1$$

↓  
stride

$$f_m = \left[ \frac{128 - 3 + 0}{2} \right] + 1 = 126$$

So, the size of resulting tensor for 16 convolution filter of size 3x3x32:

$$\Rightarrow [126 \times 126 \times 16]$$

Solution 2) h) (with stride=2)

$$f_m = \left[ \frac{128 - 3}{2} \right] + 1 = 64$$

$$= (128/2) + 1 = 62 + 1$$

$$f_m = 63$$

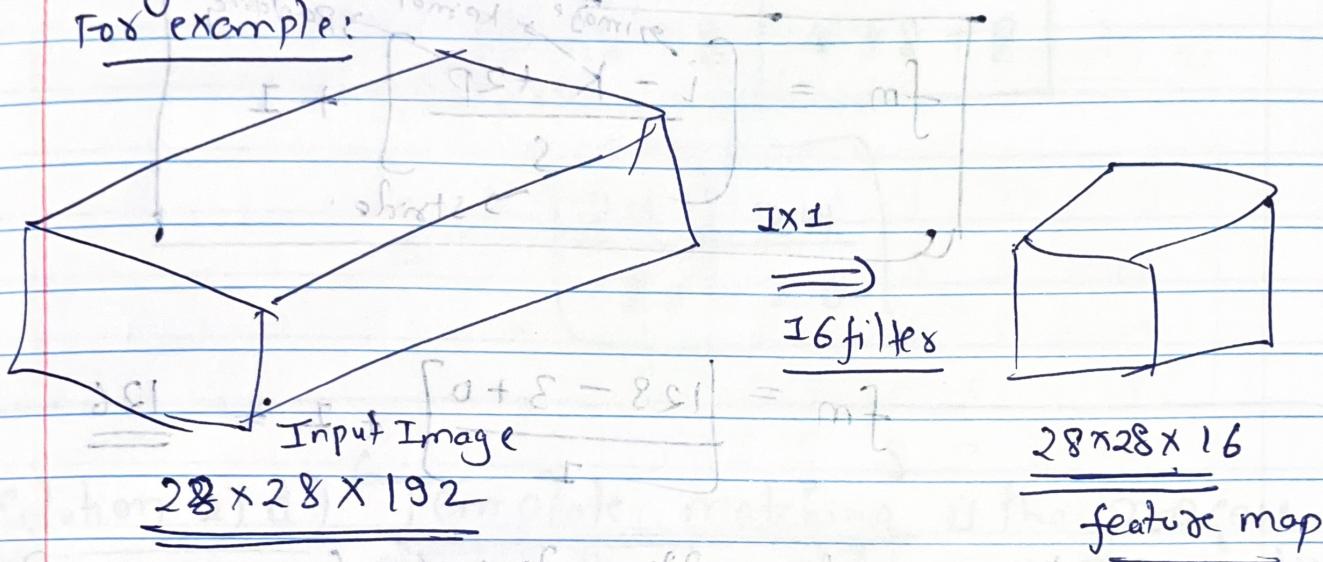
So, the size of resulting tensor:

$$\Rightarrow [63 \times 63 \times 16]$$

Solution 2) i)  $\boxed{1 \times 1}$  Convolution operation involves convolving the input with the filter

of size  $1 \times 1$ .

For example:



So, as we can see, in above operation using  $1 \times 1$  filters (16) the depth of input image

reduced / shrinked from 192 to 16.

So,  $1 \times 1$  is used for Dimension Reduction & in making computation less expensive.

$$I + S - 2S + 1 = I + (1 \times 1) =$$

$$\boxed{1} = mI$$

: cannot write  $1 \times 1 \times 1$  as  $1$

$$\boxed{1 \times 1 \times 1} \in$$

(E) (F) (I)

Solution 2) j) A convolutional layer, generally used with ConvNet, is a class of neural networks that specializes in processing data that has a grid like topology, such as an image.

Convolution layers is the core building block of CNN. It performs a dot product b/w 2 matrices, where one matrix is the set of learnable parameters known as kernel & other position is restricted portion of the receptive field.

Difference between early & deep layers:

Major difference is that early convolution layers can learn features like lines, edges etc. whereas;

Deep layers can do for more complex features like face detection, object recognition etc.

Solution 2) K)

Purpose of pooling

Pooling = down sampling (spatial dimensions)  
(depth unchanged)

- ① So, the main purpose of is it helps in deducing the number of network coefficients.
- ② Also, Pooling supports multiple scale analysis.
- ③ It makes the operations / convolutions computationally less expensive.

The most common approach used in pooling is max pooling.

Solution 2) L)

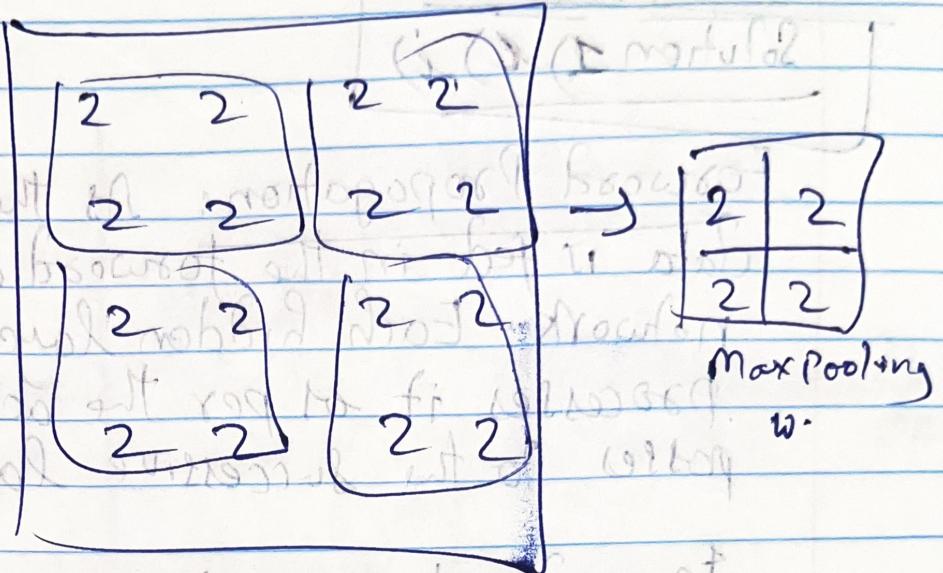
R Channel

1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

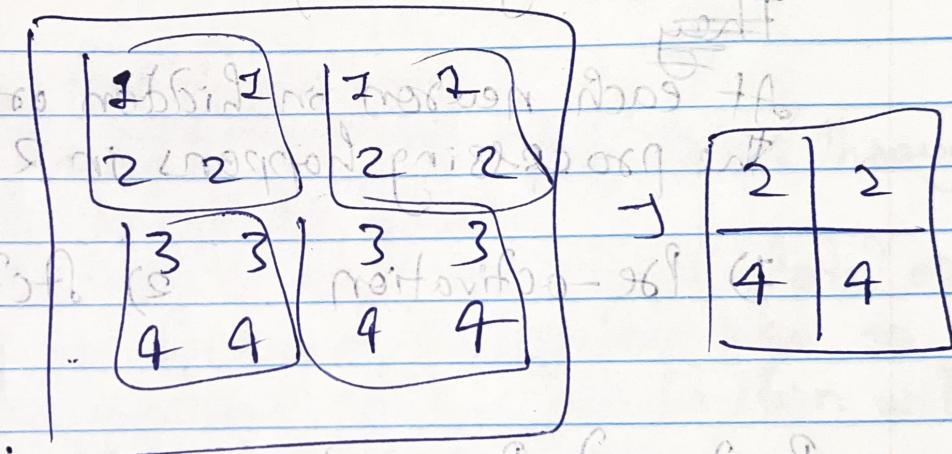
Max Pooling  
with stride = 2

1	1
1	1

G channel



B channel rich Ground truth 2d bands



cropping against Ground truth