

Kinjal Kachi

Assignment - 3

Computer Vision
fall 2020

ANumber A20449343

Q1)

Computer Vision

i) Neural Networks

q) Template matching interpretation.

- The linear classifier is given as the function of weights and input as :-

$$\hat{y} = f(x, \Theta) = \underbrace{\Theta^T x}_{k \times n} + \underbrace{\Theta_0}_{n \times 1}$$

- Template matching is used to find similarity between the feature vector and templates

- Rows at Θ^T are templates. If we have k different classes we have k templates representation as $\Theta^T_{k \times n}$.

- $\Theta^T x$ measures how well x matches with the k templates by performing dot product.

- The higher the values of $\Theta^T x$ with a template from all k -templates, higher the membership in k^{th} class.

The matrices are shown as follows :

$$\hat{y} = \begin{bmatrix} n_0 \\ \vdots \\ n_n \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} d_0 \\ \vdots \\ d_k \end{bmatrix} = \begin{bmatrix} n_1 \cdot x \\ \vdots \\ n_k \cdot x \end{bmatrix} + \begin{bmatrix} d_0 \\ \vdots \\ d_k \end{bmatrix}$$

Decision Boundary Interpretation.

- Here we calculate the distance from the decision boundary to classify the feature.

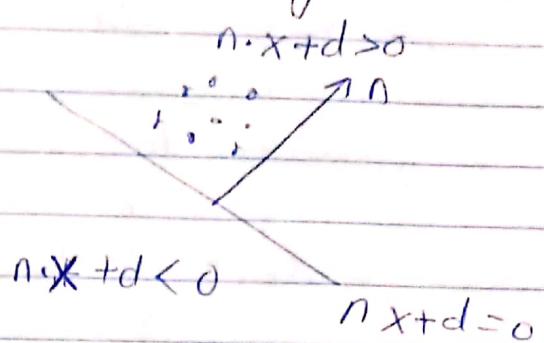
- Rows of Θ^T are parameters of k linear discriminant

function

- Each linear discriminant separates one class from "all others"

- Positive values for the function says that the examples belongs to that class and negative value suggests that it doesn't

$$\hat{y} = \begin{bmatrix} n, x \\ \vdots \\ n_k, x \end{bmatrix} + \begin{bmatrix} d_0 \\ \vdots \\ d_k \end{bmatrix}$$



n = normal

d = distance from origin

(b) Converting similarity score to probability

- Two class classification.

Here we use sigmoid function to classify the example into one of the binary class. It takes the features as input & outputs the probability

$$\hat{y}_j^{(i)} = P(y=j | x^{(i)}) = \text{sigmoid}(s_j^{(i)})$$

$$\text{sigmoid}(x) = \frac{1}{1+e^{-x}}$$

Here use softmax activation function.

$$\hat{y}_j^{(i)} = P(y=j | x^{(i)}) = \frac{\exp(s_j^{(i)})}{\sum_{l=1}^k \exp(s_l^{(i)})}$$

(c) Loss functions:-

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

where,

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

Huber loss:

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

o Cross Entropy loss:

$$L(\theta) = -\sum_{i=1}^m \sum_{j=1}^k y_j^{(i)} \log(\hat{y}_j^{(i)})$$

d) Regularization:

We add regularization to the loss function in order to get a simple solution i.e. lower weights. It helps the model generalize better.

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m L_i(f(x_i; \omega), y_i) + \lambda R(\theta)$$

Here, $R(\theta)$ is regularization term

λ = weight of regularization
(hyperparameter)

(e) In gradient descent, we try to find the minima of the function.

We take the gradient to decide the direction in which we need to move forward.

- The gradient of the function point into the direction where the function changes the most.
- In order to find minima, we move in opposite direction of the gradient

- (f) Gradient Descent can be slow when there are lot of parameters and examples in data
- It updates the new parameters by looking at each example for every iteration
 - Update is more frequent in gradient descent.
 - On the other hand, stochastic gradient descent (SGD) uses different approaches.
 - It looks at each eg. and updates the parameters on the go.
 - SGD takes less time to converge than gradient descent.
 - The problem with SGD is that as it is taking each example & updating parameters, it would not converge in the right values.

SGD

Update after
every example

Minibatch

Update after each
minibatch

GD

Update after
looking at all
examples

(g) learning rate:

- learning rate values are hyperparameters for gradient descent & can be made smaller as iteration progresses. It is called 'learning rate decay'. We can use following approaches.

(i) Step decay.

With every iterations k , we reduce the learning rate by some specified function.

(2) Exponential decay

The learning rate is decayed exponentially with each iteration.

$$n_t = n_0 e^{-kt}$$

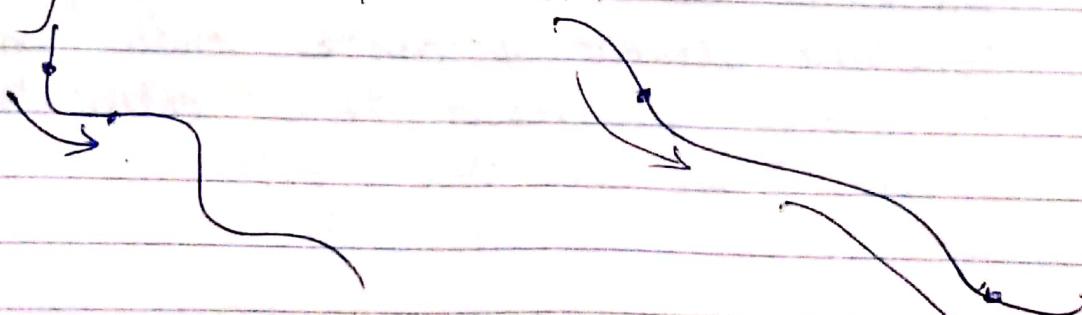
(3) Fractional decay

$$n_t = n_0 / (1 + kt)$$

With each iteration, reduce the learning rate by some fraction.

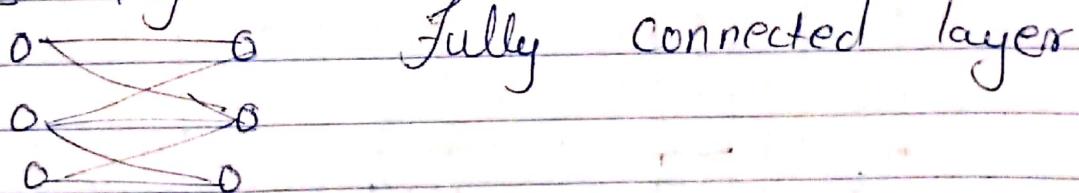
(ii) Momentum.

- We can add momentum on top of gradient descent to obtain better results.
- If the momentum is not added, GD could get stuck at the local minimum of coverage at that point.
- Adding momentum helps the algorithm understand to keep moving in the direction of descent even if the optimal values according to the equation are found.
- This phenomenon is similar to a ball rolling down the hill.

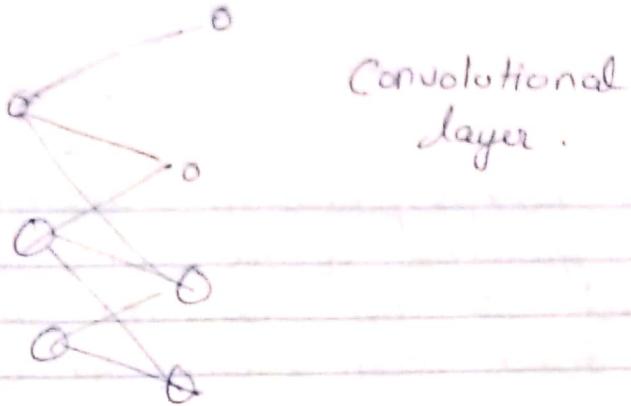


- (ii) The backpropagation algorithm calculated the gradients of the loss function with respect to the network's weights.
- We start with some initial guess for the weights and then update these weights using gradient descent for each node.
 - Forward Pass
In forward pass, we simply push the input to compute all the intermediate values at the nodes. By doing this, we have all the required values to calculate the gradients using chain rule for each node.
 - Backward Pass
In backward pass we push the gradient towards the beginning starting from the end node. We calculate the derivative at each node wrt loss using chain rule. In the backward pass we update the weights using gradients to minimize the loss function.
 - In backward pass, gradients are propagated towards the beginning from the end node in order to update the networks node parameters.

(j) In a fully connected layer, each neuron receives input from every element of the previous layer.



In a convolutional layer, neuron ~~re~~ receives input from only a restricted subarea of previous layer.



(k) Dropout

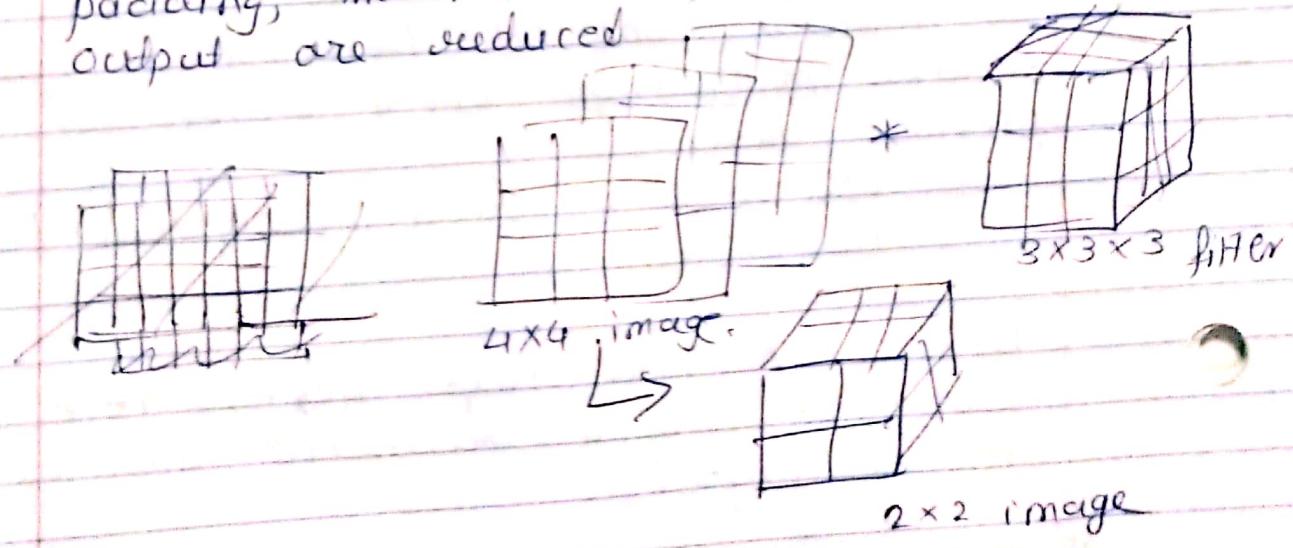
- Dropout means we remove some units in the layers of neural network during training randomly.
- At each training stage, individual nodes are either dropping out of the network with probability $1-p$, so that the reduced network is left.
- We use dropout in order to prevent overfitting & avoid co-dependency among neurons, improve training speed.

q2] Convolutional Neural Networks

(a) I is a 4×4 RGB image

filter is $3 \times 3 \times 3$ matrix with all 1s

- Convolute without zero padding.
When we apply the filter without zero padding, the number of dimensions of output are reduced



As we convolve the filter with image, we get $2 \times 2 \times 3$ matrix as output. 2×2 is result at filter with the respective channel of the input image.

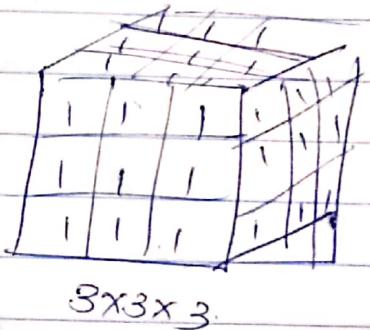
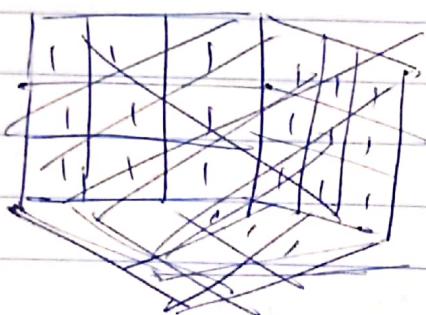
$$R = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

$$G = \begin{bmatrix} 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \end{bmatrix}$$

$$B = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \\ 4 & 4 & 4 & 4 \end{bmatrix}$$

All Is

filter



$3 \times 3 \times 3$

The results for each channel, RGB are.

	R	G	B	final
(1,1)	9	18	18	45
(1,2)	9	18	18	45
(2,1)	9	18	27	54
(2,2)	9	18	27	54

$$\begin{bmatrix} 9 & 9 \\ 9 & 9 \end{bmatrix} + \begin{bmatrix} 18 & 18 \\ 18 & 18 \end{bmatrix} + \begin{bmatrix} 18 & 18 \\ 27 & 27 \end{bmatrix} \Rightarrow \begin{bmatrix} 45 & 45 \\ 54 & 54 \end{bmatrix}$$

Output w/o zero padding

(b) Convolution with zero padding for R channel looks like.

R	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

Other channels (G and B) will look similar with shape 6×6 .

Similarly apply the filter as we did in Q 2(a), we get results for each channel as follows

$$R = \begin{array}{|c|c|c|c|} \hline 4 & 6 & 6 & 4 \\ \hline 6 & 9 & 9 & 6 \\ \hline 6 & 9 & 9 & 6 \\ \hline 4 & 6 & 6 & 4 \\ \hline \end{array} + G = \begin{array}{|c|c|c|c|} \hline 8 & 12 & 12 & 8 \\ \hline 12 & 18 & 18 & 12 \\ \hline 12 & 18 & 18 & 12 \\ \hline 14 & 21 & 21 & 14 \\ \hline \end{array}$$

Adding up values of each channel to get final output

$$\begin{bmatrix} 18 & 27 & 27 & 18 \\ 30 & 45 & 45 & 30 \\ 36 & 54 & 54 & 36 \\ 26 & 39 & 39 & 26 \end{bmatrix}$$

Output with p zero padding

(c) Atrous convolution with dilation rate of 2.
Let's do this for each channel & then add them to get final output

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

when convolving with filter with all ones,
we get,

$$\text{Iteration 1: } (0x1) + (0x1) + (0x1) + \\ (0x1) + (1x1) + (1x1) + \\ (0x1) + (1x1) + (1x1) = 4$$

$$\text{Iteration 2: } 4$$

$$\text{Iteration 3: } 4$$

$$\text{Iteration 4: } 4$$

Similarly :-

$$R = \begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix} \quad G =$$

$$G = \begin{bmatrix} 8 & 8 \\ 8 & 8 \end{bmatrix}$$

$$B = \begin{bmatrix} 12 & 12 \\ 8 & 8 \end{bmatrix}$$

$$\text{final output} = R + G + B$$

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

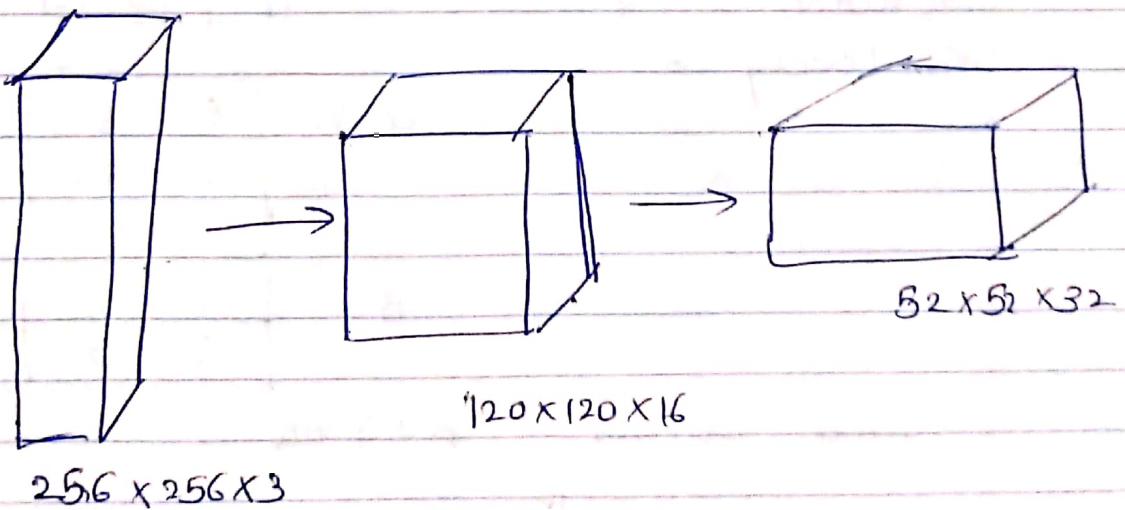
Q2d) Template matching is the technique used to find small parts of an image which matches template image. It slides the window across image that will provide a percent match with the template. When image is convoluted with the filter, it is expected to give portions of the image that match the particular template to give a higher response.

Q2 e) Multi-scale analysis can be achieved with a fixed window size by reducing an image size by continuously convolving it with filter. This will lead to a larger receptive field in the convoluted layer., better feature extraction and accounts for differently scaled objects.
eg: We start with 8×8 size image and fixed window size of 2×2 . After first convolution, we get 4×4 output, running another convolution gives an image of 2×2 . As the scale of image reduces the objects in the image get more receptive area and would actually lead to feature extraction.

(f)

When applying convolution to an image, the size of the image is reduced the spatial dimension due to pooling or stride

- If we set the spatial dimension reduces we lose information from the image.
- In order to do so, we increase the depth of the image to compensate for spatial resolution decrease.
- This is done by applying number of filters to get the depth of n.
- As the depth increases, we have more no. of coefficients which represent the information in the image at each channel.
- Each filter we apply store a specific information in the image & these n filters stacked as whole give the required information with reduced spatial dimension



(g) Input = $128 \times 128 \times 32$ No. of filters = 16
filter size = $3 \times 3 \times 3$

The output size of tensor O after convolution is given by:-

$$O = \frac{I - f + 2P}{S} + 1$$

I = Input size of input

f = filter size

P = Padding

S = stride

$$\therefore O = \frac{128 - 3 + 2(0)}{2} + 1 = 126$$

∴ Size of output = $126 \times 126 \times 16$
 \therefore (As we use 16 filters)

(h) With stride of 2, $O = \frac{128 - 3 + 2(0)}{2} + 1$
 $= 63.5$

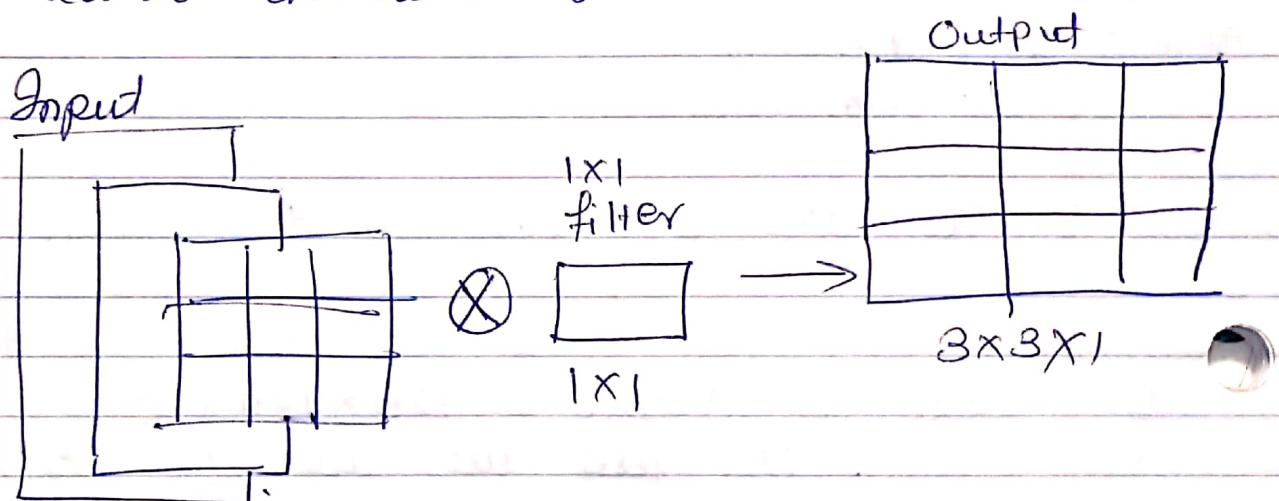
Size of output tensor = $64 \times 64 \times 32$

(i) 1×1 convolution

We apply 1×1 convolution to reduce the number of channels in the image.

A 1×1 filter will only have a single parameter of weight for each channel in the input and like the application of any filter results a single output value.

- This allows 1×1 filter to act like a neuron with an input from the same position across each of the feature maps in the input.
- A single neuron can be convolved along the image resulting in a feature map with same width & height as input. & reduced channel to 1.



(k) Pooling

- Pooling layer is added in the neural network to downsample spatial dimensions of the input from previous layer.
- The pix pooling operator is always almost of size 2×2 . It is applied with stride of 2 pixels.
- The result at using a pooling layer of creating downsampled features is a summarized versions of the features detected in the input.
- Pooling supports multiple scale analysis & it helps in reducing the number of coefficients.

- This in turn results in fewer calculations in training phase to find optimal weights for each neuron.

(j) Max. Pooling on image I .

Max. pooling filter is of dimension 2×2 & applied with stride of 2.

We will apply this filter with each channel respectively.

$$R = \begin{array}{|c|c|c|c|} \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline \end{array} \text{ Max. } \rightarrow \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 1 & 1 \\ \hline \end{array}$$

$$G = \begin{array}{|c|c|c|c|} \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 1 \\ \hline \end{array} \text{ Max. } \rightarrow \begin{array}{|c|c|} \hline 2 & 2 \\ \hline 2 & 2 \\ \hline \end{array}$$

$$B = \begin{array}{|c|c|c|c|} \hline 1 & 1 & 1 & 1 \\ \hline 2 & 2 & 2 & 2 \\ \hline 3 & 3 & 3 & 3 \\ \hline 4 & 4 & 4 & 4 \\ \hline \end{array} \text{ Max. } \rightarrow \begin{array}{|c|c|} \hline 2 & 2 \\ \hline 4 & 4 \\ \hline \end{array}$$

Now we have values for each channel, we take max. value from it to get the final output

$$R \quad G \quad B \quad \text{Max pooling output} \\ \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} \quad \begin{array}{|c|c|} \hline 2 & 2 \\ \hline 2 & 2 \\ \hline \end{array} \quad \begin{array}{|c|c|} \hline 2 & 2 \\ \hline 4 & 4 \\ \hline \end{array} \quad \Rightarrow \quad \begin{array}{|c|c|} \hline 2 & 2 \\ \hline 4 & 4 \\ \hline \end{array}$$