

NPTEL Week 11 Live Sessions

on Deep Learning (noc24_ee04)

A course offered by: Prof. Prabir Kumar Biswas, IIT Kharagpur

- Quiz 10 Solution
- Practice Problems for week 11

Week 8
Q9 → Softmax
Leaky.



By

Arka Roy

NPTEL PMRF TA

Prime Minister's Research Fellow

Department of Electrical Engineering, IIT Patna

Web: <https://sites.google.com/view/arka-roy/home>

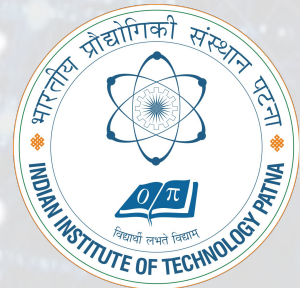
Triplet loss
one shot learning
Segmentation.

Powered by:



PMRF

Prime Minister's Research Fellows
Ministry of Education
Government of India

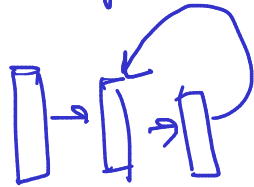


① Which of the following is False regarding batch-normalization??

- True a. Batch normalization prevents the problem of exploding
- b. For an input x , Batch Normalization produces same output regardless of in training mode or inference/test mode. \rightarrow False.
- c. Batch normalization may prevent covariate shift
- d. All of them

i) Vanishing Gradient ($\nabla_w L(w) \rightarrow 0$): $w_{n+1} \leftarrow w_n - \eta \nabla_w L(w) \rightarrow 0$ \rightarrow Regularization, Dropout, Early stopping.

ii) Exploding gradient $\nabla_w L(w) \rightarrow \text{high} \rightarrow \text{inf.}$ \rightarrow NaN

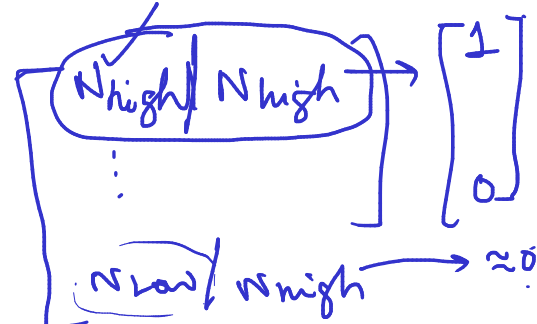


N_{high}

Brought

Normalization

down to a relatively low value.



Which of the following is False regarding batch-normalization??

- ☒ a. Batch normalization prevents the problem of exploding True.
- ☒ b. For an input x , Batch Normalization produces same output regardless of in training mode or inference/test mode.
- ☒ c. Batch normalization may prevent covariate shift True
- ☒ d. ~~All of them~~

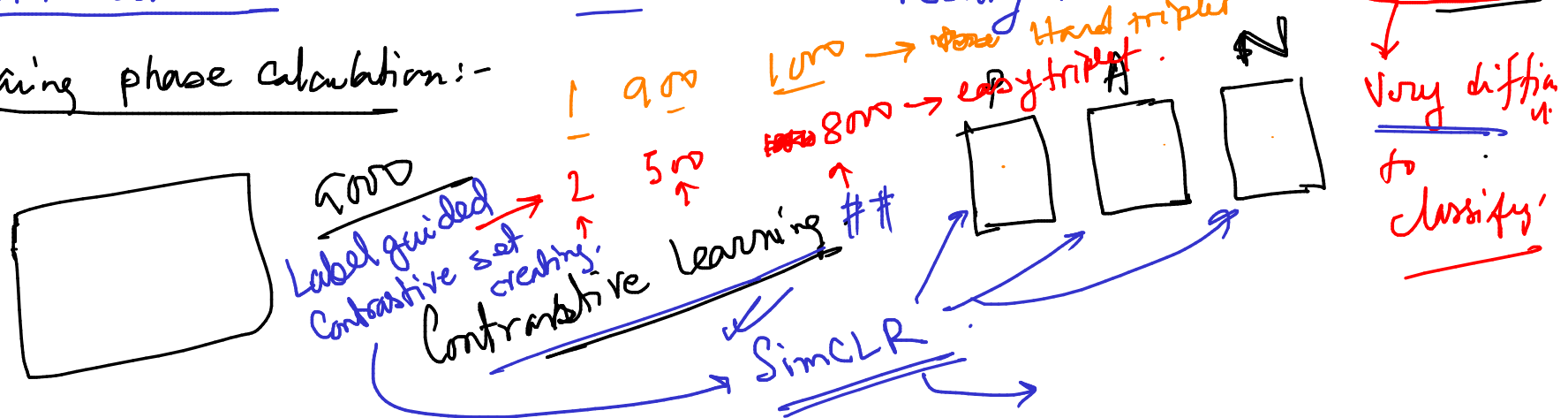


Which of the following is False regarding batch-normalization??

- a. Batch normalization prevents the problem of exploding
 - b. For an input x , Batch Normalization produces same output regardless of in training mode or inference/test mode.
 - c. Batch normalization may prevent covariate shift
 - d. All of them
- Explain*
Easy Triple
Some hard

Batch Normalization: - Mean, Standard dev. $\begin{cases} \nearrow \text{Training phase} \\ \searrow \text{Testing phase} \end{cases}$

Training phase calculation:-



7

A neural network has 3 neurons in a hidden layer. Activations of the neurons for three batches

are $\begin{bmatrix} 4 \\ 4 \\ 3 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 2 \\ 5 \end{bmatrix}$, $\begin{bmatrix} 8 \\ 9 \\ 1 \end{bmatrix}$ respectively. What will be the value of mean if we use batch normalization in

this layer?

- a. $\begin{bmatrix} 4 \\ 5 \\ 3 \end{bmatrix}$
- b. $\begin{bmatrix} 2 \\ 2 \\ 5 \end{bmatrix}$
- c. $\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$
- d. $\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$

$$\begin{bmatrix} 4+0+8/3 \\ 4+2+9/3 \\ 3+5+1/3 \end{bmatrix} = \begin{bmatrix} 4 \\ 5 \\ 3 \end{bmatrix}$$

$$A = \begin{bmatrix} 4 \\ 4 \\ 3 \end{bmatrix}$$

$$A^T = \begin{bmatrix} 4 & 4 & 3 \end{bmatrix}_{1 \times 3}$$

Numpy :-

$$W^T X$$

$$y = X^T W$$

How can we prevent underfitting? (MCO) MSA

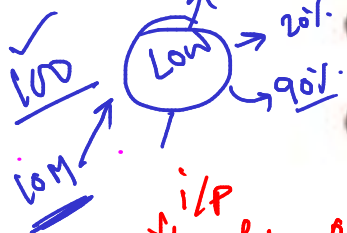
Deep architectures

- Increase the number of data samples
- Increase the number of parameters
- Decrease the number of parameters
- Decrease the number of data samples

② Tripel oneshot

Partially

Correct



Depthwise separable conv.
MobileNet
 Low computation less hardware

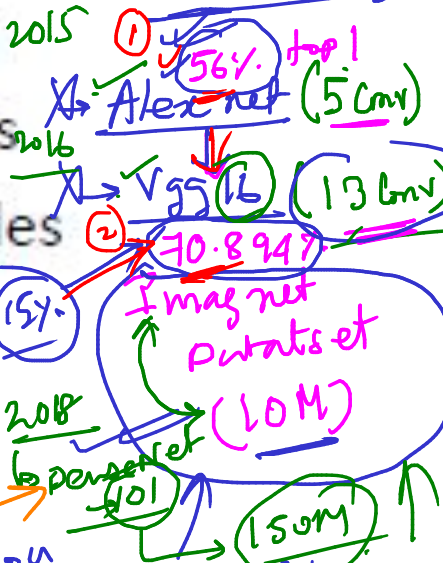
As you increase the number of layer → increase no of Parameters.

Capacity of model increased.

Whether you are learning more something efficient or not?

→ yes

Machine learning



① PASCAL

Pretrained VGG16
Model complexity
Transfer learning - Fine tune

Less data

How do we generally calculate mean and variance during testing?

a. Batch normalization is not required during testing

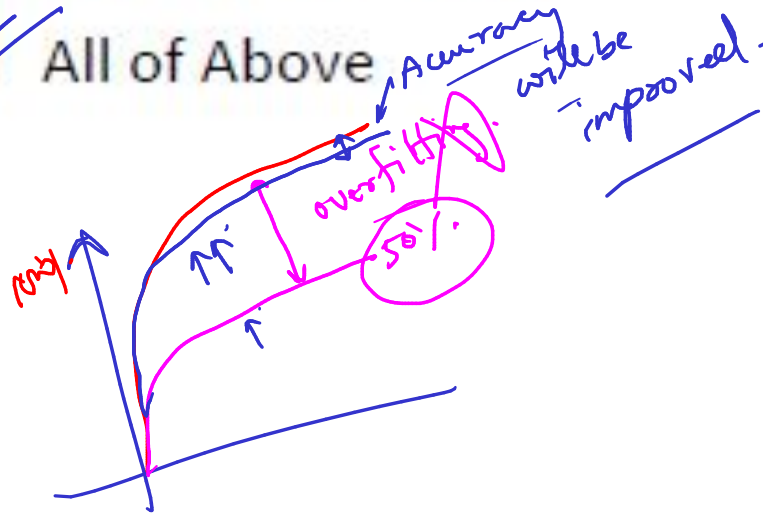
b. Mean and variance based on test image

☒ c. Estimated mean and variance statistics during training

d. None of the above

Which one of the following is an advantage of dropout?

- ☒ a. Regularization
- ☒ b. Prevent Overfitting
- ☒ c. Improve Accuracy
- ☒ d. All of Above



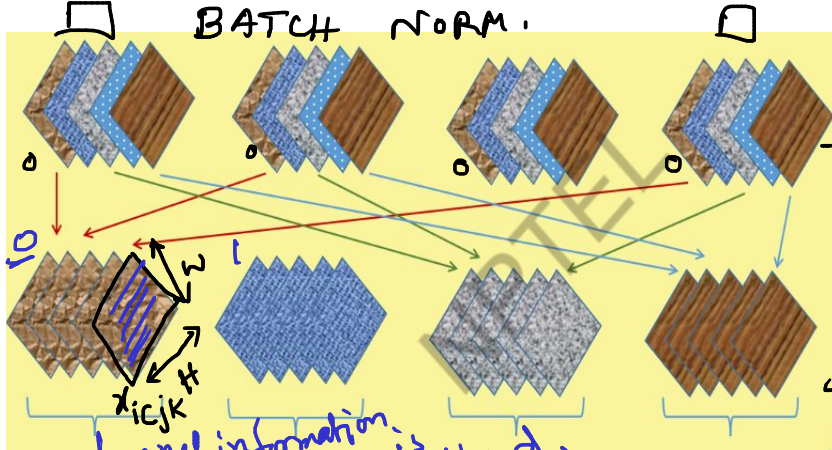
reduce overfitting.
regularization.

① ~~6~~ ②
→ 0
~~6~~ ~~6~~
→ 0
6 6
} Regularizing
Technique

Which of the following is True regarding layer normalization and batch normalization?

- a. ☒ Layer normalization normalizes over spatial and channel dimension whereas Batch normalization normalizes over spatial and batch dimensions
- b. Layer normalization normalizes over **spatial** and **channel** dimension whereas Batch normalization normalizes over **batch**, **spatial** and **channel** dimension
- c. Batch normalization normalizes over spatial and ~~channel~~ dimension whereas Layer normalization normalizes over **batch**, **spatial** and **channel** dimension
- d. None of these

BATCH NORM.



channel information is fixed.

$$x \in \mathbb{R}^{N \times C \times W \times H}$$

$$\mu_C = \frac{1}{NWH} \sum_{i=1}^N \sum_{j=1}^W \sum_{k=1}^H x_{icjk}$$

Mean

$$\sigma_C^2 = \frac{1}{NWH} \sum_{i=1}^N \sum_{j=1}^W \sum_{k=1}^H (x_{icjk} - \mu_C)^2$$

$$\hat{x} = \frac{x - \mu_C}{\sqrt{\sigma_C^2 + \epsilon}}$$

Training \neq Testing.

Tensor

Exponentially decaying average based on

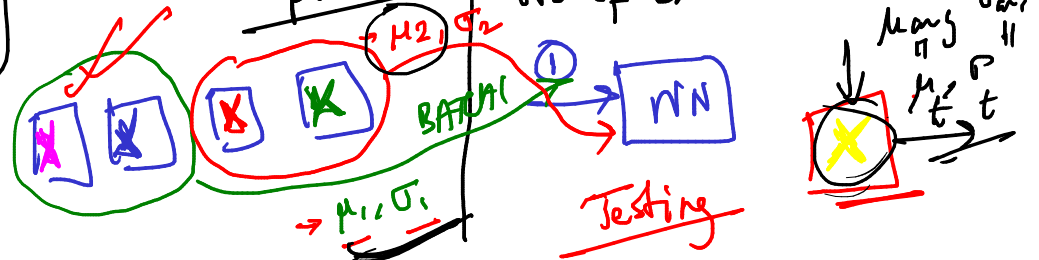
Spatial dimension (height width of Tensors)

Batch information.

Testing: \rightarrow Avg of all the μ and σ evaluated in the training phase

Estimated statistics

Training Mean and Variance extraction process



DATASET: 10,000 sample.

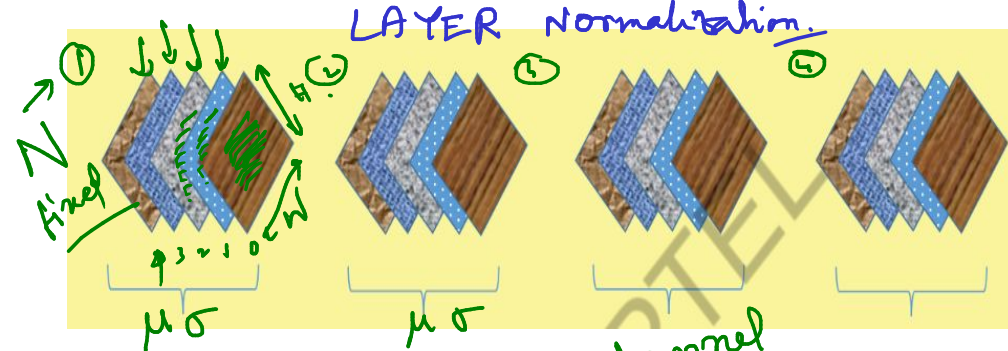
9000 samples \rightarrow Training
1000 " \rightarrow Testing.

Batch Size = 10.

No of batches = 900

Mean μ_t
Variance σ_t^2

LAYER Normalization



$$x \in \mathbb{R}^{N \times C \times W \times H}$$

$$\mu_N = \frac{1}{CWH} \sum_{i=1}^C \sum_{j=1}^W \sum_{k=1}^H x_{Nijk}$$

$$\sigma_N^2 = \frac{1}{CWH} \sum_{i=1}^C \sum_{j=1}^W \sum_{k=1}^H (x_{Nijk} - \mu_N)^2$$

$$\hat{x} = \frac{x - \mu_N}{\sqrt{\sigma_N^2 + \epsilon}}$$

Channel dimension.

spatial width / height.

Which one of the following regularization methods induces sparsity among the trained weights?

$$\nabla \frac{\partial L(w)}{\partial w}$$

Ridge

a. L_1 regularizer Lasso Absolute value

b. L_2 regularizer

c. Both L_1 & L_2

d. None of the above

$$\text{Loss}(w) = \text{Data Loss} + \text{Regularizer.}$$

$$\Rightarrow L(w) = \text{Data Loss} + \text{L}_1 \text{ regularizer of weight}$$

$$= g(w) + \text{L}_1 \text{ regularizer.}$$

$$W(m+1) \leftarrow W(m) - \eta \nabla_w L(w).$$

$$|-2| = 2$$

$$|2| = 2.$$

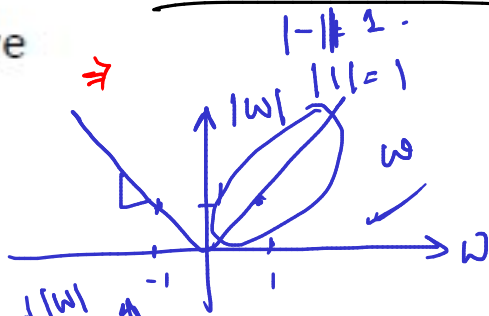
$$L(w) = g(w) + |w|$$

$$L(w) = g(w) + \lambda \sum_{i=1}^n |w_i|$$

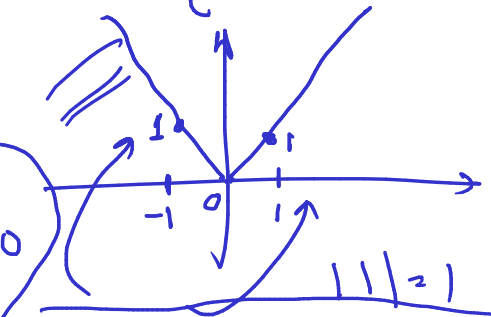
strength of regularizer

$$\nabla_w L(w) = g'(w) + \lambda \text{sgn}(w)$$

Data loss differentiation with respect to w .



$$\frac{dy}{dx} = \frac{y_2 - y_1}{x_2 - x_1} = \frac{1-0}{-1-0} = -1$$



$$\frac{dy}{dx} = 1 \cdot \frac{y_2 - y_1}{x_2 - x_1} = \frac{1-0}{1-0} = 1$$

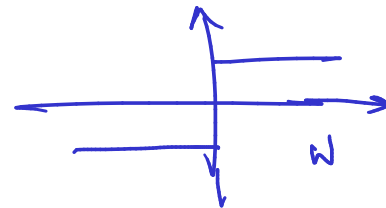
$$w_{n+1} \leftarrow w_n - \nabla_w L(w)$$

$$w_{n+1} \leftarrow w_n - [g'(w) + \lambda \text{sgn}(w)]$$

$w > 0$,

$$w_{n+1} \leftarrow \underbrace{w_n}_{\text{iteration } m} - \underbrace{g'(w)}_{\text{shrinking}} - \underbrace{\lambda}_{\text{shrinking}}$$

→ Reducing / shrinking / $\rightarrow 0 \sqrt{a^2}$
iteration 1, 2, 3, ...



$$\|w\|_2^{(2)} \rightarrow \sqrt{a^2}$$

$$L(w) = g(w) + L^2 \text{reg}(w)$$

$$L(w) = g(w) + \lambda \sum_{i=1}^N |w_i|^{(2)} \rightarrow \lambda \cdot w_i^2$$

$$\frac{\partial L(w)}{\partial w} = g'(w) + 2\lambda w_i$$

$$L^1 \text{reg.} \rightarrow 0$$

$$\begin{bmatrix} 50 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{sparsity}$$

$$w_{n+1} \leftarrow w_n - \nabla_w L(w)$$

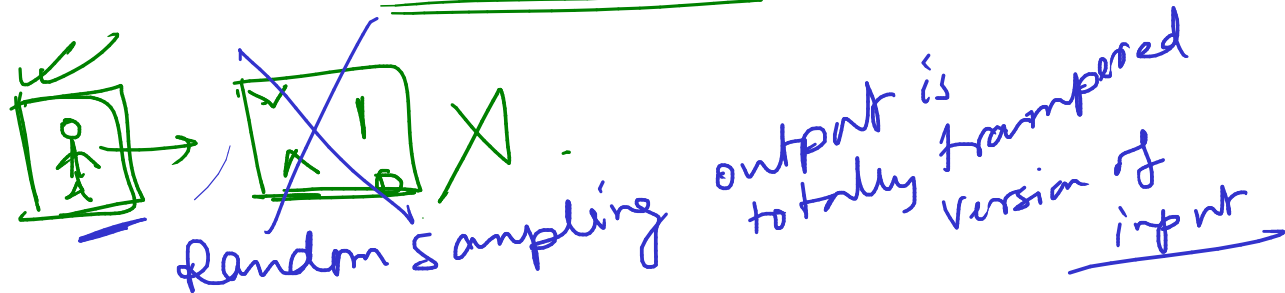
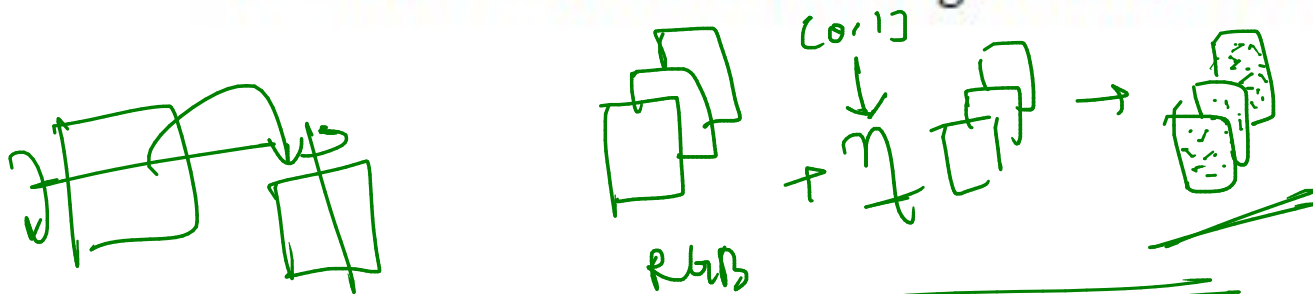
$$w_{n+1} \leftarrow \underbrace{w_n}_{\text{iteration } m} - \underbrace{g'(w)}_{\text{shrinking}} - \underbrace{2\lambda w_i}_{\text{shrinking}}$$

either they will be of very small in nature or they will be rejected. Reduce the value of weight.

$$\sqrt{w_1^2 + w_2^2}$$

Which among the following is NOT a data augmentation technique?

- a. Random horizontal and vertical flip of image
- ☒ b. Random shuffle all the pixels of an image
- c. Random color jittering
- d. All the above are data augmentation techniques



Which of the following is true about model capacity (where model capacity means the ability of neural network to approximate complex functions)?

- a. As number of hidden layers increase, model capacity increases
- b. As dropout ratio increases, model capacity increases
- c. As learning ~~rate~~ increases, model capacity increases
- d. None of these

Two variant training schedule samples its minibatches in the following manner

Training Schedule 1

{ Mini batch 1=[Image1, Image2, Image3]
Mini batch 2=[Image4, Image5, Image6]

Training Schedule 2

{ Mini batch 1=[Image1, Image4, Image3]
Mini batch 2=[Image2, Image5, Image6]

images inside the minibatch is different \rightarrow their statistics μ, σ will also differ \rightarrow Activation will also differ.

$\hat{x} = \frac{x - \tilde{\mu}}{\tilde{\sigma}}$

The output activations of each corresponding image is compared across Training schedule 1 and Training schedule 2 for a CNN with batch norm layers. Choose the correct statement

- a. Activation outputs of corresponding image will be same across Training schedule 1 and Training schedule 2
- ☒ b. Activation outputs of corresponding image will be different across Training schedule 1 and Training schedule 2
- c. Some activations outputs of corresponding images will be same but some will be different
- d. None of these.

Suppose you have a 1D signal $x = [5, 4, 3, 2, 1]$ and a filter $f = [1, 2, 3, 4]$, and you perform stride 2 transpose convolution on the signal x by the filter f to get the signal y . What will be the signal y if we don't perform cropping?

a. $y = [1, 2, 5, 8, 9, 14, 13, 20, 19, 26, 3, 4]$

~~b. $y = [5, 10, 19, 28, 15, 22, 11, 16, 7, 10, 3, 4]$~~

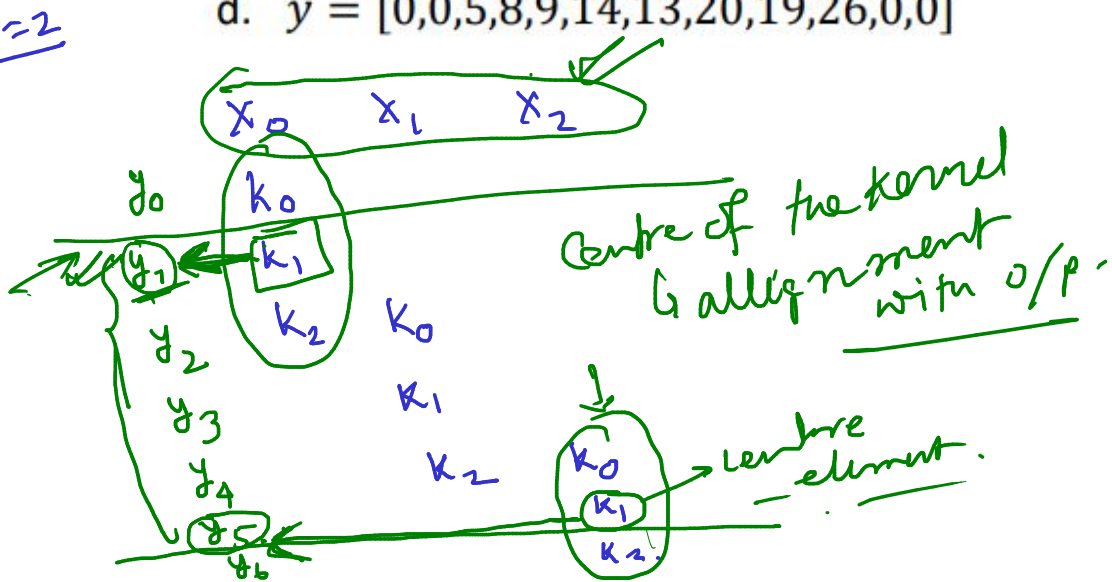
c. $y = [1, 2, 5, 8, 9, 14, 13, 20, 17, 26, 15, 20]$

d. $y = [0, 0, 5, 8, 9, 14, 13, 20, 19, 26, 0, 0]$

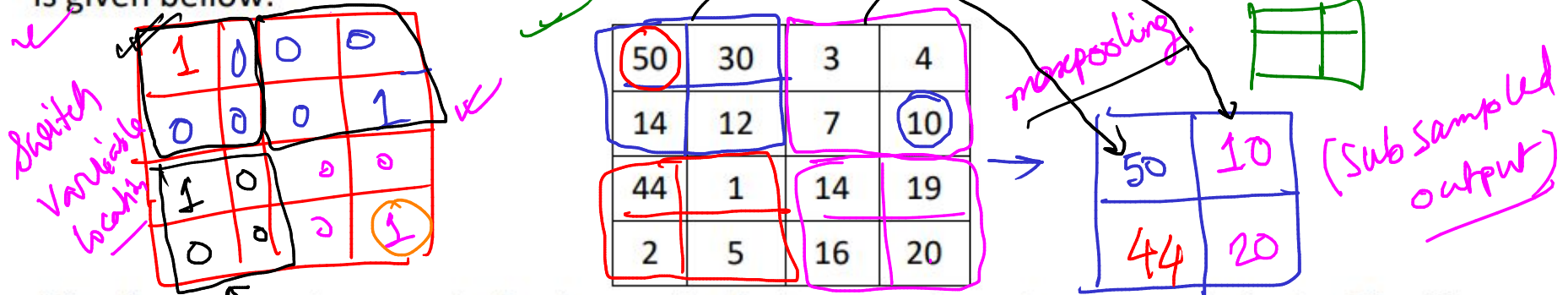
$[5 \quad 4 \quad 3 \quad 2 \quad 1]$

Handwritten calculation for the output signal y using the filter $f = [1, 2, 3, 4]$ and stride 2:

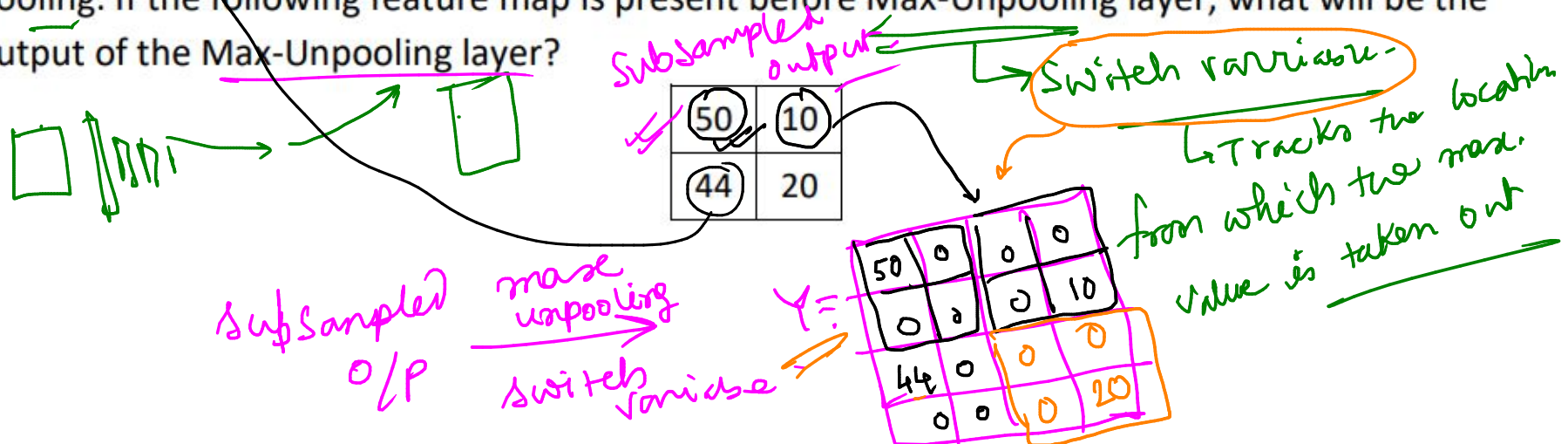
$$\begin{array}{l}
 y_0: [1] \rightarrow 5 \\
 y_1: [2] \rightarrow 10 \\
 y_2: [3, 1] \rightarrow 15 + 4 = 19 \\
 y_3: [4, 2] \\
 y_4: [3, 3] \\
 y_5: [4, 4, 2] \\
 y_6: [3, 4] \\
 y_7: [1, 2] \\
 y_8: [3, 3] \\
 y_9: [4, 4] \\
 y_{10}: [1, 2] \\
 y_{11}: [3, 4]
 \end{array}$$



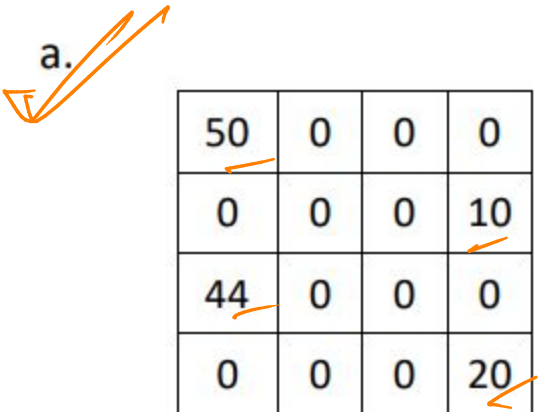
In a Deep CNN architecture, the feature map before applying a max pool layer with (2x2) kernel is given bellow.



After few successive convolution layers, the feature map is again up-sampled using Max Unpooling. If the following feature map is present before Max-Unpooling layer, what will be the output of the Max-Unpooling layer?



a.



50	0	0	0
0	0	0	10
44	0	0	0
0	0	0	20

b.

50	0	10	0
0	0	0	0
44	0	20	0
0	0	0	0

c.

50	0	0	10
0	0	0	0
44	0	0	20
0	0	0	0

d.

0	0	0	0
0	50	0	10
0	0	0	0
0	44	0	20

~~✓~~

<u>0.10</u>	<u>0.47</u>	0.88	0.25
0.65	0.11	0.66	1
0.91	0.00	0.99	0.29
0.78	1	0.40	0.99

~~✓~~

A

Ground truth mask

<u>0.01</u>	<u>0.44</u>	0.48	0.18
0.60	0.09	0.97	0.99
0.87	0.02	0.98	0.20
0.77	0.99	0.40	0.99

~~✓~~

B

Prediction mask

~~Dice Coefficient~~ $(A \cap B)$

=

$$\sum_{i=1}^n A \times B$$

$$\sum_{i=1}^n A + \sum_{i=1}^n B$$