

NPTEL Week 12 Live Sessions

on Deep Learning (noc24_ee04)

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

- Quiz 11 Solution
- Practice Problems for week 12



By

Arka Roy
NPTEL PMRF TA

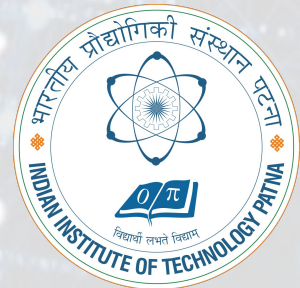
Prime Minister's Research Fellow
Department of Electrical Engineering, IIT Patna
Web: <https://sites.google.com/view/arka-roy/home>

Powered by:



PMRF

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



Which of following can be a target output of semantic segmentation problem with 4 class?

a.

0	1	0
0	1	0
1	0	0

I

1	1	0
0	0	0
0	0	0

II

0	0	1
1	0	0
0	0	0

III

0	0	0
1	0	0
0	0	0

IV

b.

0	1	0
0	1	0
1	0	0

I

1	0	0
0	1	0
0	0	0

II

0	0	1
1	0	0
0	0	0

III

0	1	0
0	0	0
0	1	1

IV

c.

0	1	0
0	1	0
1	0	0

I

1	0	0
0	0	1
0	0	0

II

0	0	1
1	0	0
0	0	0

III

0	0	0
0	0	0
0	1	1

IV

d.

0	1	0
0	1	0
1	1	0

I

1	0	0
1	0	0
0	1	0

II

0	0	1
1	0	0
0	0	1

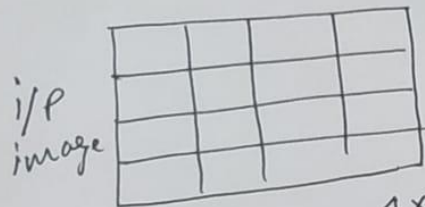
III

0	0	0
0	0	0
0	1	1

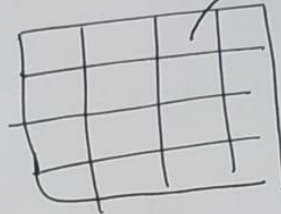
IV

- ① In case of image segmentation, the input image shape should be equal to output image shape.

But every pixel inside the output image will have the onehot vectors.



4x4



o/p image.

4x4.

$[0 \ 1 \ 0 \ 0]$

Let's say we have four classes.

in case of onehot vector only one element is meant to have 1 (hot) rest are zero (cold)

Therefore in the options we have to search for that every element in a particular pixel point should have the same property of one hot vector.

What will be the dice coefficient of following two one hot encoded vector? ($|A|$ =no of 1 bit)

A	1	0	1	0	0	0	1	1	1	0	0	1	0	1
B	1	0	0	0	0	1	1	1	0	0	0	1	0	0

- a. 0.83
- b. 0.41
- c. 0.67
- d. 0.90

②

1	0	1	0	0	1	1	1	0	0	1	0	1
1	0	0	0	0	1	1	1	0	0	0	1	0

$|A| = 7$ $|B| = 5$

$|A \cap B| = \text{Multiplication between } A \text{ and } B$
 $= |1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0|$
 $= 4$

Dice Coeff = $\frac{2|A \cap B|}{|A| + |B|}$
 $= \frac{2 \times 4}{7 + 5} = \frac{2 \times 4}{12} = \frac{2}{3} = 0.67$

What will be the value of dice coefficient between A and B?

0.01	0.03	0.02	0.02
0.05	0.12	0.09	0.07
0.89	0.85	0.88	0.91
0.99	0.97	0.95	0.97

A

0	0	0	0
0	0	0	0
1	1	1	1
1	1	1	1

(Consider, $|A|$ = sum of all elements

- a. 0.23
- b. 0.77
- c. 0.11
- d. 0.93

$$\begin{aligned} \textcircled{3} \quad |A| &= 0.01 + 0.03 + 0.02 + 0.02 + 0.05 + 0.12 \\ &\quad + 0.09 + 0.07 + 0.89 + 0.85 + 0.88 + 0.91 \\ &\quad + 0.99 + 0.97 + 0.95 + 0.97 \\ &= 7.82 \end{aligned}$$

$$|B| = 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 = 8$$

$$|A \cap B| = \sum_i A(i) \cdot B(i) = 7.42$$

$$\text{Dice coeff} = \frac{2 \times |A \cap B|}{|A| + |B|}$$

Suppose you have a 1D signal $x = [1, 2, 3, 4, 5]$ 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. $y = [1, 2, 3, 4, 5, 4, 3, 2, 1]$
- c. $y = [1, 2, 5, 8, 9, 14, 13, 20, 17, 26, 15, 20]$

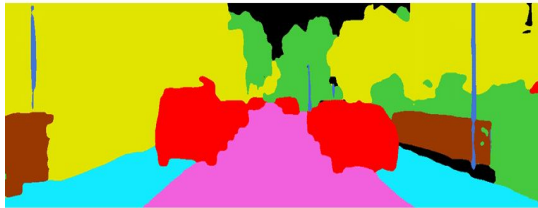
④


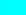


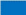
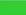


	$x = [1 \quad 2 \quad 3 \quad 4 \quad 5]$					
k	$\begin{bmatrix} 1 \end{bmatrix}$					$f(0) = 1$
e	$\begin{bmatrix} 2 \end{bmatrix}$					$f(1) = 1 \times 2 = 2$
y	$\begin{bmatrix} 3 \end{bmatrix}$	$\begin{bmatrix} 1 \end{bmatrix}$				$f(2) = 3 + 2 = 5$
4	$\begin{bmatrix} 4 \end{bmatrix}$	$\begin{bmatrix} 2 \end{bmatrix}$				$f(3) = 4 + 4 = 8$
e		$\begin{bmatrix} 3 \end{bmatrix}$	$\begin{bmatrix} 1 \end{bmatrix}$			$f(4) = 6 + 3 = 9$
1		$\begin{bmatrix} 4 \end{bmatrix}$	$\begin{bmatrix} 2 \end{bmatrix}$			$f(5) = 8 + 6 = 14$
			$\begin{bmatrix} 3 \end{bmatrix}$	$\begin{bmatrix} 1 \end{bmatrix}$		$f(6) = 9 + 4 = 13$
			$\begin{bmatrix} 4 \end{bmatrix}$	$\begin{bmatrix} 2 \end{bmatrix}$		$f(7) = 12 + 8 = 20$
				$\begin{bmatrix} 3 \end{bmatrix}$	$\begin{bmatrix} 1 \end{bmatrix}$	$f(8) = 12 + 5 = 17$
				$\begin{bmatrix} 4 \end{bmatrix}$	$\begin{bmatrix} 2 \end{bmatrix}$	$f(9) = 16 + 10 = 26$
					$\begin{bmatrix} 3 \end{bmatrix}$	$f(10) = 15$
					$\begin{bmatrix} 4 \end{bmatrix}$	$f(11) = 20$

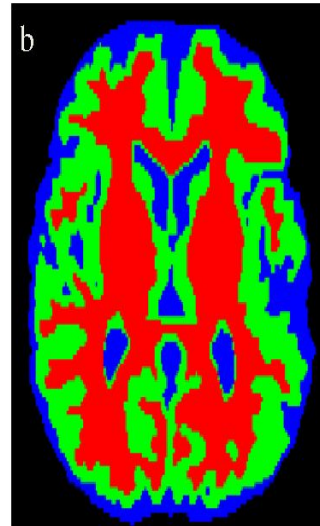
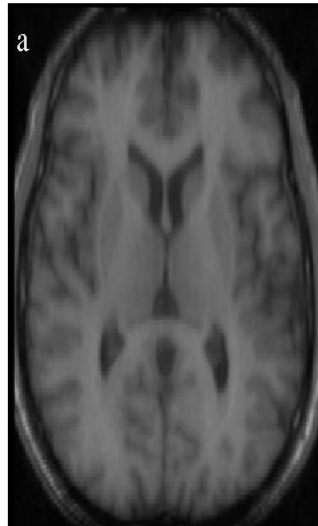
$$y = [1, 2, 5, 8, 9, 14, 13, 20, 17, 26, 15, 20]$$

Which of the following is true for semantic segmentation?

- a. Semantic Segmentation can be considered as pixel wise classification problem.
- b. Semantic Segmented output has same dimension as the input image dimension.
- c. It has application in Autonomous driving, Industrial inspection, and Medical imaging analysis.
- d. All of the above



 Road	 Sidewalk	 Building	 Fence
 Pole	 Vegetation	 Vehicle	 Unlabel



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

2	30	3	14
14	12	7	10
4	1	14	19
2	5	16	2

After few successive convolution layers, the feature map is again up-sampled using Max Un-pooling, what will be the output of the Max-Unpooling layer?

a.

2	30	3	14
14	12	7	10
4	1	14	19
2	5	16	2

b.

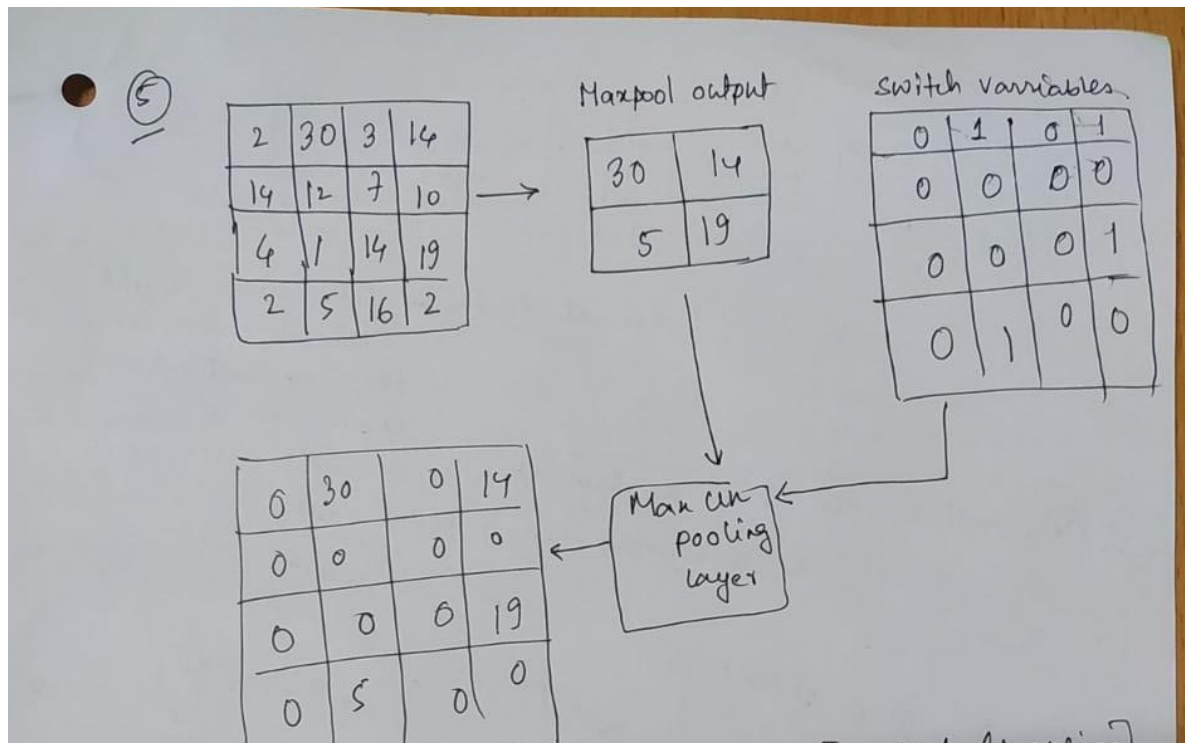
0	0	0	14
14	0	0	0
0	0	0	0
0	5	0	2

c.

0	30	0	14
0	0	0	0
0	0	0	19
0	5	0	0

d.

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



Which of the following operation reduces spatial dimension of features?

- a. Max un-Pooling
- b. Convolution with 3×3 Kernel, Stride=2, Padding all sides = 1
- c. Convolution with 3×3 Kernel, Stride=1, Padding all sides = 1
- d. Transposed convolution

⑥ Let's take the image shape $I_s = x \times x$ [spatial dimension].
or $= \underbrace{x \times x \times 3}_{\text{channel dimension}}$

① 3×3 kernel, stride = 1, Pad = 1 \Rightarrow

$$\text{o/p dimension} = \frac{x - \text{kernel} + (2 \times \text{padding})}{(\text{stride})} + 1$$
$$= \frac{x - 3 + 2}{1} + 1$$
$$= x - 3 + 3 = x.$$

i/p dim = o/p dim

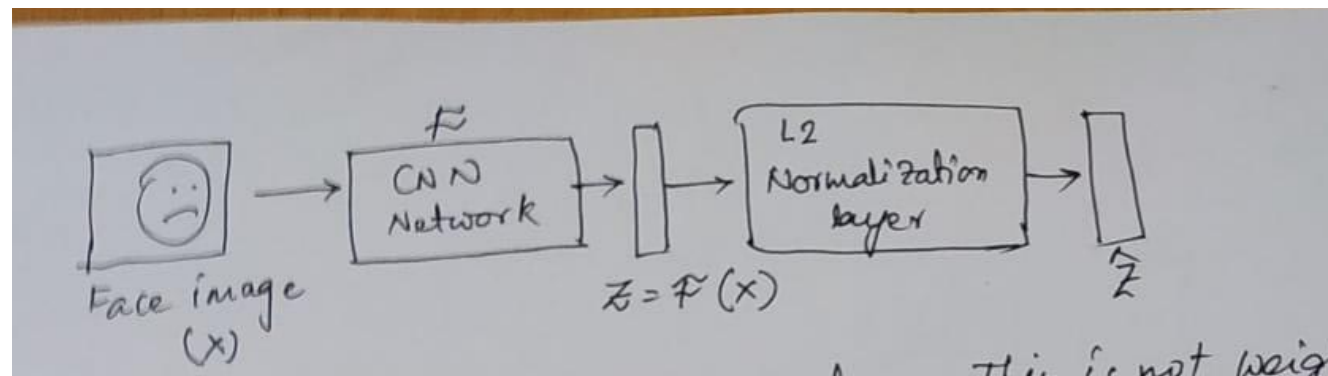
② stride = 2, $\text{o/p dim} = \frac{x - 3 + 2 \times 1}{2} + 1$

$$= \frac{x - 1}{2} + 1$$

??

In FaceNet, why the L2 normalization layer is used?

- a. To constrain the embedding function in a d-dimensional hyper-sphere.
- b. For regularization of weight vector, i.e. L2 regularization.
- c. For getting a sparse embedding function.
- d. None of the above.



\tilde{x} is L^2 Normalized \rightarrow This is not weight normalization as the normalization is applied on \tilde{x} .

* As well as this is not a regularization of weight vector as it is not L^2 regularizer that is being asked.

$$\begin{array}{ccc} \boxed{\tilde{x}} & \xrightarrow[\substack{L^2 \text{ normalization} \\ \text{layer } (f(\cdot))}]{\tilde{x} \in \mathbb{R}^{d \times 1}} & \boxed{\hat{z}} \\ \tilde{x} & & \hat{z} \end{array} \quad \hat{z} = f(\tilde{x} \in \mathbb{R}^{d \times 1})$$

$$= \sqrt{\tilde{z}_1^2 + \tilde{z}_2^2 + \tilde{z}_3^2 + \dots + \tilde{z}_d^2} = 1.$$

$$\therefore \hat{z} = \sqrt{\sum_{i=1} \tilde{z}_i^2} = 1.$$

$$\sum_{i=1} \tilde{z}_i^2 = 1. \quad \text{if it is a 3d vector then } \rightarrow$$

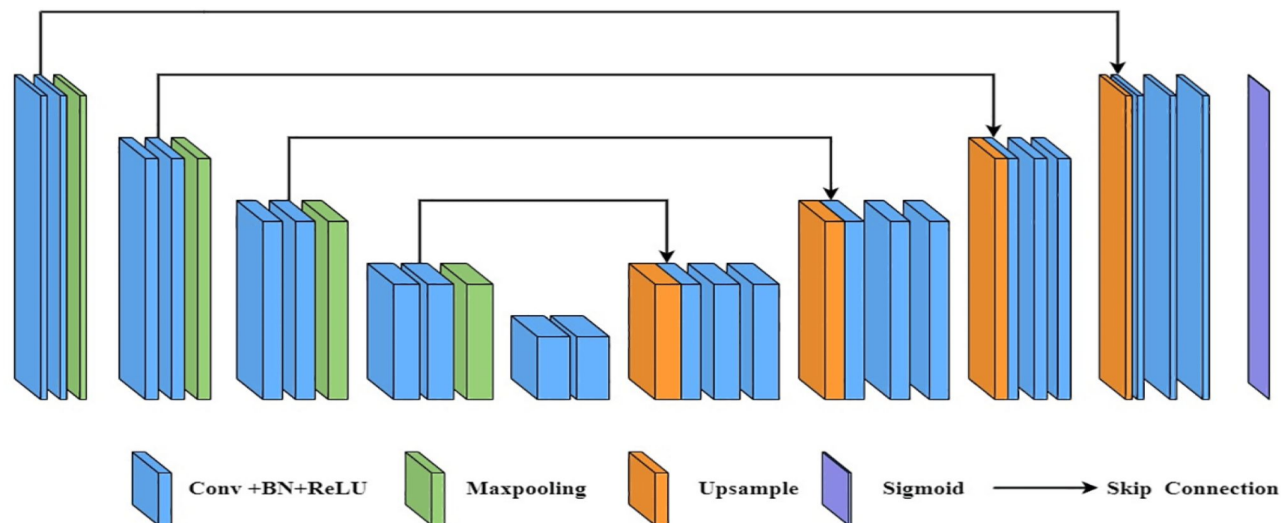
$$\tilde{z}_1^2 + \tilde{z}_2^2 + \tilde{z}_3^2 = 1 \rightarrow \text{sphere.}$$

if $d > 3$ then $\underbrace{\tilde{z}_1^2 + \tilde{z}_2^2 + \tilde{z}_3^2 + \dots}_{\rightarrow \text{High dimensional sphere.}} = 1.$

\rightarrow Hyper sphere.

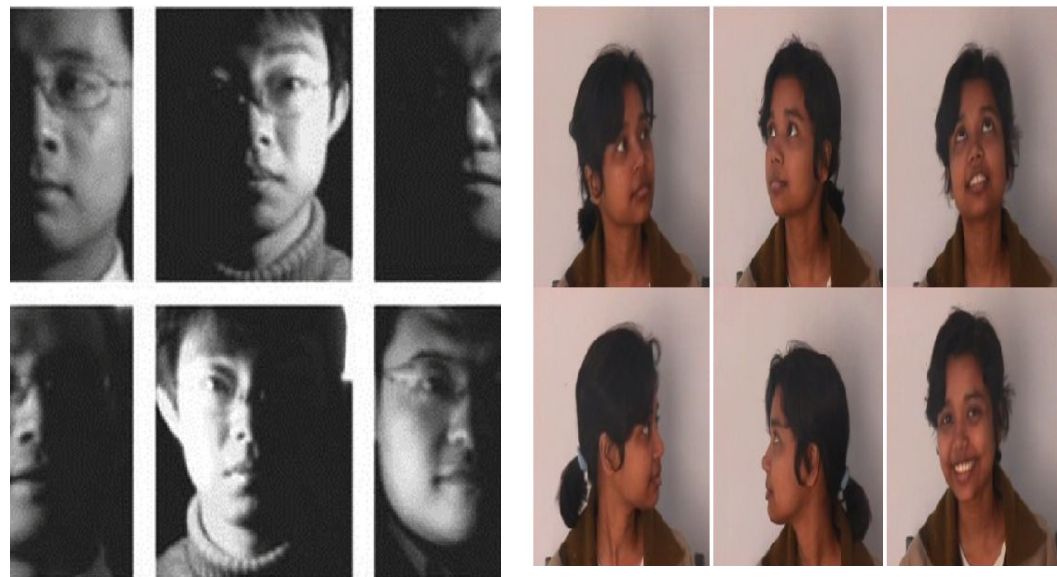
What is the use of Skip Connection in image denoising networks?

- a. Helping de-convolution layer to recover an improved clean version of image.
- b. Back propagating the gradient to bottom layers, which makes the training easy.
- c. To create the direct path between convolution layer and the corresponding mirror de-convolution layer.
- d. All of the above.



What are the different challenges one face while creating a facial recognition system?

- a. Different illumination condition
- b. Different pose and orientation of face images
- c. Limited dataset for training
- d. All of the above



Learning Face Recognition from Limited Training Data using Deep Neural Networks

Xi Peng¹
Department of Computer Science
Rutgers University
Piscataway, New Jersey 08854

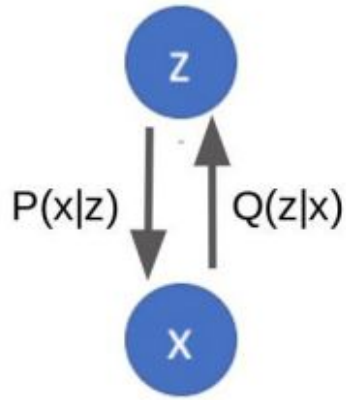
Nalini Ratha
IBM Thomas J. Watson
Research Center
Yorktown Heights, New York 10598

Sharatchandra Pankanti
IBM Thomas J. Watson
Research Center
Yorktown Heights, New York 10598

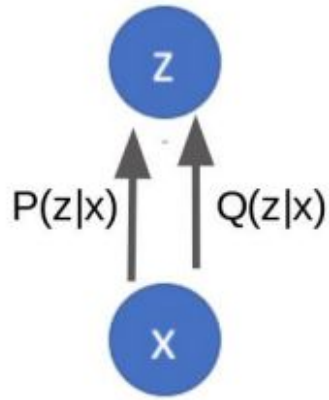
Multi-Angled Face Segmentation and Identification using Limited Data

Dane Brown
Department of Computer Science
Rhodes University
Grahamstown, South Africa
d.brown@ru.ac.za

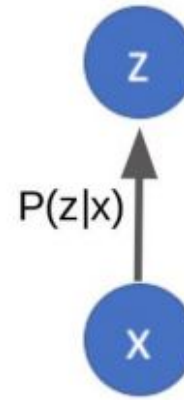
Which one of the following graphical models fully represents a Variational Auto-encoder (VAE) realization?



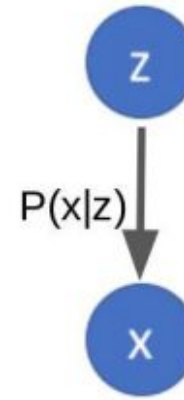
(a)



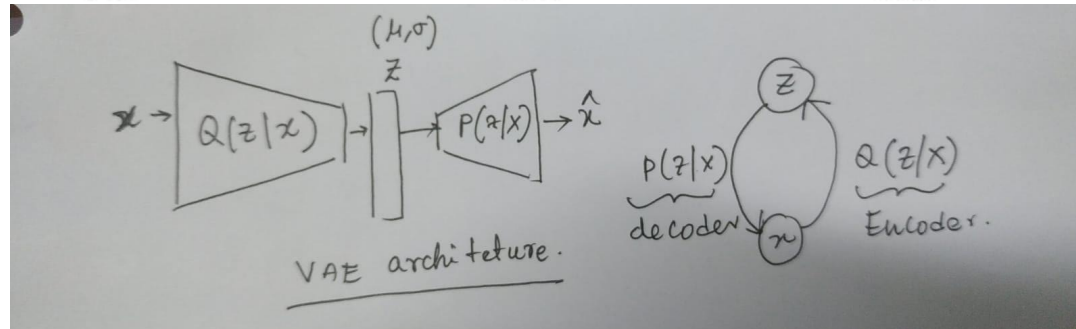
(b)



(c)



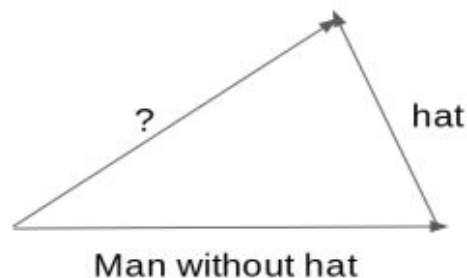
(d)



Which of the following is an INVALID activation function inside a neural network?

- a. $f(x) = \max(0, 2x)$
- b. $f(x) = \min(0, 2x)$
- c. $f(x) = \tanh(x)$
- d. None of the above

Figure shows latent vector addition of two concepts of “man without a hat” and “hat”. What is expected from the resultant vector?



- a. Hat without man
- b. Man with hat
- c. Woman with hat
- d. Woman without hat