

28/12/25

Autoencoder (questions) has one hidden layer

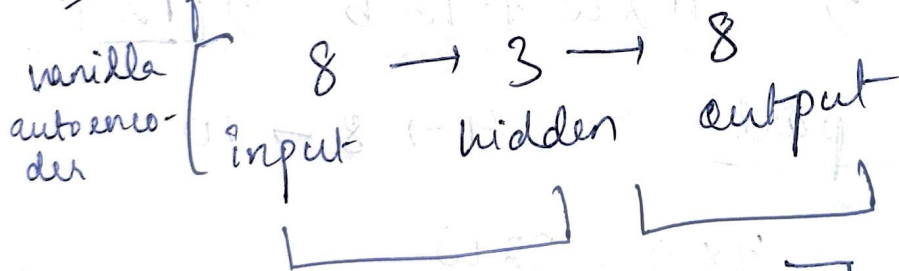
Q.1. An auto encoder [ vanilla autoencoder ], where ANN is used for encoder and decoder architecture]

Given  $\rightarrow$  Input layer  $\rightarrow$  has 10 neurons  
 $\hookrightarrow$  Hidden layer  $\rightarrow$  has 4 neurons

Then, in output layer  $\rightarrow$  how many neurons will be there?

Output layer  $\rightarrow$  10 neurons (ans)

2. If the auto encoder architecture is:



$$\begin{aligned} 8 \times 3 &= 24 \text{ weights} \\ + 3 \text{ bias} \\ &= \underline{\underline{27}} \end{aligned}$$

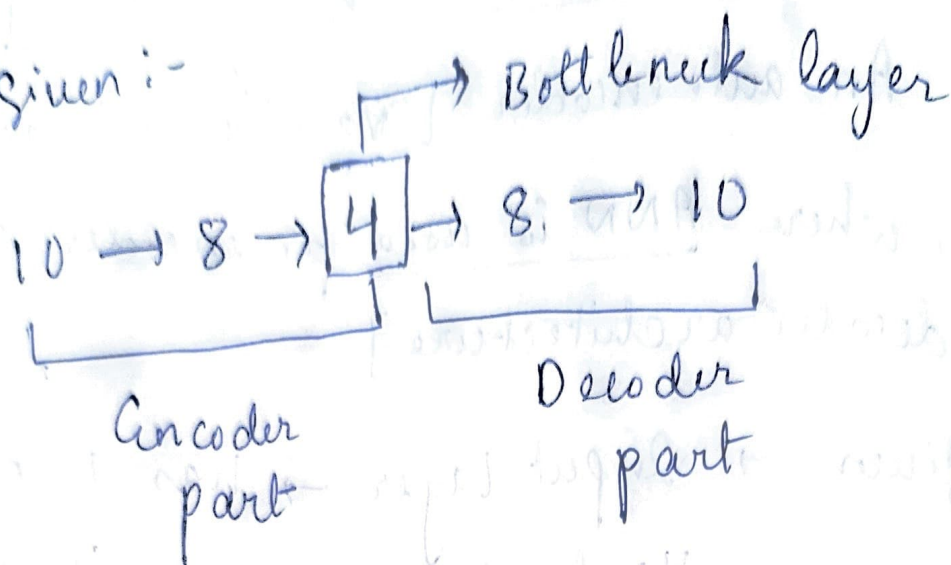
$$\begin{aligned} 3 \times 8 &= 24 \text{ weights} \\ + 8 \text{ bias} \\ &= \underline{\underline{32}} \end{aligned}$$

How many parameters are there in this autoencoder? (simple vanilla autoencoder)

Total parameters  $\rightarrow 27 + 32 = \boxed{59}$  (ans)

### 3. Stacked autoencoder / deep autoencoder

Given:-



Total parameters = ?

i)  $10 \rightarrow 8 \rightarrow 4$  (Encoder)

$$10 \rightarrow 8 \Rightarrow 10 \times 8 = 80 \text{ weights}$$

$$8 \rightarrow 4 \Rightarrow 8 \times 4 = 32 \text{ weights}$$

$$10 \rightarrow 8 \rightarrow 4 \Rightarrow 80w + 32w + 8b + 4b$$

$$\Rightarrow 112w + 12b \Rightarrow \boxed{124}$$

ii) Decoder part :-  $4 \rightarrow 8 \rightarrow 10$

$$\Rightarrow 4 \times 8 = 32w$$

$$\Rightarrow 8 \times 10 = 80w$$

$$\Rightarrow 8 + 10 = 18b$$

$$\Rightarrow 112w + 18b \Rightarrow \boxed{130}$$

Total parameters  $\Rightarrow$

$$124$$

$$+ 130$$

$$\hline 254$$



$$254$$

parameters

(ans)



4. Suppose <sup>in</sup> an autoencoder:-

Input neurons = 50, hidden neurons = 80

Then, what kind of autoencoder is it?  
(what type)

Ans: There are two types:-

✓  
undercomplete  
autoencoder

↓  
overcomplete  
autoencoder

⇒ here, it is an overcomplete autoencoder  
[where the latent vector or the hidden neurons  
consist of more number of neurons in  
comparison to the input neurons]

→ hidden neurons  $\geq$  (greater than or equal to) input neurons  $\Rightarrow$  overcomplete

→ hidden neurons  $<$  input neurons  $\Rightarrow$  undercomplete.

⇒ Training process or Forward propagation of a vanilla autoencoder:

a. Binary autoencoder (forward pass)

Given:- Binary input,  $x = [1, 0]$

Encoder weight  $\rightarrow W_e$

$$w_e = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix}, \quad b_e = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

↓  
bias of  
encoder

Decoder weight  $\rightarrow w_d = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix},$

$$b_d = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Q.  $\rightarrow$  what is the structure/architecture of the auto encoder based on given information.

ans  $\rightarrow$  2 input neurons, 2 hidden + 2 output:  $2 \rightarrow 2 \rightarrow 2$

What type of auto encoder is this?

$\rightarrow$  Overcomplete [input  $\leq$  hidden]

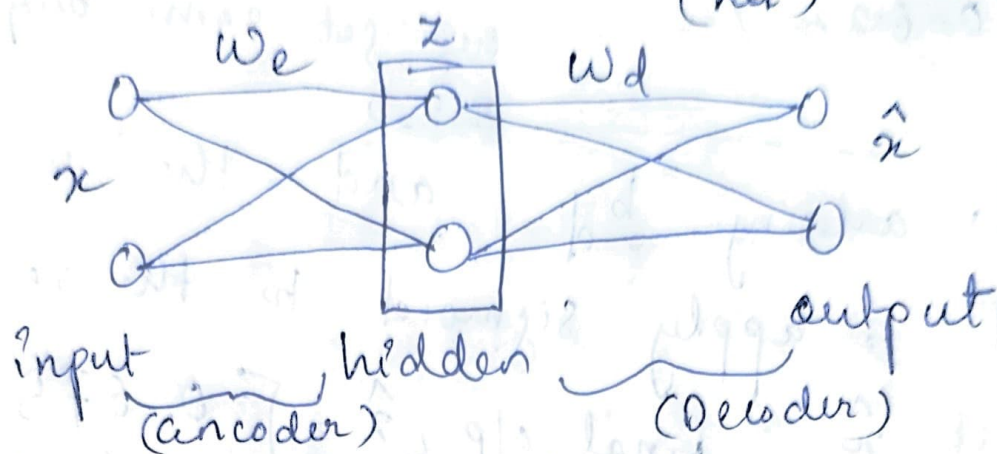
Q. what task is it performing?

[It copies the content of the input, directly to the decoder neurons. This architecture can perform a trivial copy operation. Input information  $\rightarrow$  passed to hidden layer  $\rightarrow$  decoder to produce a reconstructed output]

$\rightarrow$  overcomplete auto encoder



⇒ given activation function: Sigmoid (let)



① Equation of Encoder → sigmoid

Linear Transformation:  $\rightarrow \sigma(w_e \cdot x + b_e)$   
 $= z$  ;  $z \Rightarrow$  o/p of encoder. bias of encoder

② Equation of Decoder :-

$z$  is input to decoder (o/p of encoder)

$$\hat{x} = \sigma(w_d \cdot z + b_d)$$

Now,  $w_e \cdot x = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$   
 $2 \times 2 \cdot 2 \times 1$

$$w_e x + b_e \Rightarrow \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$z = \sigma \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right) = \begin{bmatrix} \sigma(1) \\ \sigma(0) \end{bmatrix} = \begin{bmatrix} 0.73 \\ 0.5 \end{bmatrix}$$

↳ encoder output ←

$$w_d \cdot z = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \times \begin{bmatrix} 0.73 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.675 \\ 0.622 \end{bmatrix}$$

$$\rightarrow \begin{pmatrix} 0.675 \\ 0.622 \end{pmatrix}$$

we get some output

adding  $b_d$  and then we need to apply sigmoid to the result to get  $\hat{x}$ . Final op,  $\hat{x} = \begin{bmatrix} 0.675 \\ 0.622 \end{bmatrix}$

$\hat{x} = [0.68, 0.62] \rightarrow$  reconstructed output

Reconstruction Error

$$\hat{x} = [0.68, 0.62], x = [1, 0]$$

Loss function used  $\rightarrow$  Binary Cross Entropy

$$L = -\frac{1}{n} \sum [x \log \hat{x} + (1-x) \log (1-\hat{x})]$$

$$= -x \log \hat{x}$$

Mam's steps that I'm unsure about:-

$$-x \log \hat{x} = -1 \cdot \log(0.68) = \underline{0.136}$$

$$-(1-0.62) \log(1-0.62) = \underline{0.257}$$

Ignore this and calculate in your own way to check with the ans. [Find yourself]

$$L = \frac{0.136 + 0.257}{2}$$

[Not sure of this process]

⇒ If input is given as Real values to the autoencoder, what is the Loss function is used?

ans → MSE (Mean Squared Error). Let:-

original i/p =  $x = [3, 5, 2]$  → 3 neurons

reconstructed o/p =  $\hat{x} = [2, 4, 3]$  / in i/p layer &

3 neurons in o/p

Then, MSE:-

$$= \frac{1}{3} [(3-2)^2 + (5-4)^2 + (2-3)^2]$$

$$= \frac{1}{3} \{ 1^2 + 1^2 + (-1)^2 \} = \underline{\underline{1}} \text{ (ans)}$$

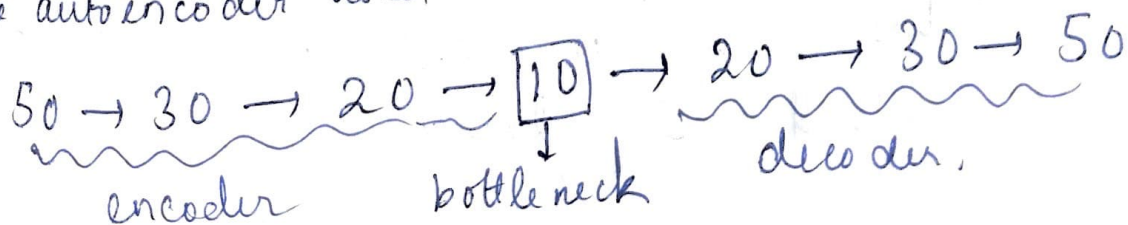
Now ⇒ Autoencoders use Unsupervised Learning Algorithm. why ~~not~~ so?

ans → We are NOT using any Target values or labeled data.

→ If we don't use any labeled data, then for a given dataset [let 50 features

100 samples. ⇒ 50 i/p neurons. let

the autoencoder used here ⇒ architecture is:-





Here, we have used a stack autoencoder (deep autoencoder) as we have more no.

of features and we need to find out / capture the inherent patterns or features in a multiple level of hierarchy: 1st hidden layer → finds out primitive patterns

2nd hidden layer → less complex patterns  
\* data is more complex. ]

↓

[ ] → 50 ifp neurons  
(for one sample at a time)  
100 samples in the batch

objective → reconstruct the features, when we feed entire data at a time. No

labeled data required or true target attr. value. [no need to check if true target value = predicted value. No need.]

→ All data samples is used for training purpose

↓

This is why we call autoencoder as unsupervised learning.



a. How can an autoencoder be used for classification task?

ans → objective of an autoencoder → reconstruct the data. Such that  $\hat{x} \approx x$

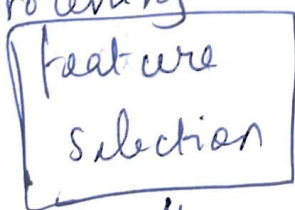
[if cannot reconstruct the data 100% accurately as it loses certain amt. of data during compression.]

Still, this compressed data (bottleneck) has all the inherent (inherent) patterns of all the i/p data so that we can reconstruct that particular input approximately equal to i/p given to the autoencoder]

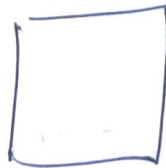
\* which part of autoencoder is used as part of classification task?

classification task: feature selection → classifying the data.

↳ classification preprocessing



training



classify on unknown data to its true class.

↓  
reduces the dimension of data



$n \times m$



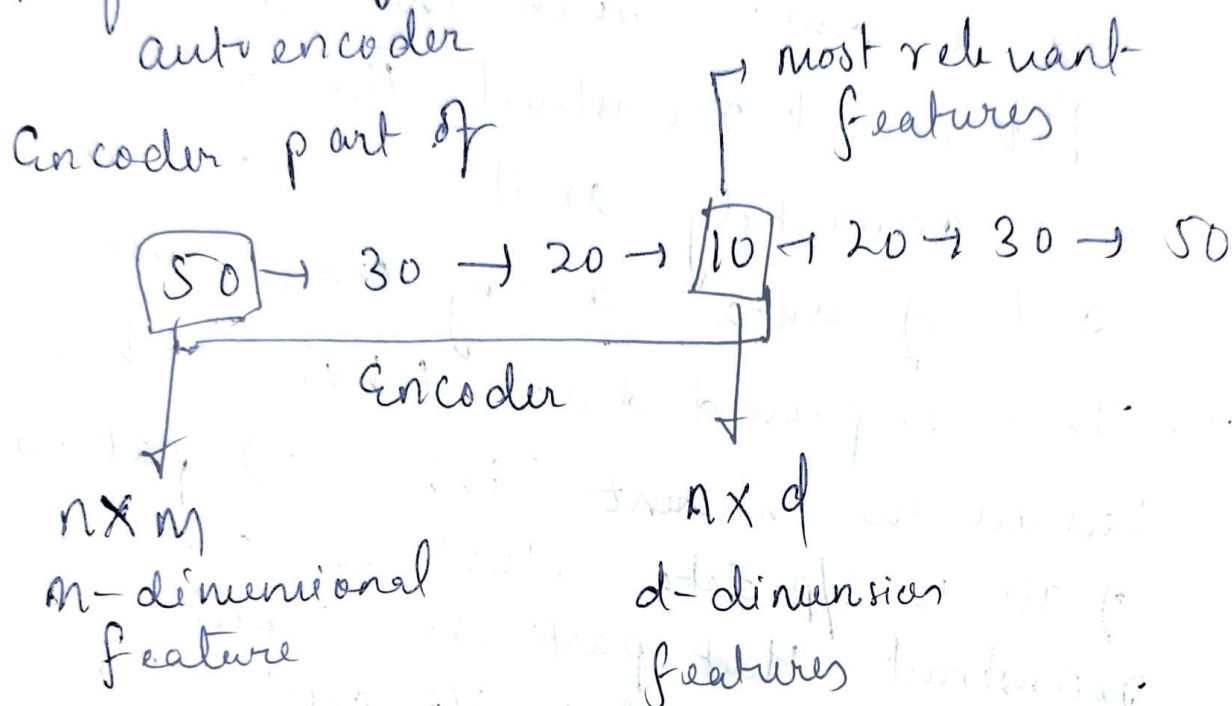
feature selection



$n \times d$

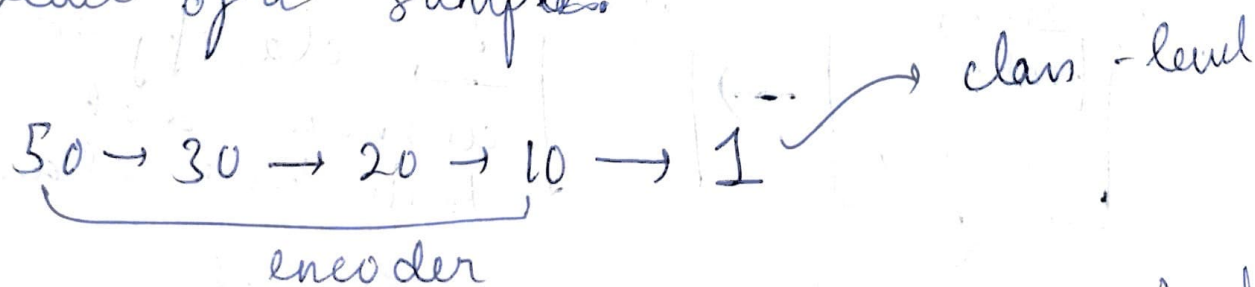
then  $d < m$

Now, we observe the feature selection part of classification is already ~~per~~ being performed by the encoder part of auto-encoder



bottleneck  $\rightarrow$  MOST Relevant features as features we can reconstruct the entire i/p data almost accurately from these features.

Now, encoder part taken individually and we add to it a single neuron (let) for classification purpose and apply an activation function  $\rightarrow$  it can give the class level of a sample.



\* Encoder is used for classification task.

Autoencoder can be used for :-

① Image reconstruction

② Text reconstruction

[We cannot use vanilla autoencoder as it uses ANN]. We use variants of autoencoder.

Encoder & decoder architecture

Image reconstruction →

Text Reconstruction →

CNN

RNN/LSTM/GRU

Same training process for CNN for  
Encoder & Decoder

—X—