

28/10/22

Autoencoders.

⇒ No target variable
→ unsupervised learning

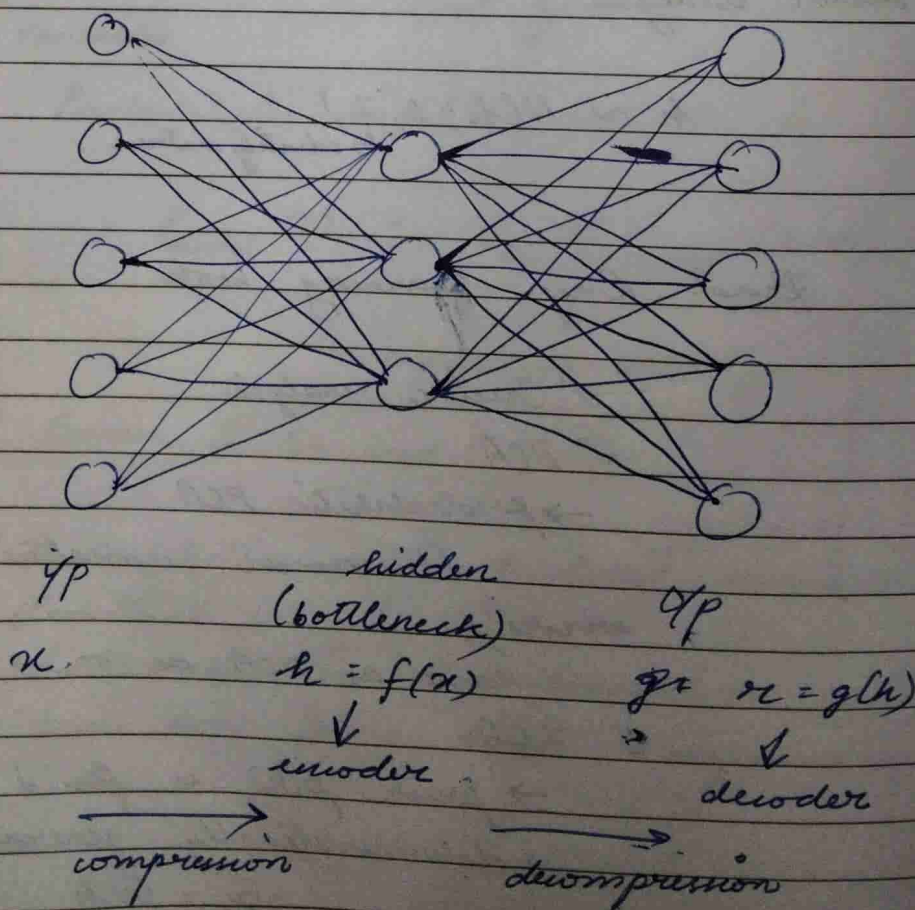
⇒ Tries to reconstruct the input
↓
image, text, ...

Input is $x \in \mathbb{R}^D$

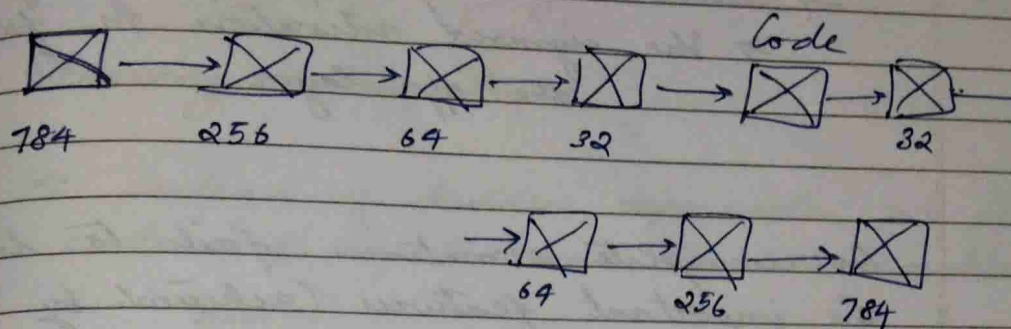
⇒ Output is also $x \in \mathbb{R}^D$
(same dimensions)

stands for reconstruction

Compression followed by decompression
(Input layer → hidden layer) (hidden layer \rightarrow output layer)



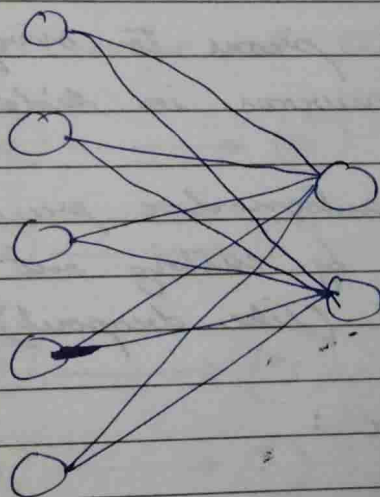
Decoder network \Rightarrow reverse of encoder



Undo

Undercomplete autoencoder :

Bottleneck has smaller dimension than ip layer

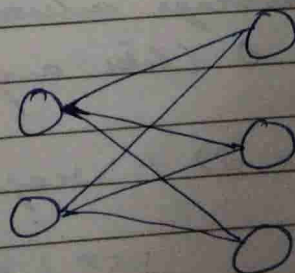


Overcomplete autoencoder :

Bottleneck dimension

\geq input dimension

(or) previous layer's dimension?



Output - binary

- ⇒ autoencoder o/p must be binary
- ⇒ Use sigmoid activation fn. for the o/p layer

Overcomplete sometimes fails to learn the important features (achieved by compression) because it only as good as simply copying the given image (not learning much).

- ⇒ can be useful if there are more latent variables

Overcomplete prone to overfitting as it uses more neurons in hidden units.

Regularized autoencoder prevents this overfitting by cutting out certain hidden units (like dropout).

Sparse autoencoder :

Here we penalize the activation of hidden layers (instead of weights as in regularization)

- ⇒ only few nodes will be active

Assumption : average activation is nearly 0 (like a sparse matrix)

Objective fn = $L(x, \hat{x}) + \text{regularization} +$
loss

hyperparameter $\leftarrow \lambda \sum_i |a_i^{(h)}|$

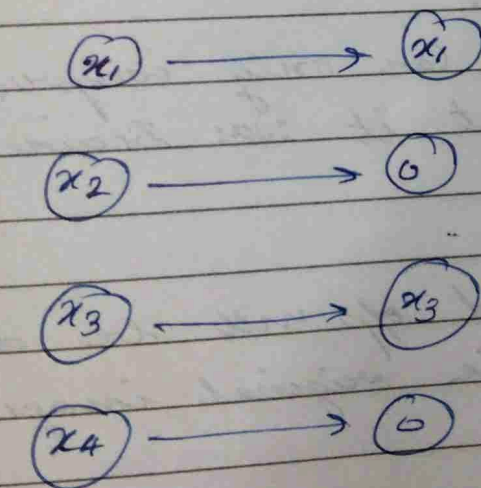
So we try to force $\sum |a_i^{(h)}|$ close to 0 to minimize the objective fn.

Denoising autoencoder:

→ Instead of hidden ~~inp~~ nodes, we force certain ϕ input nodes to be inactive

^{similar to}
Idea: ~~passing~~ noisy image as input
(achieved by ~~by~~ certain input features) ~~such~~
that ~~noises are filtered out~~
Output image - free from noise

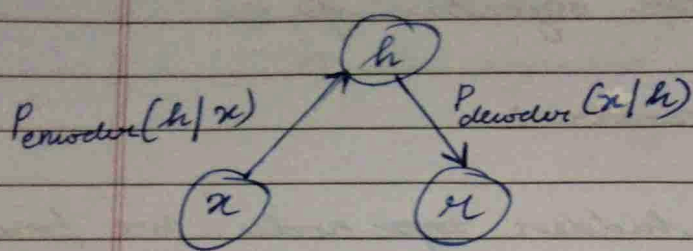
Loss computed between output and original image not the noisy image



orig ip

noisy inp
after dropping
50% of the
ip nodes.

Stochastic Encoders and Decoders



Was a probabilistic approach
encoding function $f(x)$ is replaced
with a probability distribution
 $P_{\text{encoder}}(h/x)$.

Similarly for $g(h)$ also $\rightarrow P_{\text{decoder}}(x'/h)$.

Properties of autoencoder

1) Data-specific :

Can only compress data
similar to what it was trained on

2) Lossy :

decompressed app will be degraded
compared to the original image

3)

Applications :

- ✓ Anomaly detection
- ✓ Object detection segmentation
(UNet, VNet, TransNet all use
autoencoders)

29/10/22

classmate

Date _____

Page _____

1 LSTM neuron ~~with~~ [^] \rightarrow 4 weights

utf-8 encoding \Rightarrow plaintext input

Preprocessing steps (for text)

1. Convert to lowercase
2. Remove numbers
3. Map every char to int

W_c

\hookrightarrow candidate

(new candidate value to be stored in memory).

RNNSequence to sequence
processing

⇒ there is dependency

No parallel computation

so slower

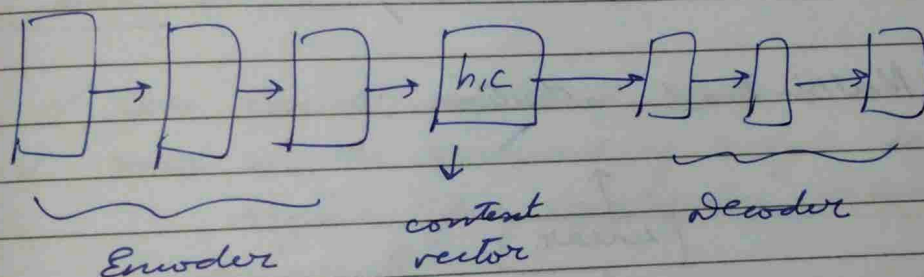
⇒ Difficult to process
longer sequences

Language translation:

TransformerProcesses all sequence
at once(not word by word)
⇒ there is no dependency

Parallel computation

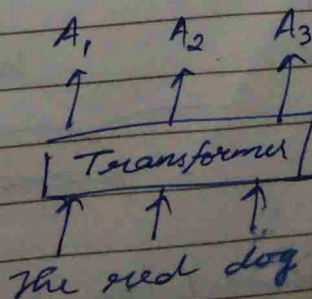
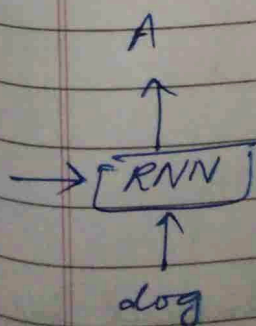
Faster

Easier to process long
sequences (as there is
no dependency)

Encoder-decoder models ⇒ self-supervised

Transformer uses 2 new layers

- ✓ Attention
- ✓ Positional encoder - captures content info



1. Attention module

Need to preserve semantics

(Eg): How are you?

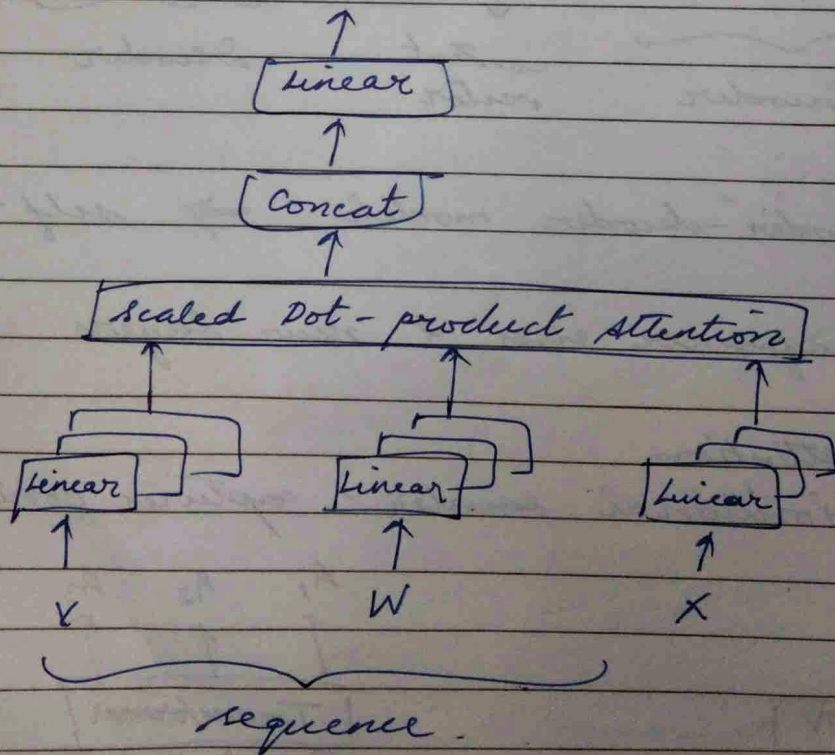
How - ~~is~~ ~~are~~ ~~you~~
are - ~~is~~
you - ~~is~~

But after putting them together the sentence becomes

~~is~~ ~~are~~ ~~you~~

(These semantics are preserved by using Multi-head attention).

Multi-head attention:



Positional encoding adds contextual info to ~~embed~~ a ~~sequence~~ embedding word

Fire tune - calculate error & backpropagate

Transfer learning :

2 steps :

1. Pre-train a model
2. Fire tune

Transfer learning for NLP :

✓ BERT → Bidirectional Encoder Representations from Transformers

✓ GPT

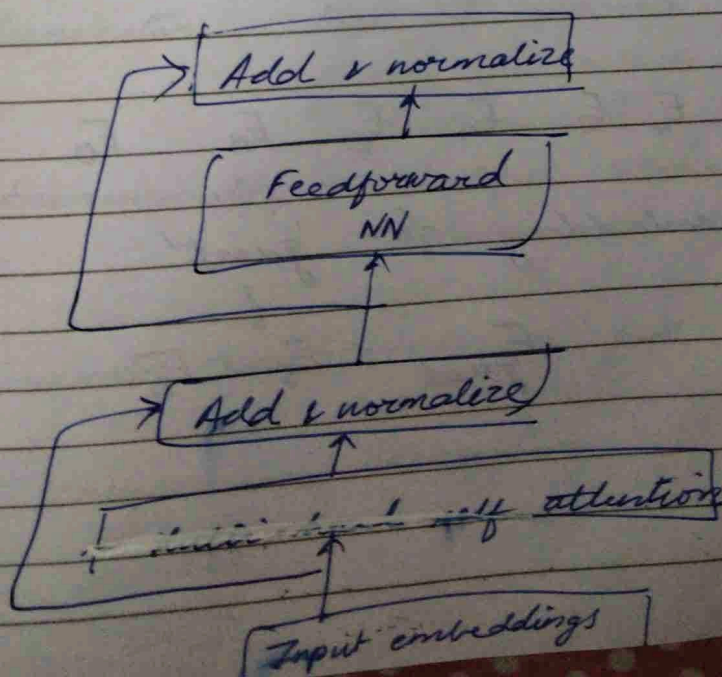
✓

BERT :

✓ Can handle sentiment analysis

✓ used in the backend of Google search engine for identifying relevant documents

✓ similar questions displayed on Quora



SEP \rightarrow used to separate sequences

CLS \rightarrow classification
used at beginning of input
to identify the task

(Ex): CLS My cat is great SEP she loves food SEP

BERT uses 3 ^{input} embeddings

✓ Token
✓ Position
✓ Sentence
} Input embedding

CLS My cat is great SEP she loves food SEP

Token E_{CLS} E_{My} E_{cat} E_{is} E_{great} E_{SEP} E_{she} E_{loves} E_{food} E_{SEP}
 \downarrow
used to
id the
tokens

Position E_1 E_2 E_3 E_4 E_5 E_6 E_7 E_8 E_9 E_{10}
 \downarrow
used to
id
position of
tokens

Sentence E_A E_A E_A E_A E_A E_A E_B E_B E_B E_B

Input embedding for 'great'
 \downarrow
 $E_{great} + E_5 + E_A$

BERT Pre-training

1. Predict the mask input ↑ allows BERT to better understand words.
2. Check whether subsequent sentences follow the previous sentence or not

My cat is great. It will rain today

completely unrelated to sentence 1.

Image data - VGN transformer

14/11/22

Generative Adversarial Networks (GANs)

like autoencoder

⇒ generates i/p again at o/p

2 neural networks

→ Generator

→ Discriminator

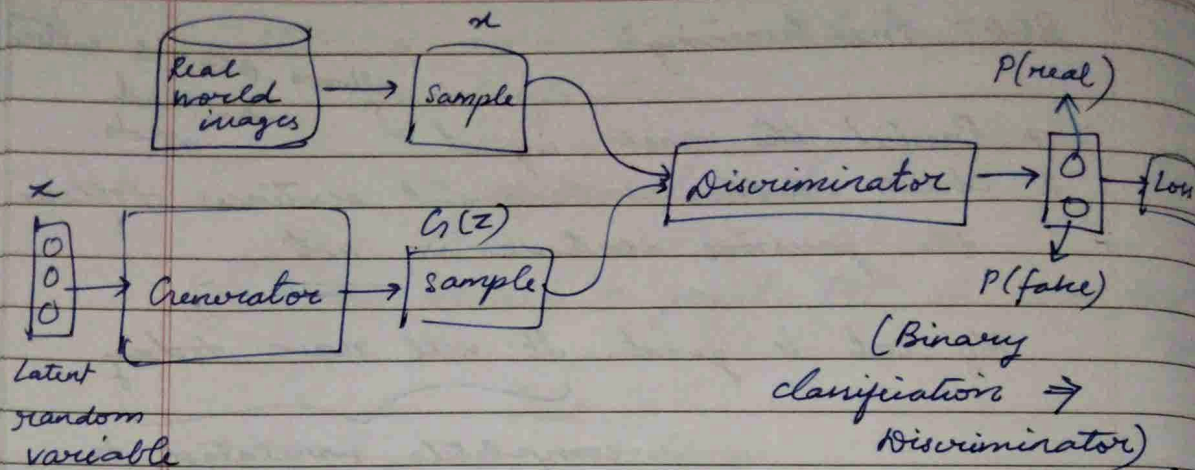
} zero-sum game
(one wins, other ~~loses~~ loses)

Generator - tries to generate fake i/p from i/p which is passed to discriminator

Discriminator - identifies which inputs are fake

Generator exists only to fool the discriminator

⇒ Unsupervised learning



$Z \rightarrow$ random noise (Gaussian/uniform).

Loss

1. If generator's output is correctly classified then by discriminator then generator has to be improved (update by backprop) ∇_v but generator's weights are frozen discriminator's

2. If generator's output is incorrectly classified by discriminator then discriminator has to be improved (generator's weights - kept const \Rightarrow not updated).

2 types of loss for:

Discriminator \Rightarrow Cross-entropy loss

$$\min_G \max_D V(D, G) = E_x [\log(D(x))] + E_z [\log(1 - D(G(z)))]$$

$$G(z) = \begin{cases} 1 & \text{if image is generated by generator} \\ 0 & \text{otherwise} \end{cases}$$

$$D(G(z)) \rightarrow 0 \text{ if } G(z) = 1$$

(classify as fake if img is generated by generator).

Correctly
classifying
wrong
data

$$D(x) \rightarrow 1$$

Correctly
classifying
real
data

Objective of generator is to minimize the

Objective of discriminator is to maximize

