

Lect 33: Attention Model & Transformers

review of previous lecture RNN & its types

can be

The Seq2Seq Model

model with RNNs / Seq2Seq

model in NLP used to
convert seq of type A to
seq of type B

early

classifiers

seq \rightarrow label

auto-encoder (Sequence - Sequence)

earlier with fully connected layer

now we can use different networks of our-choice

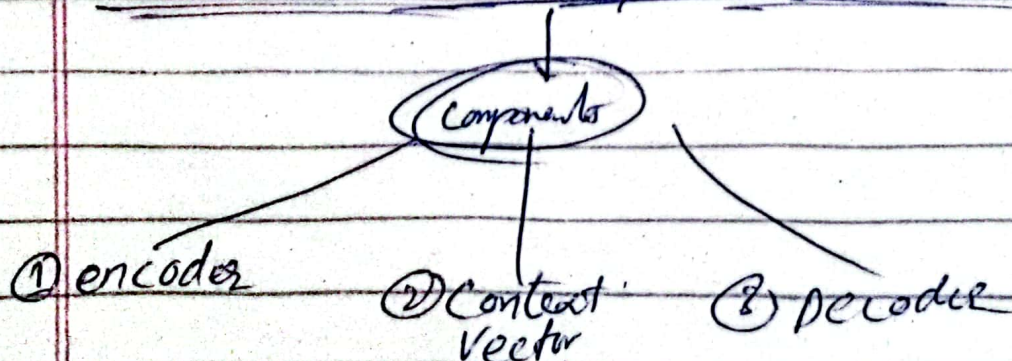
auto-encoder

Encoder/decoder model (will ^{be} using RNN)

example

(i) translation of english sentence to Spanish

Encoder-decoder Network / Seq2Seq Model Architecture



encoders are good at storing represent of any info. it could be audio/video/image/text

encoder → Takes input as sequence and output as fixed representation n-d
• read variable-length input

Decoder → Take context vector / ^{latent representation} hidden representation and give sequence as output.

where this used

(i) Language related tasks

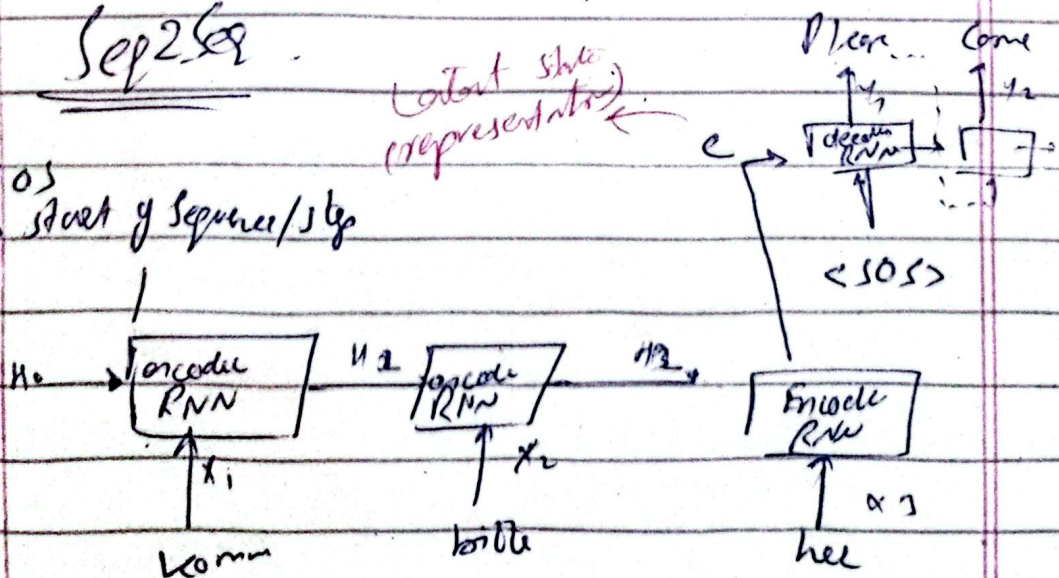
- summarization
- Translation
- Question Answering

Context Vector

- (fixed dimension representation of input sequence)
- hidden state which is a vector that represents the learned features
- output of encoder

Seq2Seq

∴ SOS
↳ start of sequence/stp



- a, b as input (different of sequence) - (different time steps)

information (input) + summary \Rightarrow output

Cell can be either vector, LSTM. It can be any.

input $\Rightarrow x_1, x_2, x_3$ \rightarrow variable set
output $\Rightarrow y_1, y_2, y_3$ \rightarrow variable set

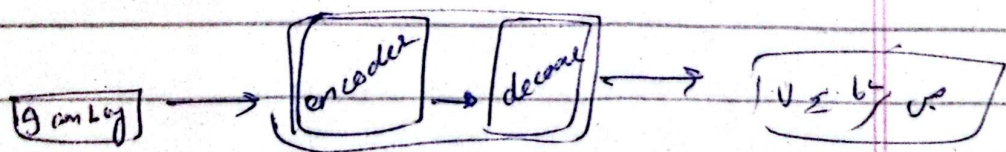
applications & popularity

① Sentiment Analysis

② Question Answering

③ Machine Translation

\Rightarrow Translation



Neural Machine Translation

Seq-2-Seq Model

\rightarrow it learns this due to training.

Train together

& then they can work independently.

Problems with RNN

Time step

① They process input data sequentially one after the other.

So not make use of modern GPU, which were designed for parallel computation. So slow training

② Long-term dependency problem

یعنی ایک نقطہ اور (sequence) کے شروع میں آیا ہے یا نہ

(Sequence) کے آخر میں آیا ہے یا نہ (relationship) کے ساتھ

- quite ineffective when elements are distant from one another.

- due to fact, information is passed at each step

- longer the chain is, the more probable the info gets lost along chain.

③ bottleneck problem

- due to fixed-size context vector

- it may fail to capture all relevant details leading to information loss.

- So can not adequately represent long-range dependencies within input seq.

کچھ نہ (fixed dimens) میں (features) کے ساتھ

One = 1

Problems with Traditional learning model

↳ struggle with remembering long stories or instructions

example

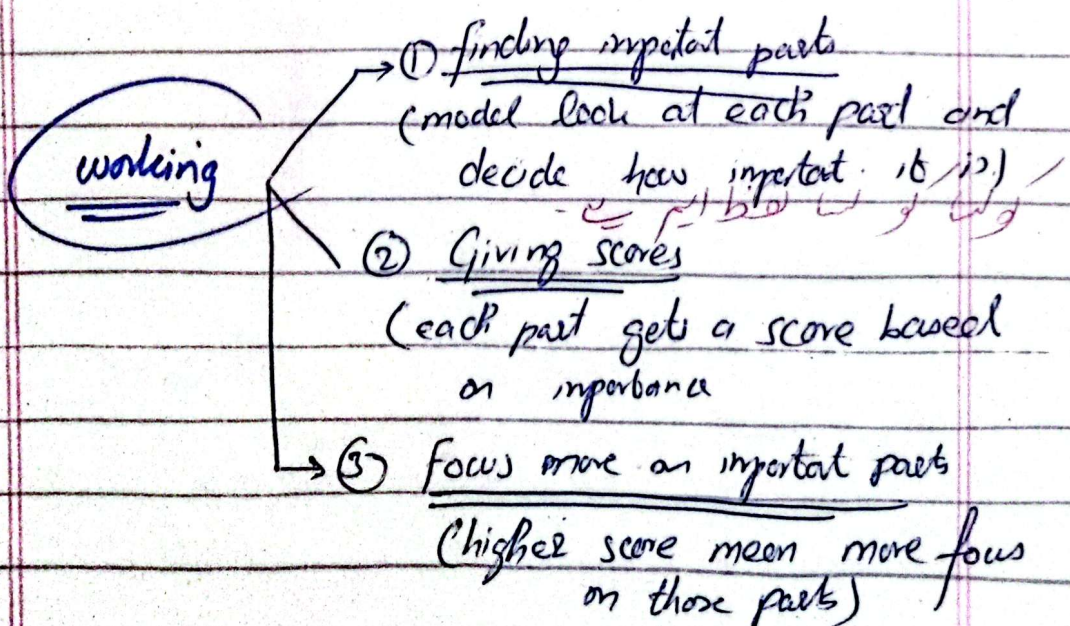
reading a long book & forgetting key details

Concept of Attention

like highlighter - it highlights important information and then generates output considering that highlighter (Attention)

Such similar seq-seq model are called attention models

$w' b' w' d' \rightarrow (\text{importance}) w' z, y'$
 $- e^{-\beta}$



Attention in action

Translating a sentence

with attention: Computer focus on key words for better translation

Example

"The Professor assigned homework on differential equations."

incorrect translation without attention:

ustad nay forg kay eqvator par homework dya

correct translation with attention

Professor ne differential eqvator par homework dy.

Attention

focus on important parts

helps model make better decisions

Can be applied b/w
different input &
different part of
same input

Solve the problem gNNs or LSTM

$(x \in \mathbb{R}^n, y \in \mathbb{R}^m)$ is $\{x' \in \mathbb{R}^n\}$ (decide) $\text{at } \omega$
(Attention)

- all neural network use weights.
- weights are random at start.

- when we train the model weights are trained.

- So they update values

$z \rightarrow vt$

- allerta weights are random in start

- during train time the weights of allerta
