

Week 3.

Various Sequence to Sequence Architectures

Video 1: Basic Models

$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$ $x^{(5)}$
 Jane visite l'Afrique en septembre

$y^{(1)}$ $y^{(2)}$ $y^{(3)}$ $y^{(4)}$ $y^{(5)}$ $y^{(6)}$
 Jane is visiting Africa in September

Many to many models are used in machine translation & speech recognition

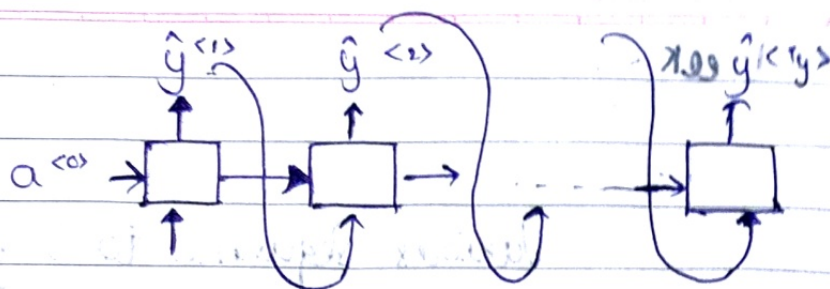
encoder - RNN/LSTM/GRU

~~inputs~~ ~~the~~ takes as input the whole sentence and outputs a vector that should represent the whole input.

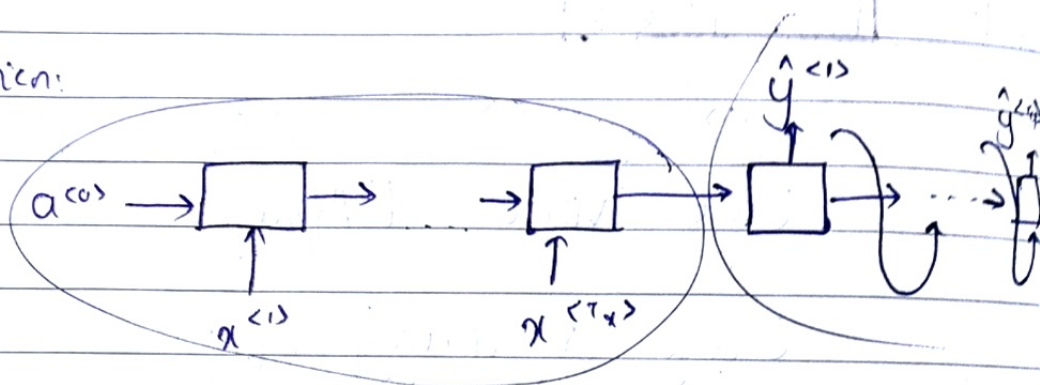
decoder - takes in sequence built by encoder and outputs new sequence

Video 2: Picking the most likely sentence

Language Model:



Machine Translation:



language model helps us estimate the probability of a sentence & generate sentences too

$$P(y^{(1)}, \dots, y^{(T_y)})$$

Machine Translation:

$$P(y^{(1)}, \dots, y^{(T_y)} | x^{(1)}, \dots, x^{(T_x)})$$

Example:

x : "Jane visite l'Afrique en Septembre"

y may be:

- Jane is visiting Africa in September
- Jane is going to visit Africa in September
- In September Jane will visit Africa.

So for best output,

$$\underset{y^{(1)}, \dots, y^{(T_y)}}{\operatorname{argmax}} (P(y^{(1)}, \dots, y^{(T_y)} | x))$$

Greedy Search

first word - most likely first fr word according to conditional language model.

↓

second most likely

↓

third most likely

X working

example "Jane is" most likely beginning
after is going is most common word
∴ "Jane is going to be visiting Africa in
September". $P(\text{Jane is going} | x) > P(\text{Jane is visiting} | x)$

~~Jane is~~ In September, Jane will visit Africa.

Her African friend welcomed Jane in September

Video 3: Beam Search

$B = 3$

↑

beam width

input french sentence → encoder → ~~softmax~~ decoder
~~softmax~~ top 3 possibilities ← Softmax output
 (10000 possibilities)

In } most likely beginning words
 June
 September

Then for each of the above words the most likely 2nd word will be generated according to $P(z)$.
 ∴ for 'In' there will 3 most likely upcoming words & so on for the rest.

Out of these again top 3 best possibilities are selected.

↳ so on

Video 4: Refinements in Beam Search

length normalisation

$$\arg \max \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

$$\begin{aligned} \text{i.e. } & P(y^{<1>} \dots y^{<T_y>} | x) \\ &= P(y^{<1>} | x) \cdot P(y^{<2>} | x, y^{<1>}) \end{aligned}$$

Probabilities are too small

∴ we can log

$$\arg \max \sum_{y=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

$$\frac{1}{T_y} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

$\alpha = 1$ (full normalisation)

$\alpha = 0$ (no normalisation)

$\alpha = 0.7 \Leftarrow$

$B \uparrow$, time \uparrow

Videos: Error Analysis in Beam Search

To segregate if error is due to B / RNN

$x =$ 'Jane visite l'Afrique en Septembre'

$y^* =$ "Jane visits Africa in September" \Leftarrow

$\hat{y} =$ "Jane visited Africa last September" \Leftarrow

Case ①

$$P(y^* | x) > P(\hat{y} | x)$$

beam search chose \hat{y} although y^* has higher $P(y|x)$

beam search is at fault.

Case ②

$$P(y^* | x) \leq P(\hat{y} | x)$$

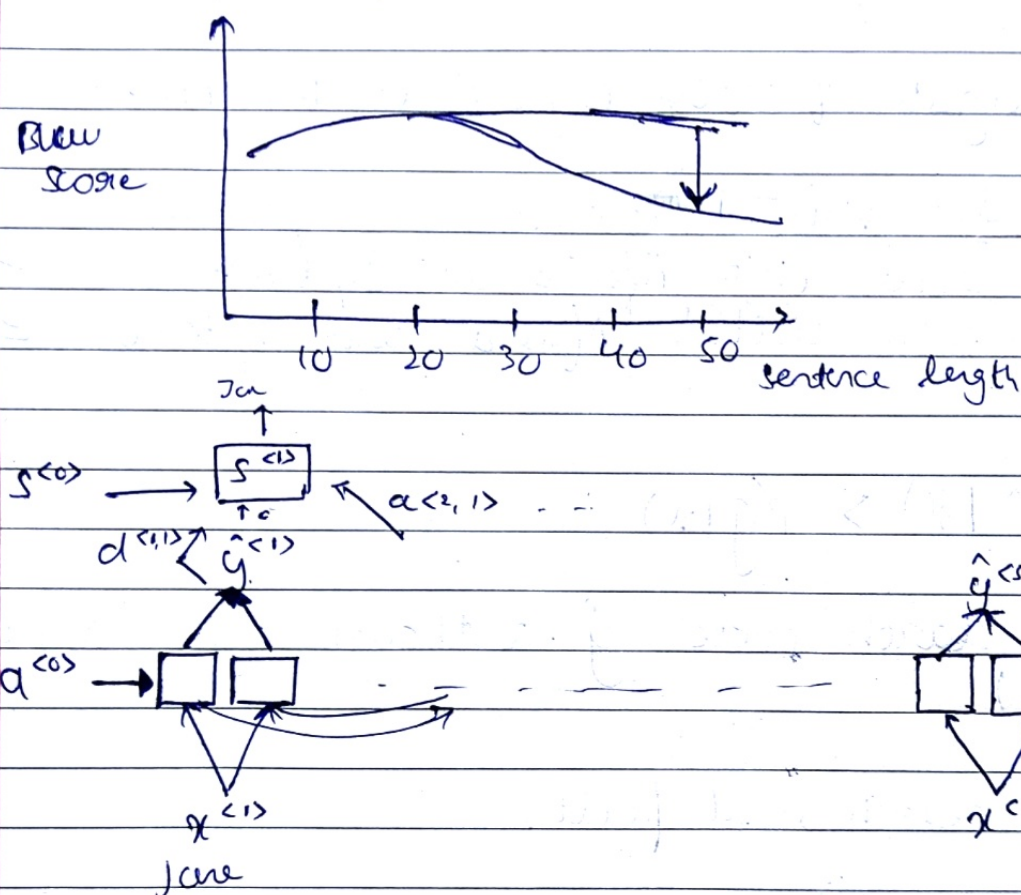
RNN predicted $P(\hat{y} | x) > P(y^* | x)$

RNN - faulty

Video 6: Attention Model Intuition

Given a long sentence and the job of translation the humans wouldn't memorise the sentence and translate it but they would translate in parts.

Model's performance ↓ if sentence is long



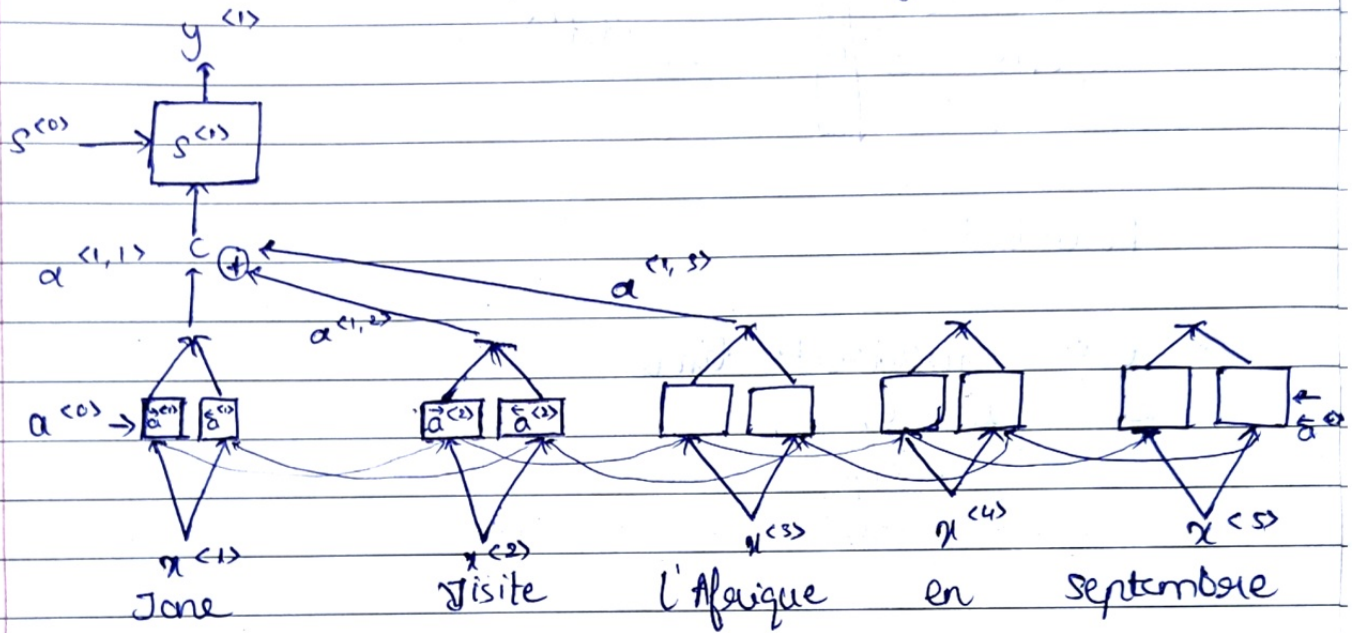
encoder is bidirectional RNN

it generates a vector representing inputs.

Video 7:

Attention Model

- * input sentence is processed by bidirectional RNN/GRU/LSTM to compute features on every word.



- * α tells us how much the context would depend on the features we are getting or the activations we are getting from different time steps.

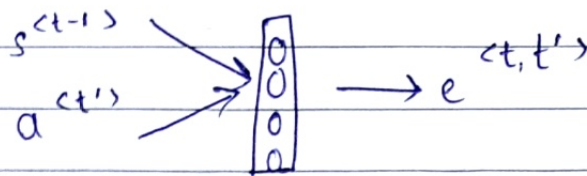
$$\alpha^{(t,t')} = [\vec{a}^{(t,t')}, \bar{a}^{(t,t')}]$$

$$\sum_{t'} \alpha^{(t,t')} = 1$$

$$c^{(1)} = \sum_{t'} \underset{\substack{\uparrow \\ \text{weights}}}{\alpha^{(1,t')}} \underset{\substack{\uparrow \\ \text{activation}}}{a^{(t')}} a^{(t')}$$

$\alpha^{(t,t')}$ = Amount of attention $y^{(t)}$ should pay to $a^{(t')}$

$$d(t, t') = \frac{\exp(e(t, t'))}{\sum_{t'=1}^T \exp(e(t, t'))} = 1$$



Drawback :

- takes quadratic time

eg. $7x - 4p$

$7y - 4p$

time $= 7x \cdot 7y$

Speech Recognition - Audio Data

Video Speech Recognition

