

Week 3 : Sequence Models &

Attention Mechanism

1] Sequence to sequence architectures

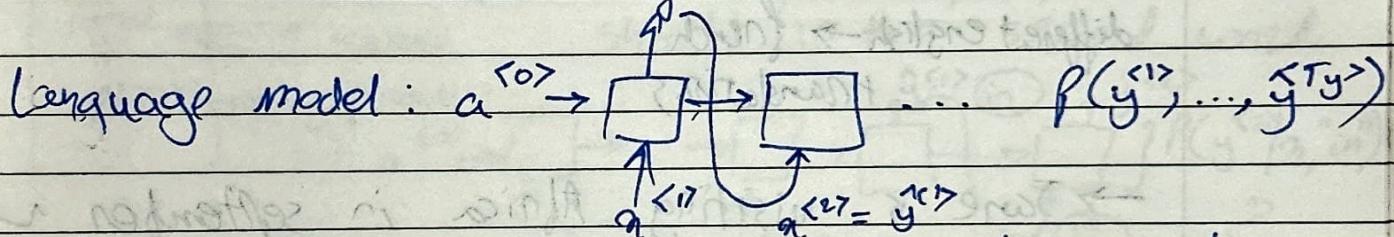
1) Basic Models

→ Machine translation & image Captioning

The input sentence is encoded using a NN (GRU or LSTM) then passed to a decoder which generates the output sentence. Similar methods using CNNs are applied for image captioning which will provide a short description of the input image.

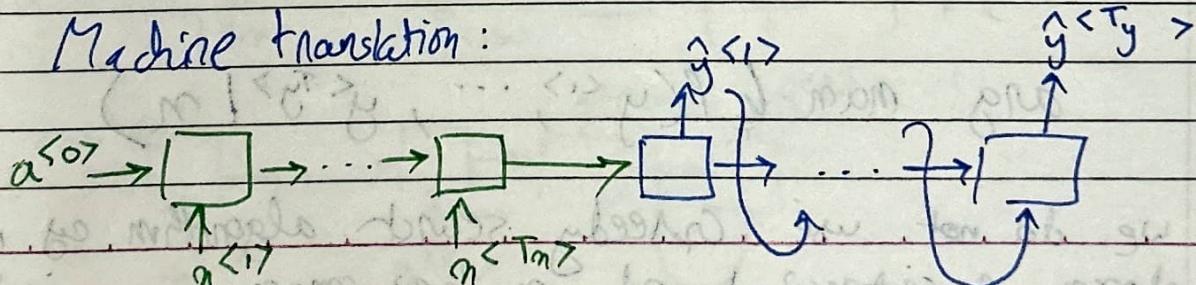
2) picking the Most likely sentence

Machine translation by building a conditional language model



Used to estimate the probability of a sentence or to generate novel sentences.

Machine translation:



Instead of vectors of zeroes that are passed to the NN, it has an encoded network that figures out some representation for the input sentence rather than with representation of all zeros.

→ it is "Conditional language model"

because it is modeling $P(y^{<1>}; \dots; y^{<Ty>} | n^{<1>}; \dots; n^{<Ty>})$

Probability of output English translation given (conditioned on) some french sentence.

Finding the most likely translation:

a) Jane visite l'Afrique en septembre

$$P(y^{<1>}; \dots; y^{<Ty>} | n)$$

different english → french
translations

→ Jane is visiting Africa in september ✓

→ Jane is going to be visiting Africa in september.

→ In September, Jane will visit Africa

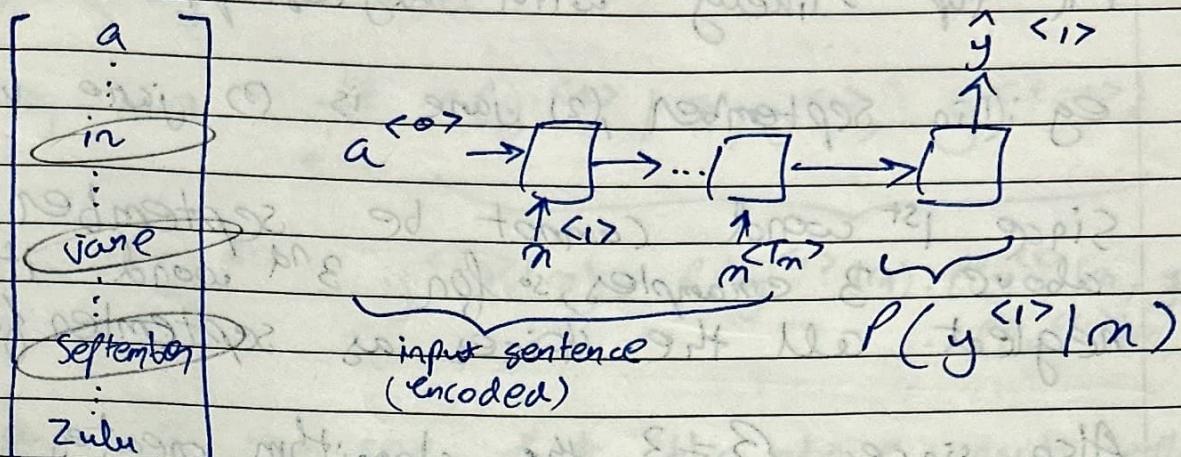
→ Her African friend welcomed Jane in September

$$\arg \max \{ P(y^{<1>}; \dots; y^{<Ty>} | n) \}$$

We do not use greedy search algorithm as it chooses a sentence based on how common a word is than appropriate

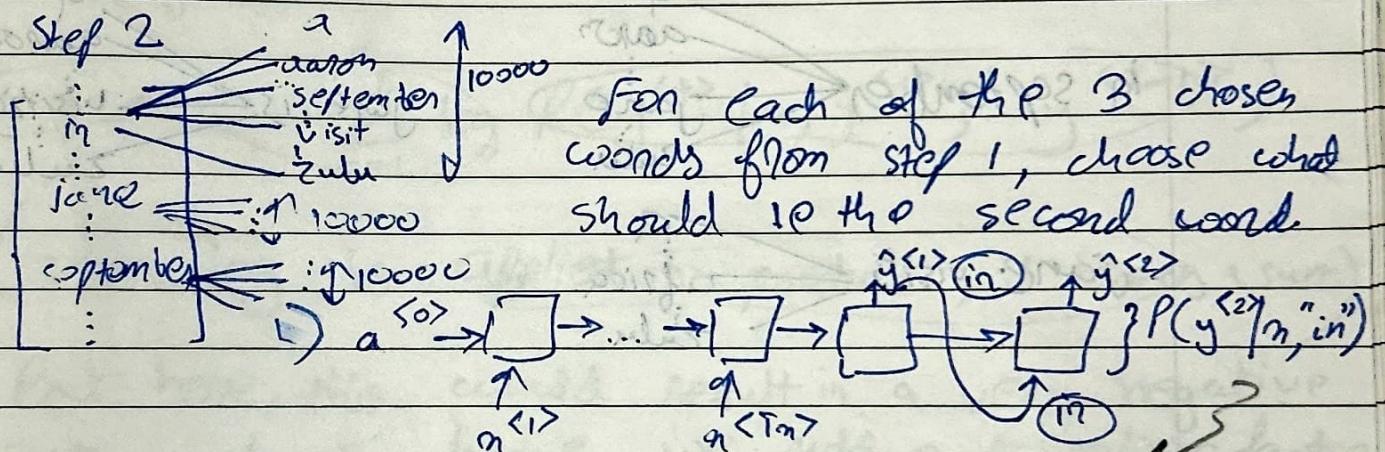
3) Beam Search $\rightarrow B$ (beam width)

Step 1

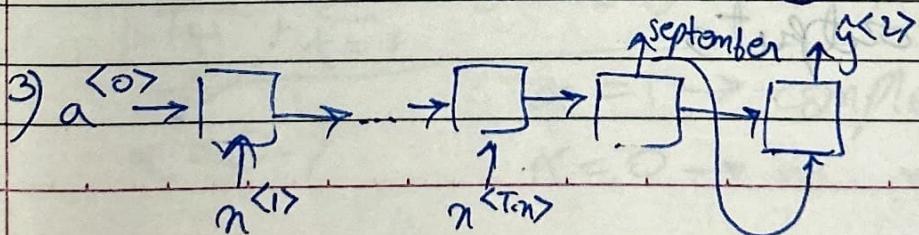
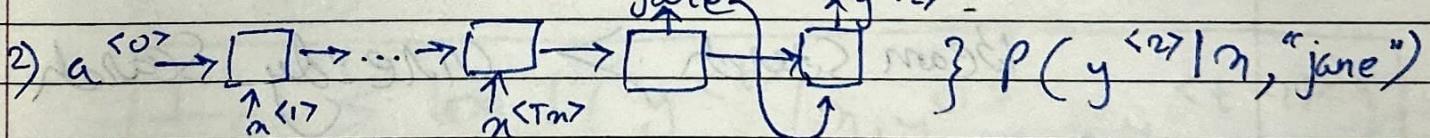


$\hat{y}^{<1>}$ is a softmax layer overall outputting 10000 possibilities and keep the top 3 words in the memory given that $B=3$.

Step 2



$$\rightarrow P(y^{<1>}, y^{<2>} | n) = P(y^{<1>} | n) * P(y^{<2>} | n, y^{<1>})$$



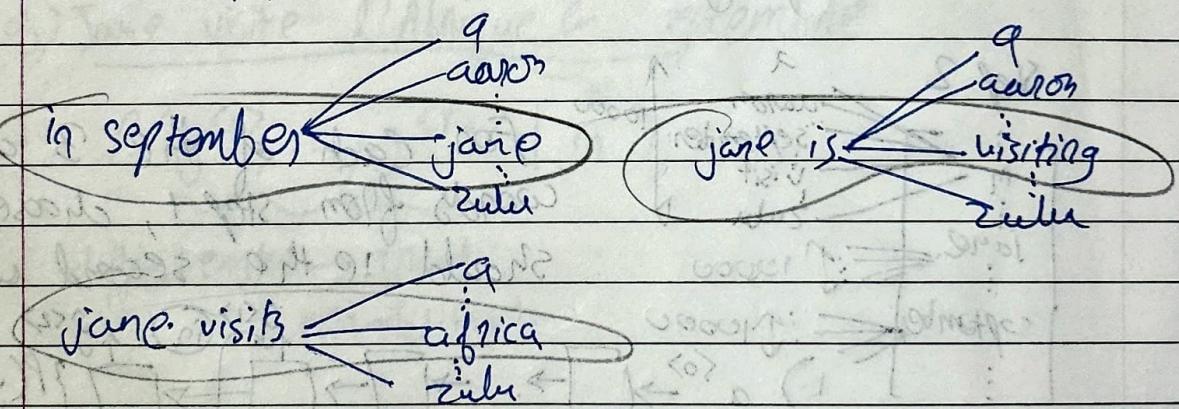
Total $3 \times 10000 = 30000$ possibilities.
 Evaluate all the 30000 options and pick the top 3 likely with highest probabilities.

Eg : (1) in september, (2) jane is (3) jane visits

Since 1st word cannot be september acc. to above 3 examples, so for 3rd word it will neglect all the choices as september for 1st word.

Also since $B=3$, the algorithm creates 3 copy of different $\hat{y}^{(1)}$ which are further used to predict the 2nd word.

Next Step,



Note : If $B=1$, this algo becomes a greedy search algo.

Beam Search > Greedy Search

→ Better outputs.

5) Refinements to Beam Search (Length Normalisation)

Previously we used:

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

$$P(y^{<1>} \dots y^{<T_y>} | m) = P(y^{<1>} | m) \cdot P(y^{<2>} | m, y^{<1>}) \cdot \dots \cdot P(y^{<T_y>} | m, y^{<1>} \dots y^{<T_y-1>})$$

as we can see that the above result will be a product of probabilities ~~here~~ and each prob. being < 1 will result in a number very close to zero ~~hence~~^{and since} the small value cannot be stored in the memory, we shall add log to the function.

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

(when log is applied to product, it becomes a sum)

But here, this would result in a very negative sum
1 hence we add a normalizing factor.

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

$\alpha = 1 \rightarrow$ completely normalized length
 $\alpha = 0.7$ $\alpha = 0 \rightarrow$ no normalization as $(T_y = 1)$

Beam width B :

- \rightarrow large B : better results, but slower (computationally)
also memory requirements \uparrow
- \rightarrow small B : worse results, (faster computationally)

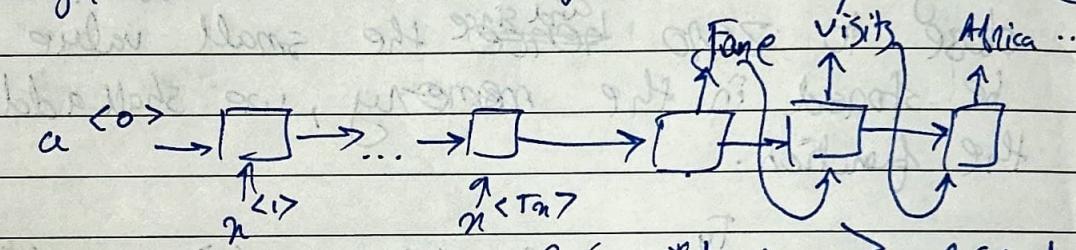
5) Error analysis in Beam Search

~~RNN computes $P(y|n)$~~

Jane visite l'Afrique en septembre.

Human: Jane visits off Africa in September. (y^*)

Algorithm: Jane visited Africa last September (\hat{y})



To see which one is bigger to clearly associate/identify and attribute a specific error to a cause on source within the process.

Case 1: $P(y^*|n) > P(\hat{y}|n)$

Beam search chose \hat{y} . But y^* attains higher $P(y|n)$
Conclusion: Beam search is at fault.

Case 2: $P(y^*|n) \leq P(\hat{y}|n)$

in ~~our~~ y^* is a better translation than \hat{y} . But RNN predicts $P(y^*|n) \leq P(\hat{y}|n)$

Conclusion: RNN model is at fault.

Error Analysis process

Human Algorithm $P(y^*|n)$ $P(\hat{y}|n)$ At fault?

$$2 \times 10^{-10} \quad 1 \times 10^{-10}$$

Beam

RNN

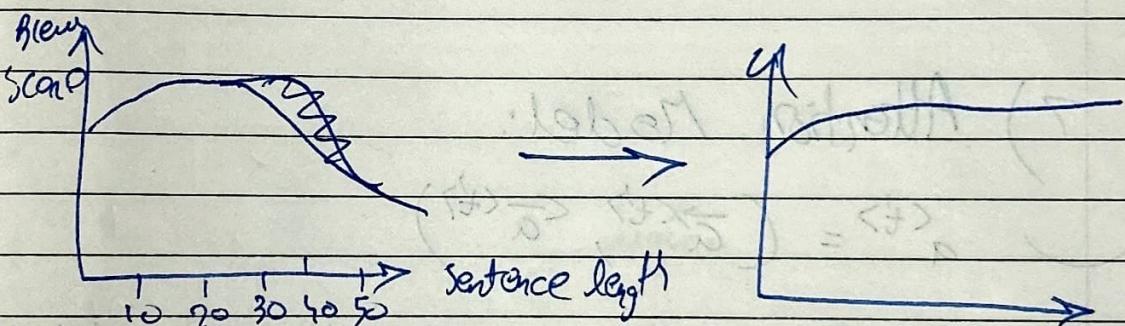
Beam

→ figure out what fraction of errors are "due to" beam search vs RNN model

i.e. beam search at fault

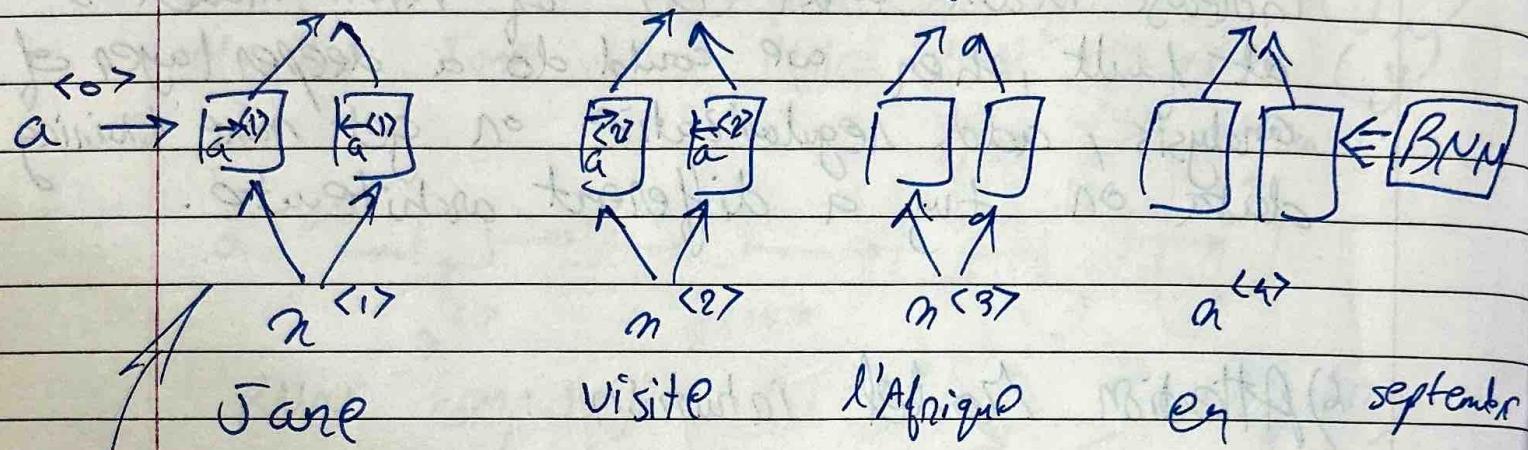
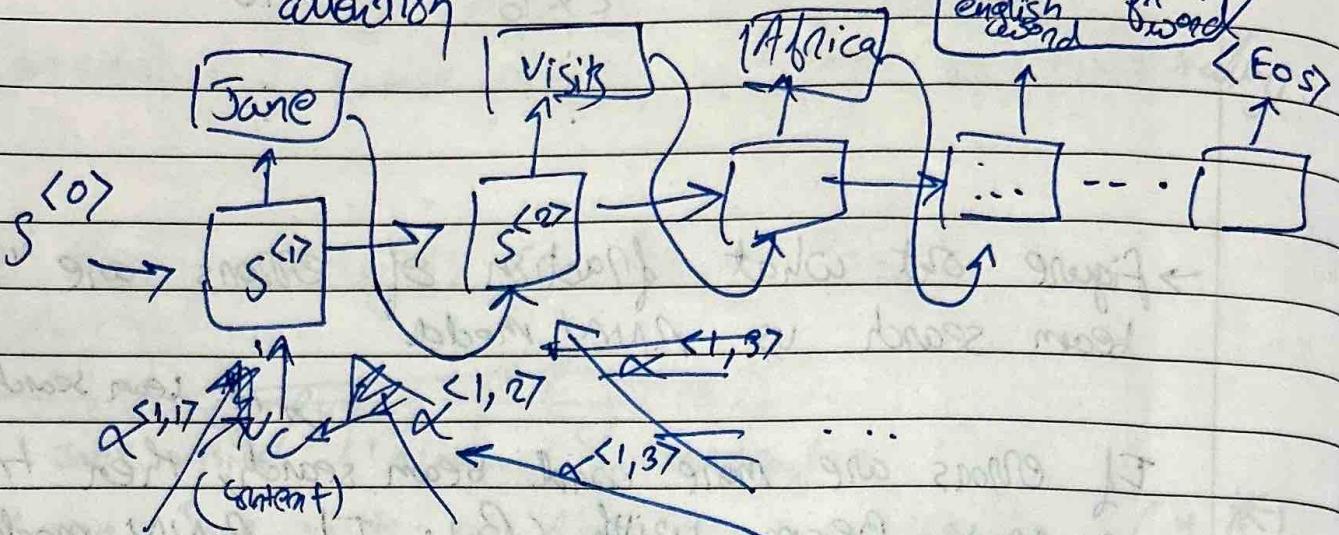
If errors are more w/ beam search, then try to increase Beam width (B). If RNN model is at fault, then we could do a deeper layer of analysis, add regularization or get more training data or try a different architecture.

b) Attention Model intuition:



This score goes down as the model cannot ~~not~~ remember sentences of high lengths at once. Hence, attention model is used which will translate \mathbb{P} the sentence part-by-part just like human translations do.

$\alpha \rightarrow$ weights
attention



7) Attention Model:

$$\alpha^{(t)} = (\bar{a}^{(t)}, \underline{a}^{(t)})$$

α tells us how much the context could depend on the features/activations we're getting from different time steps.

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

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

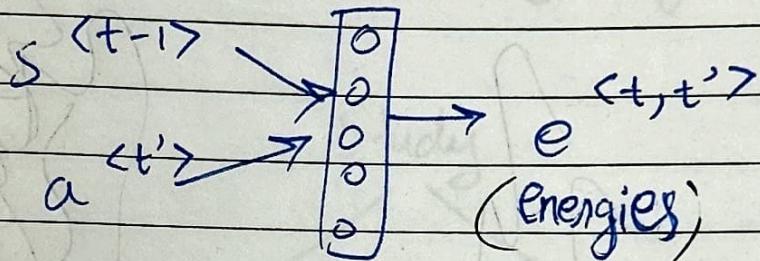
$\alpha^{(t, t')}$

$\alpha^{(t, t')}$ = amount of "attention" $y^{(t)}$ should pay
to $a^{(t')}$. i.e. importance of that input word.

 $\alpha^{(t, t')}$

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

Can be computed using a small NN:



$s^{(t-1)}$ \rightarrow hidden/state
of post Bi-RNN
 $a^{(t')}$ \rightarrow hidden state
of pre-attention

2] Trigger word detection

→ Alexa, Google Home