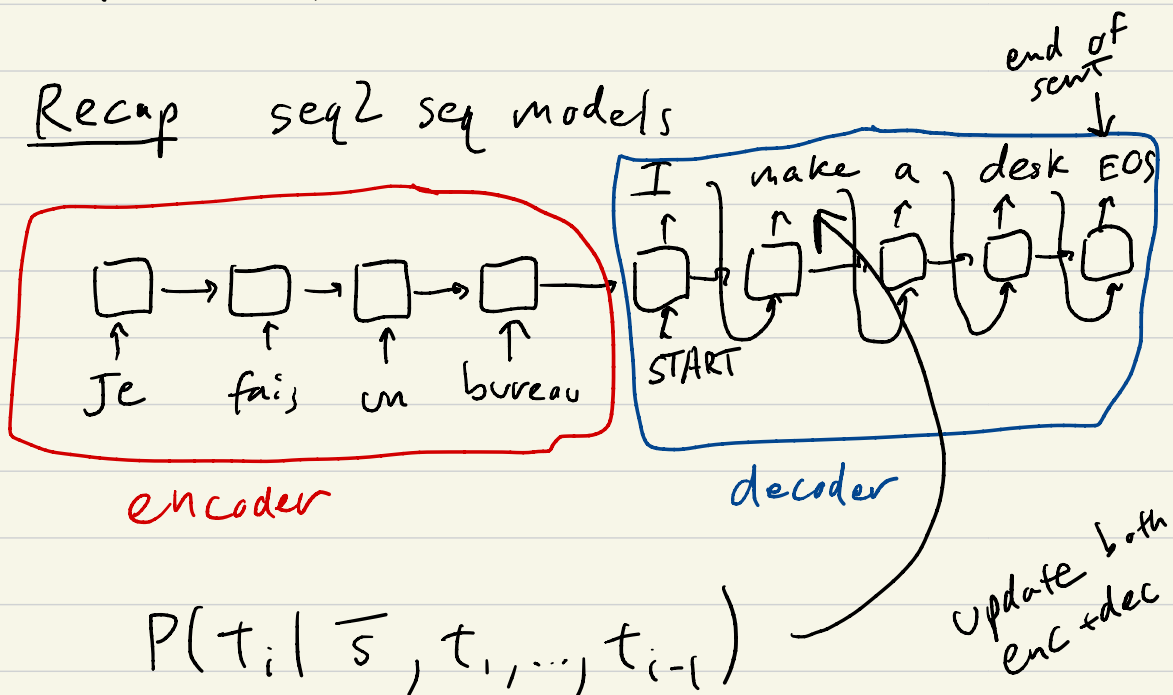


CS 378 Lecture 20

Today

- Attention in seq2seq models
- Neural machine translation

Recap seq2seq models



Training
teacher forcing

Assume everything correct through $i-1$, maximize log prob of word t_i

$$\text{loss} = \sum -\log P(t_i^* | \bar{s}, t_1^*, \dots, t_{i-1}^*)$$

Announcements

- A4 due
- A5 out tonight (due in 1 week)

Problems with seq2seq models

① Model repeats itself in a loop

Je fais ... \Rightarrow I make a desk a desk
a desk a desk ...

Why didn't this happen in phrase-based?
We had a notion of coverage!

RNN doesn't track "progress"

② Fixed vocabulary

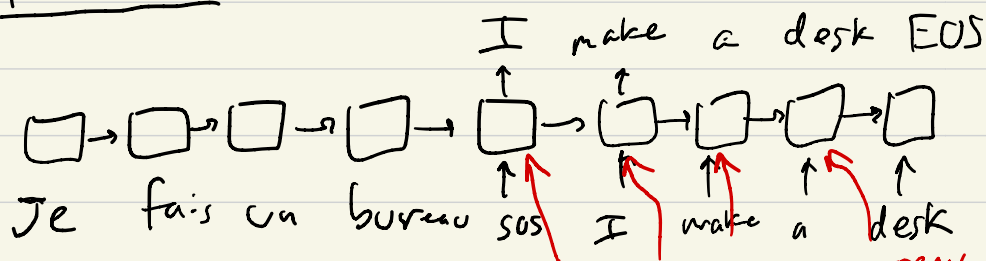
Elle est allée à Pont-de-Buis \Rightarrow She went to
UNK

③ Bad at long sentences



LSTMs have fixed hidden state, 50+ time steps is hard

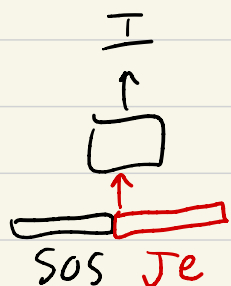
Attention



Requires modifying LSTM to take 2 inputs

Suppose it was always a word-by-word translation in order

Decoder:



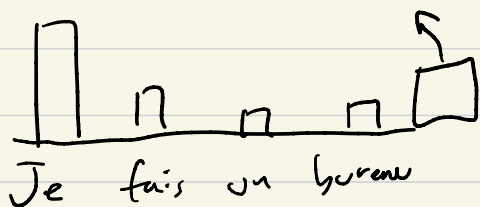
RNN just needs to map J_e to I

Can learn this more easily
+ with fewer params than a
normal seq2seq.

(could even delete encoder)

This is too rigid. Instead we want
the decoder to softly pick where it
looks in the input.

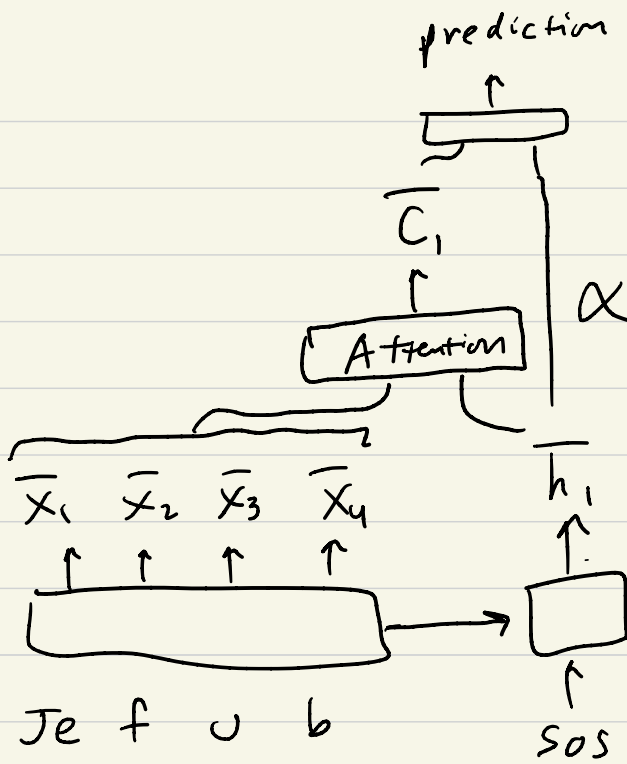
Attention:



distribution over input
positions

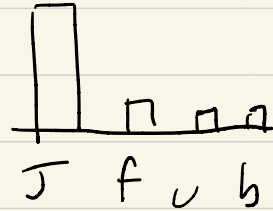
→ weighted average
of input

→ used for prediction



"loop over all inputs
 i , compute this
 score for \bar{x}_i , softmax"

$$\alpha = \text{softmax}_i(\bar{h}_i^T W \bar{x}_i)$$



form $\sum_{i=1}^n \alpha_i \bar{x}_i$

$$\bar{c}_i = \sum_{i=1}^n \alpha_i \bar{x}_i$$

$$P(\tau, \bar{s}) = \text{softmax}_{(\text{vocab})} \left(V \begin{bmatrix} \bar{h}_i \\ \bar{c}_i \end{bmatrix} \right)$$

vector
concat

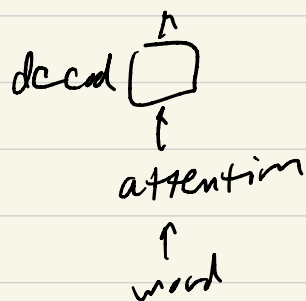
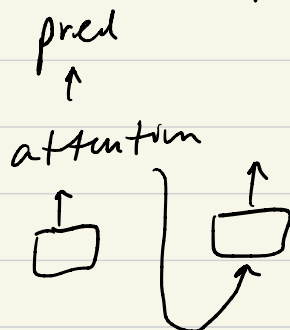
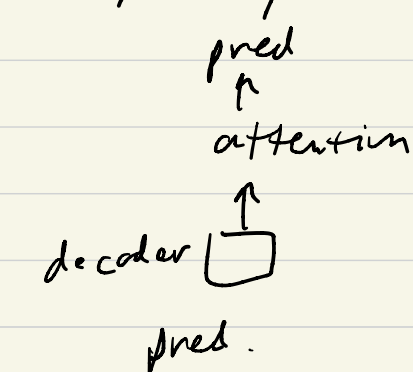
$$\bar{c}_i \approx "Je"$$

Do this at every timestep!

Training: with backprop

Details:

- Many ways to set this up



$$\alpha = \text{softmax}(f(\bar{h}_i, \bar{x}_i))$$

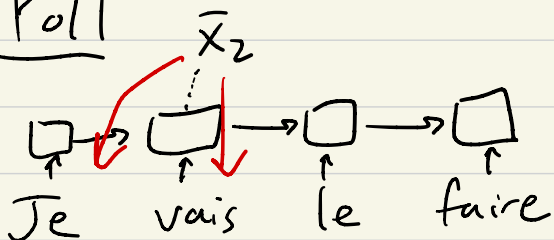
f : dot product,
 $\bar{h}_i^T W \bar{x}_i$

one-layer NN

⋮

Overall idea: form context vector \bar{c}
 that captures input directly

Poll \bar{x}_2



I am going to do it

↑ ? look back at the verb
look back at position 2 in the
(index) source

- ① am : look at v_{ais}
- ② look at Je (subject)
look at other stuff?
(other dependencies of the verb)
- ③ Encoder tracks: content

position

Decoder timestep 1: model wants to attend to enc-timestep 1