

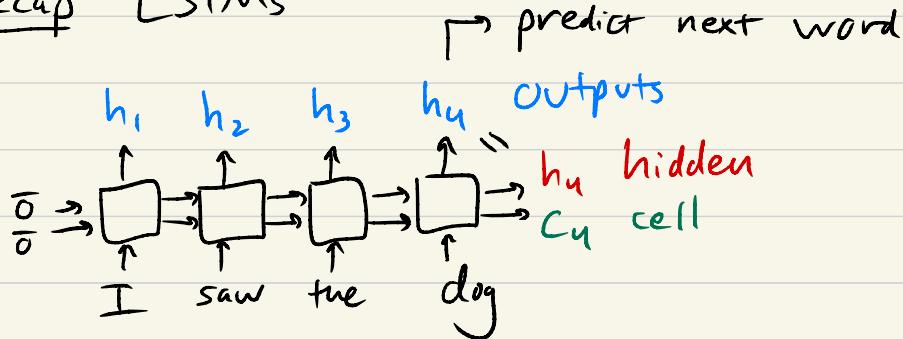
Announcements

- A4 due in 1 week



Star Wars The Third Gathers: The Backstroke of the West
(subtitles machine translated from Chinese)

Recap LSTMs



Machine Translation

Phrase-based MT: 2000 – 2015

Neural MT: 2015 – present

This lecture: PBMT and word alignment

Thursday: finish PBMT, start NMT: seq2seq models
→ leads to attention: critical idea in modern deep learning w/roots in word alignment

Goal: understand core concepts of PBMT,
understand word alignment (AS)

Learn to translate sentence in one lang (source) to another lang (target)

French → English

Data bitext, parallel sent pairs

Je fais un bureau I make a desk

Je fais une soupe I'm making a soup

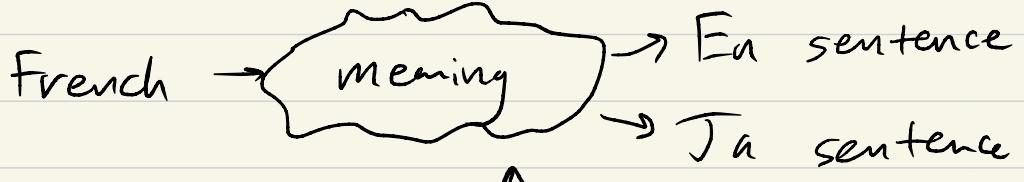
Qu'est-ce que tu fais? What are you doing?
Quelle erreur! What a mistake!

You make a mistake

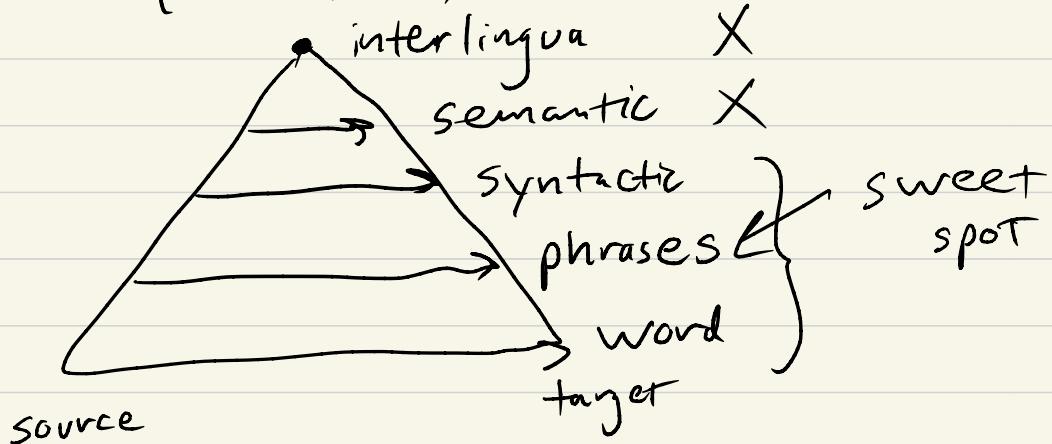
① What does "tu fais une erreur" mean?

② How do you know what *fais* means?

- Many-to-many translations
- Cooccurrence → alignment
- Phrasal translations



Vauquois (1968) levels of transfer



- ① How to get chunks/phrases?
- ② How to translate sentences? next class

Word alignment

Input: bitext

Output: Je fais un bureau

$a_1=1$ / $a_2=2$ / $a_3=2$ | $a_4=3$ | $a_5=4$

I am making a desk

1 2 2 3 4

Source = French, target = English

Alignment family of models

$P(\bar{a}, \bar{T} | \bar{S})$

\bar{T} : Eng words
 \bar{S} : Fr words
 \bar{a} : alignments

$a_i=j$: the i th En word aligns to j th Fr word

a_i takes one value
 one-to-many alignment

a_i are not necessarily monotonic

Fais - tu un bureau?

Are you making a desk?

Unsupervised learning problem

- We do not see labeled \bar{a} anywhere in our data

IBM Model 1 (1993)

- Canadian Hansards Fr-En

$$\bar{a} = (a_1, \dots, a_n) \quad \bar{t} = (t_1, \dots, t_n)$$

$\bar{s} = (s_1, \dots, s_m, \text{NULL})$ allow unaligned words

$$P(\bar{t}, \bar{a} | \bar{s}) = P(\bar{a}) P(\bar{t} | \bar{a}, \bar{s})$$

$$= \prod_{i=1}^n P(a_i) \underbrace{P(t_i | s_{a_i})}_{\text{"emission" probs}}$$

dist over $\{1, \dots, m+1\}$
uniform $\frac{1}{m+1}$ prob for
each

"emission" probs
(next page)

no parameters

$$\underline{M1}: \prod_{i=1}^n P(a_i) P(t_i | s_{a_i})$$

Model params: emissions in an HMM

$\rightarrow |V_s| \times |V_t|$ matrix of probs
 includes NULL $P(\text{target word} | \text{src word})$

$P(t_i | s_{a_i})$: look up the prob of word t_i conditioned on word s_{a_i} in the matrix

a_i is a "switch" that tells t_i what to condition on

<u>Ex</u>	<u>params</u>	Je	J1	mange	aimé	NULL	I like eat	$P(T s)$
							0.8 0.1 0.1	
							0.8 0.1 0.1	
							0 0 1.0	
							0 1.0 0	
							0.4 0.3 0.3	

Je, NULL₂

$$I \quad a_1 = \{1, 2\}$$

$$P(a_1=1, t= "I" | s= "Je")$$

$$= P(a_1=1) P(I | Je)$$

$$= 0.5 \cdot 0.8$$

Inference: $P(\bar{a} | \bar{s}, \bar{t})$ posterior over \bar{a}

From HMMs: $\bar{a} \approx \bar{y}$, $\bar{t} \approx \bar{x}$

\bar{s} just hangs out constant

doesn't depend on \bar{a}

$$\downarrow = \frac{P(\bar{a}, \bar{t} | \bar{s})}{P(\bar{t} | \bar{s})} = \frac{\prod_{i=1}^n P(a_i) P(t_i | s_{a_i})}{\underbrace{P(\bar{t} | \bar{s})}_{\text{constant}}}$$

$$P(\bar{a} | \bar{s}, \bar{t}) \propto \prod_{i=1}^n P(t_i | s_{a_i})$$

prop. to

$\bar{J}e$, $NULL_2$

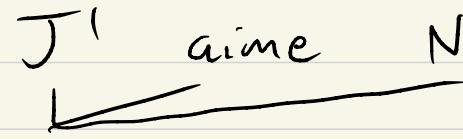
I $a_1 = \{1, 2\}$

$\begin{matrix} I \\ \parallel \\ I \end{matrix}$ $\begin{matrix} Je \\ \parallel \\ Je \end{matrix}$

$P(T, | s_1)$

$$P(a_1 | \bar{s}, \bar{T}) \propto \begin{cases} a_1 = 1 : & 0.8 \\ a_1 = 2 : & P(T, | s_2) \\ & 0.4 \end{cases} \quad NULL$$

$$P(a_1 | \bar{s}, \bar{T}) = \begin{cases} 2/3 & a_1 = 1 \\ 1/3 & a_1 = 2 \end{cases}$$

J' aime $NULL$


What is $P(a_1 | \bar{s}, \bar{T})$?

$a_1 = 1 \quad a_1 = 2 \quad a_1 = 3$

$P(I J')$	$P(I aime)$	$P(I NULL)$
0.8	0	0.4
2/3	0	1/3

Learning

No labeled \bar{a}

If we had labeled \bar{a} :

$$\max_{\bar{a}} \sum_{i=1}^D \log P(\bar{a}^{(i)}, \bar{f}^{(i)} | \bar{s}^{(i)})$$

$(\bar{s}^{(i)}, \bar{f}^{(i)}, \bar{a}^{(i)})$ labeled aligned bitext

Problem: no \bar{a} labeled

Instead: $\max_{\bar{a}} \sum_{i=1}^D \log \sum_{\bar{a}} P(\bar{a}, \bar{f}^{(i)} | \bar{s}^{(i)}) / P(\bar{f}^{(i)} | \bar{s}^{(i)})$

log marginal probability

Expectation Maximization

For Model 1: $\sum_{i=1}^D \log \prod_j P(a_j, \bar{f}_j^{(i)} | \bar{s}^{(i)})$

because of independence assumptions

$O(nm)$