

# CS 378 Lecture 18

Today

- Finish IBM Model 1
- Phrase-based MT
- Syntactic MT (briefly)
- Neural MT and seq2seq models

Recap IBM Model 1

$\bar{s} = \text{Je fais un bureau NULL}$   
 $\bar{a} = \begin{cases} a_1=1 \\ a_2=2 \\ a_3=3 \\ a_4=4 \\ a_5=5 \end{cases}$   
 $\bar{t} = \begin{cases} t_1=I \\ t_2=\text{am} \\ t_3=\text{making} \\ t_4=a \\ t_5=\text{desk} \end{cases}$

$$P(\bar{t}, \bar{a} | \bar{s}) = \prod_{i=1}^5 P(a_i) P(t_i | s_{a_i})$$

$$P(a_i) = \text{uniform dist over } \{1, 2, \dots, m+1\} \xrightarrow{\text{NULL}}$$

$$P(t_i | s_{a_i}) \quad P(I | Je) = 0.8$$

## Announcements

- Custom FP proposals returned w/comments
- Optional lecture Tuesday
- A4
- A5

## Inference in Model 1

What we care about:  $P(\bar{a} | \bar{t}, \bar{s})$

( $\underset{\bar{a}}{\operatorname{argmax}}$   $P(\bar{a} | \bar{t}, \bar{s})$ )

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

$\frac{1}{m+1}$

constant

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

$$P(a_i | \bar{t}, \bar{s}) \text{ proportional to } P(t_i | s_{a_i})$$

(model params)

$P(t|s)$

Example

I like eat

$\bar{s} = Je \text{ NULL}$

|       | I   | like | eat |
|-------|-----|------|-----|
| Je    | 0.8 | 0.1  | 0.1 |
| Tl    | 0.8 | 0.1  | 0.1 |
| mange | 0   | 0    | 1.0 |
| aine  | 0   | 1.0  | 0   |
| NULL  | 0.4 | 0.3  | 0.4 |

$a_1$

$\bar{t} = I$

$$P(a_1 | \bar{t}, \bar{s}) \stackrel{\text{prop.}}{\rightarrow} \begin{cases} a_1=1 & P(I | Je) = 0.8 \\ & Je \\ a_1=2 & P(I | \text{NULL}) = 0.4 \\ & \text{NULL} \end{cases}$$

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

Ex 2  $\text{J}^1$  are NULL

I like

$$P(a_1 | \bar{s}, \bar{t}) = \begin{cases} 0.8 \text{ J}^1 & 2/3 \\ 0 \text{ aine} \Rightarrow 0 \\ 0.4 \text{ NULL} & 1/3 \end{cases}$$

$$P(a_2 | \bar{s}, \bar{t}) = \begin{cases} 0.1 \text{ J}^1 & 1/14 \\ 1.0 \text{ aine} \Rightarrow 10/14 \\ 0.3 \text{ NULL} & 3/14 \end{cases}$$

More complex models : HMM alignment model

$$P(a_i | a_{i-1})$$

IBM Models 2-4

Learning Unsupervised learning:  $\bar{s}, \bar{t}$   
EM (Expectation Maximization)

maximize  $P(\bar{f}^{(i)} | \bar{s}^{(i)})$

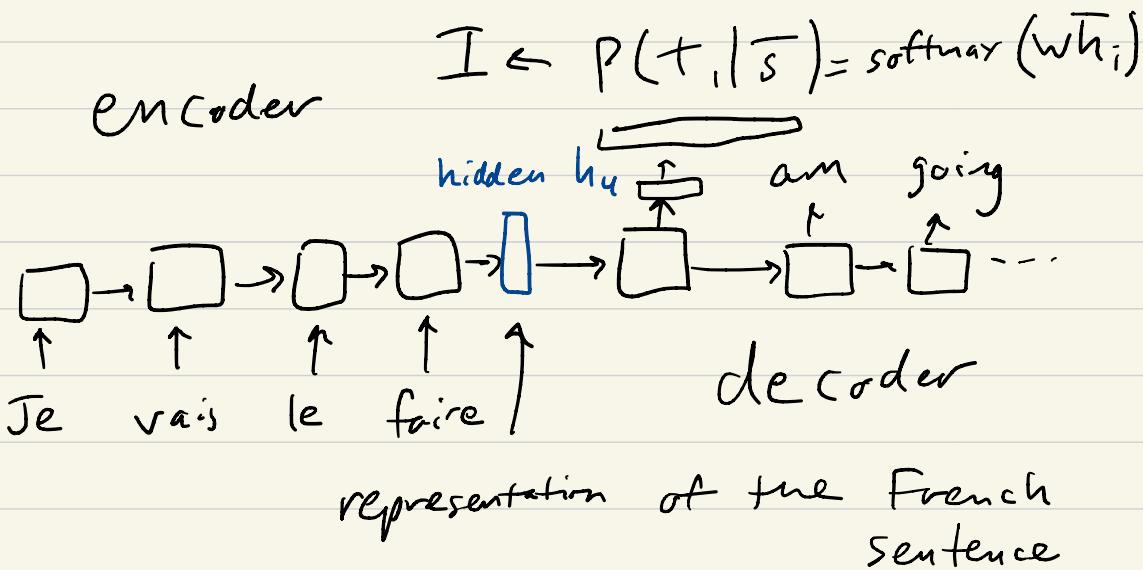
$$\sum_{i=1}^D \log \sum_{\bar{a}} P(\bar{a}, \bar{f}^{(i)} | \bar{s}^{(i)})$$

log marginal likelihood

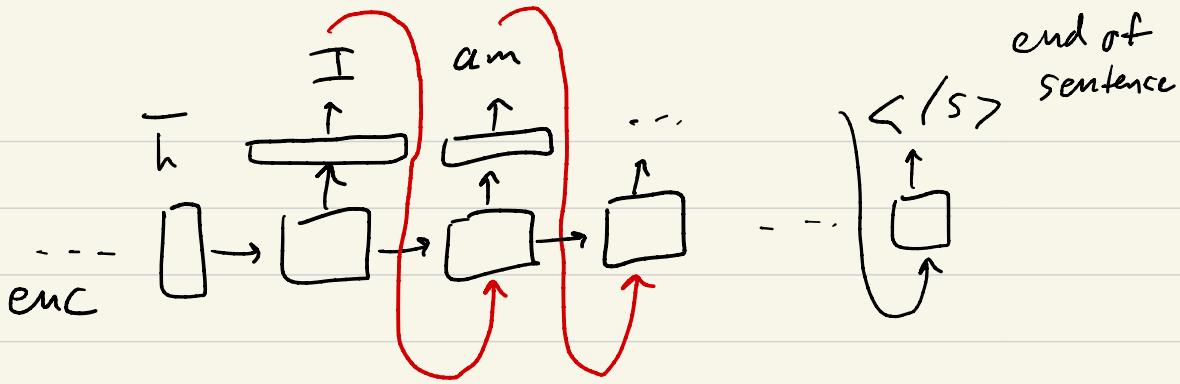
SLIDES

## Seq 2 seq models

Key idea: encode source sent w/RNN,  
"decode" target w/another RNN



$$P(\bar{F} | \bar{s}) = P(t_1 | \bar{s}) P(t_2 | \bar{s}, t_1) \dots$$



Feed output  $t_i$  into cell input for  $t_{i+1}$

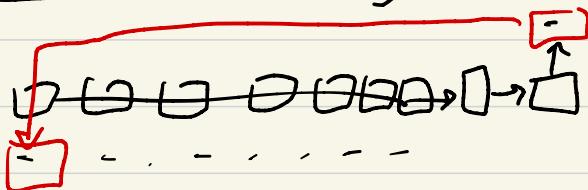
Why? Model captures word-word dependencies more easily

Training Given sequence pairs  $(\bar{s}, \bar{t})$

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

Like in language modeling, assume everything up until now matches our reference / gold

## Problem Long-range dependencies



What we want is a way to look back at the input more easily

Solution: attention