

Natural Language Processing

CSE 325/425



Sihong Xie

Lecture 24:

- Training alignment models
- Phase-based translation

learn $P(f_j | e_{a_j})$ from data
 $\{ (E_s, F_s) \}_{s=1}^{|S|}$

Training alignment models

- We have describe two alignment models

- IBM Model-1

- HMM $P(F, A|E) = P(J|I) \times \prod_{j=1}^J P(a_j|a_{j-1}, I) \boxed{P(f_j|e_{a_j})}$

translation model
parameters to be learned,

a_j comes from Alignment matrix A

and is
unknown

- How to find the probability $P(f_j|e_{a_j})$

\triangleq prob of translating e_{a_j} into f_j

where a_j = index of word in E translated

- Given an alignment A of two sentences (E, F) , this probability can be found using MLE.

from j -th
word in F

- Example: $P(verde|green)$ is estimated by the number of times *verde* is aligned with *green*, divided by number of *green*.

(E_1, F_1) \dots $(E_{|S|}, F_{|S|})$
 A_1 \dots $A_{|S|}$

EM training

$$A = \begin{matrix} & \begin{matrix} E & F & 1 & \dots & J \end{matrix} \\ \begin{matrix} E \\ F \end{matrix} & \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \end{matrix}$$

0,1 matrix $I \times I$

- However, we don't have a large corpus of aligned source and target sentences and need to estimate both $P(f_j|e_{a_j})$ and A .
 - Global model parameter: translation probability $P(f_j|e_{a_j})$
 - Local hidden (latent) variables A ($|E|, |F|$)
 - Similar to the EM training for HMM with model parameters (A, B, π) and latent tags P_b - tags of sentences.
- E-step: given translation probabilities, estimate probability of each possible alignment of each (E, F) . $P(F, A | E)$

$|E| = I$
 $|F| = J$

All possible Alignments:
 $(I+1)^J$
- M-step: given the distribution of alignments, estimate translation probabilities.

EM training

neither the Model-2, nor the HMM

- Assume a simplified model $P(F, A|E) = \prod_{j=1}^J P(f_j|e_{a_j})$
 - \uparrow
a fixed Alignment
 - critical for phase-based Translation,
- This is the common portion of both IBM Model-1 and HMM.
- Use the example training corpus of two pairs of (E, F)

E_1	green house	E_2	the house
	(A) ?		(A) ?
F_1	casa verde	F_2	la casa

- Parameters: translation probabilities initialized uniformly

Max Entropy

$P(\text{casa} \text{green})=1/3$	$P(\text{verde} \text{green})=1/3$	$P(\text{la} \text{green})=1/3$
$P(\text{casa} \text{house})=1/3$	$P(\text{verde} \text{house})=1/3$	$P(\text{la} \text{house})=1/3$
$P(\text{casa} \text{the})=1/3$	$P(\text{verde} \text{the})=1/3$	$P(\text{la} \text{the})=1/3$

} English

Spanish

EM training

- E-step 1a: ^{sub step.} for each pair and each possible alignment A , find

$$P(F, A|E) = \prod_{j=1}^J P(f_j | e_{a_j}) \quad \text{determined by } A$$

E/F

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{first iteration}$$

E e_1 e_2
 green house
 $| a_1=1$ $| a_2=2$
 F f_1 casa f_2 verde

$$P(F, A|E) = P(\text{casa}|\text{green}) \times P(\text{verde}|\text{house})$$

$$= (1/3) \times (1/3) = 1/9$$

green₁ house₂
 $a_1=2$ $a_2=1$
 casa₁ verde₂

$$P(\text{casa}|\text{house}) \times P(\text{verde}|\text{green})$$

$$= 1/9$$

the house
 $|$ $|$
 la casa

$$P(\text{la}|\text{the}) \times P(\text{casa}|\text{house})$$

$$= 1/9$$

the₁ house₂
 $|$ $|$
 la₁ casa₂

$$P(\text{la}|\text{house}) \times P(\text{casa}|\text{the})$$

$$= 1/9$$

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$P(\text{casa} \text{green})=1/3$	$P(\text{verde} \text{green})=1/3$	$P(\text{la} \text{green})=1/3$
$P(\text{casa} \text{house})=1/3$	$P(\text{verde} \text{house})=1/3$	$P(\text{la} \text{house})=1/3$
$P(\text{casa} \text{the})=1/3$	$P(\text{verde} \text{the})=1/3$	$P(\text{la} \text{the})=1/3$

EM training

- E-step 1b: for each pair and each possible alignment A , find

$$P(A|E, F) = \cancel{P(A|E)} = \frac{P(F, A|E)}{\sum_A P(F, A|E)}$$

(in HMM, $P(\theta|O)$)

$$P(A|E, F) = \frac{P(F, A|E)}{P(F|E)}$$

Prob $\frac{1}{2}$ A_1 $\begin{array}{cc} \text{green} & \text{house} \\ | & | \\ \text{casa} & \text{verde} \end{array}$ $\frac{1}{2}$ A_2 $\begin{array}{cc} \text{green} & \text{house} \\ \diagdown & \diagup \\ \text{casa} & \text{verde} \end{array}$ Prob $\frac{1}{2}$ A_1 $\begin{array}{cc} \text{the} & \text{house} \\ | & | \\ \text{la} & \text{casa} \end{array}$ $\frac{1}{2}$ A_2 $\begin{array}{cc} \text{the} & \text{house} \\ \diagdown & \diagup \\ \text{la} & \text{casa} \end{array}$

$P(F, A|E) = P(\text{casa}|\text{green}) \times P(\text{verde}|\text{house})$ $P(\text{casa}|\text{house}) \times P(\text{verde}|\text{green})$ $P(\text{la}|\text{the}) \times P(\text{casa}|\text{house})$ $P(\text{la}|\text{house}) \times P(\text{casa}|\text{the})$

$= 1/9$ $= 1/9$

- E-step 1c: find expected counts of word translation.

Expected Count (casa | green)

$= \sum \text{Prob (when "casa" is aligned with "green")}$

count(casa green)=1/2	count(verde green)=1/2	count(la green)=0	Total(green)=1
count(casa house)=1/2+1/2	count(verde house)=1/2	count(la house)=1/2	Total(house)=2
count(casa the)=1/2	count(verde the)=0	count(la the)=1/2	Total(the)=1

$$P(A_1|E, F) = \frac{(1/9)}{(2/9)} = \frac{1}{2}$$

EM training

- M-step 1: re-estimate the translation probabilities by row normalization

count(casa green)=1/2	count(verde green)=1/2	count(la green)=0	Total(green)=1
count(casa house)=1/2+1/2	count(verde house)=1/2	count(la house)=1/2	Total(house)=2
count(casa the)=1/2	count(verde the)=0	count(la the)=1/2	Total(the)=1



$P(f_i | e_{a_j})$
↘

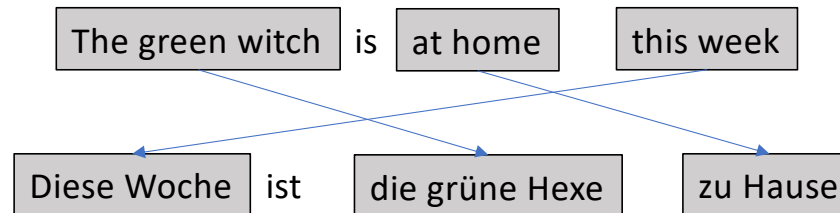
$P(\text{casa} \text{green})=1/2$	$P(\text{verde} \text{green})=1/2$	$P(\text{la} \text{green})=0$
$P(\text{casa} \text{house})=1/2$	$P(\text{verde} \text{house})=1/4$	$P(\text{la} \text{house})=1/4$
$P(\text{casa} \text{the})=1/2$	$P(\text{verde} \text{the})=0$	$P(\text{la} \text{the})=1/2$

- Observation: the unlikely translation has a smaller probability.
- Then go to E-step 2, M-step 2, E-step 3, and so on ...

How 6?

Phase-based translation

- We have been working with word-based translation models.
- Phase-based translation has some more advantages:
 - phases are constituent unit that can somehow be move freely in a sentence.



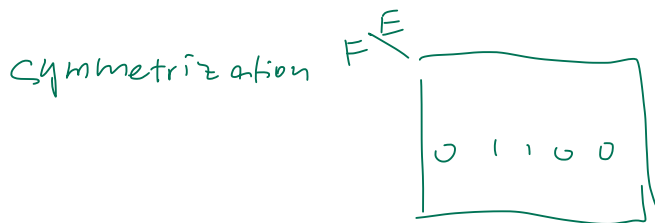
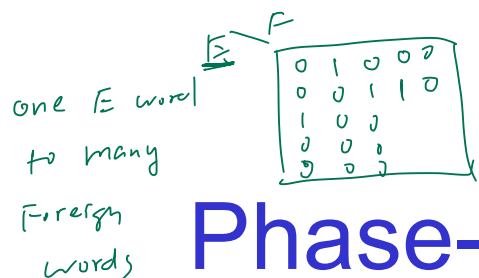
- phases have more contexts for faithful and robust translations.

Mary **didn't slap** the green witch

*Maria **no dió una bofetada** a la bruja verde*

Mary not gave a slap to the witch green

Word - word translation \Rightarrow phase - phase translation (many-to-many)
(one-to-one)



many E words
mapped to one
Foreign word

Phase-based translation

shallow parsing



(stay at the
phase level.)

- Segment F into phases $F = (\bar{f}_1, \dots, \bar{f}_I)$

- Generate $E = (\bar{e}_1, \dots, \bar{e}_J)$, with faithfulness $P(F|E) = \prod_{i=1}^I \phi(\bar{f}_i, \bar{e}_i) d(a_i - b_{i-1})$

- Phase-translation table records $\phi(\bar{f}_i | \bar{e}_i)$

- Distribution of "distortion" $d(a_i - b_{i-1})$

$d(\text{distance}) \downarrow$
if distance \uparrow

- a_i : the position of the first word of the i -th phase in E
- b_{i-1} : the position of the first word of the $(i-1)$ -th phase in E

translation of the

A translation of the

preserve "locality":

two adjacent English
phrases $(i, i-1)$ should
be translated into

Spanish phrases "close to"
each other. ?

e_0	e_1	$\hat{v}=2$ e_2	$\hat{v}=3$ e_3	$\hat{v}=4$ e_4	$\hat{v}=5$ e_5
E	Mary ₁	didn't ₂	slap ₃	the ₄	green-witch ₆
F	Maria ₁	no ₂	dió una bofetada _{3 4 5}	a la _{6 7}	bruja verde _{8 9}
b_0	a_1	a_2	a_3	b_4	a_5
"	"	"	"	"	"
0	b_1	$=b_2$	$=3$	7	8
	"	"			
	1	$=2$			

$$d(a_i - b_{i-1}) = 1 \text{ if } \hat{v}=3$$

$$Pr(F, A | E) = Pr(\text{maria} | \text{Mary}) \times d(a_1 - b_0) \times Pr(\text{no} | \text{didn't}) \times d(a_2 - b_1) \\ \times \dots \times Pr(\text{bruja verde} | \text{green, house}) \times d(a_5 - b_4)$$