

INF5820/INF9820

LANGUAGE TECHNOLOGICAL APPLICATIONS

Jan Tore Lønning, Lecture 5, 21 Sep. 2016

jtl@ifi.uio.no

Today

2

- Repetition:
 - ▣ Statistical machine translation:
 - The noisy channel model
 - IBM model 1
 - Training the intuitive way
- Training – the fast way
- Higher IBM-models

The noisy channel model

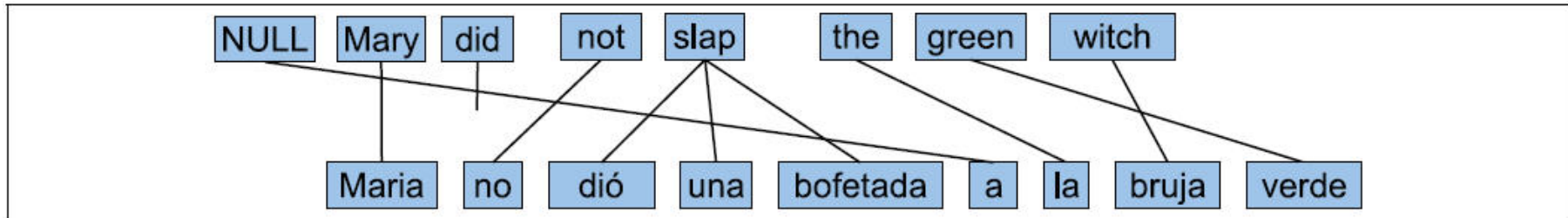
$$\begin{aligned}\hat{E} &= \arg \max_E P(E | F) \\ &= \arg \max_E \frac{P(F | E)P(E)}{P(F)} \\ &= \arg \max_E P(F | E)P(E)\end{aligned}$$



- Use n-gram language model for $P(E)$

Alignment

4



- Length of English string: k ($=7$)
- Length of foreign string: m ($=9$)
- An alignment is a vector of length m , each entry a number between 0 and k
- The example:
 - ▣ $\langle a_1, a_2, \dots, a_9 \rangle = \langle 1, 3, 4, 4, 4, 0, 5, 7, 6 \rangle$

IBM Model 1

5

- Consider all possible alignments \mathbf{a} :

$$P(\mathbf{f} | \mathbf{e}) = \sum_{\mathbf{a}} P(\mathbf{f}, \mathbf{a} | \mathbf{e})$$

- For each alignment use the simplified generative model:

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = \frac{\varepsilon}{(k+1)^m} \prod_{j=1}^m t(f_j | e_{a_j})$$

- ε is a normalisation factor
- Formula 4.7 in the SMT book
 - (The book goes $f \rightarrow e$, not $e \rightarrow f$)

Step 3.1: Collect frac. counts ctd

6

f_2, e_2 :

□ 64 alignments, with prob $1/64$

□ $f = \text{bet}, e = \text{bit}$

▣ 16 alignments connect them: $\langle x, 2, z \rangle$ for x, z in $\{0, 1, 2, 3\}$

▣ $c(\text{bet} | \text{bit}; \mathbf{f}_2, \mathbf{e}_2) = 16/64 = 1/4$

□ $f = \text{bet}, e = \text{dog}$

▣ all alignments $\langle x, 1, z \rangle$ and $\langle x, 3, z \rangle$ for x, z in $\{0, 1, 2, 3\}$

▣ $c(\text{bet} | \text{dog}; \mathbf{f}_2, \mathbf{e}_2) = 2 * 16/64 = 1/2$

e_2 : Dog bit dog

f_2 : Hund bet hund

$$c(\text{hund} | \text{dog}; \mathbf{f}_2, \mathbf{e}_2) = 1$$

$$c(\text{bet} | \text{dog}; \mathbf{f}_2, \mathbf{e}_2) = 1/2$$

$$c(\text{hund} | \text{bit}; \mathbf{f}_2, \mathbf{e}_2) = 1/2$$

$$c(\text{bet} | \text{bit}; \mathbf{f}_2, \mathbf{e}_2) = 1/4$$

$$c(\text{hund} | 0; \mathbf{f}_2, \mathbf{e}_2) = 1/2$$

$$c(\text{bet} | 0; \mathbf{f}_2, \mathbf{e}_2) = 1/4$$

Step 3.1: Collect frac. counts ctd

7

f_2, e_2 : 64 alignments, with prob $1/64$

- $f = \text{hund}, e = \text{dog}$
- 2 ways of thinking

What is the count where

- f_1 is a translation of e_1 ?
 - The alignments $\langle 1, x, y \rangle$ for any x and y
 - $16/64$
- f_1 is a translation of e_3 ?
- e_3 a translation of f_1 ?
- e_3 a translation of f_3 ?

$$4 * (16/64) = 1$$

e_2 : Dog bit dog

f_2 : Hund bet hund

For the alignment $\langle 1, 2, 3 \rangle$:

- How many times is *hund* aligned with *dog*? 2
- Similarly for $\langle x, y, z \rangle$ where x is 1 or 3 and z is 1 or 3:
 - 16 alignments
 - Frac count: $16 * 2 / 64$

For the alignments $\langle 1, y, z \rangle$ where z is 0 or 2:

- *hund-dog* aligned 1 time
- $8 * 1 / 64 = 1 / 8$
- Similarly for
 - $\langle 3, y, z \rangle$, z is 0 or 2
 - $\langle x, y, z \rangle$, for $x \in \{0, 2\}, y \in \{1, 3\}$

Fractional counts

Counting for one sentence

$$c(f | e; \mathbf{e}, \mathbf{f}) = \sum_{\mathbf{a}} (p(\mathbf{a} | \mathbf{e}, \mathbf{f}) \sum_{j=1}^m \delta(f, f_j) \times \delta(e, e_{\mathbf{a}_j}))$$

- The part $\sum_{j=1}^m \delta(f, f_j) \times \delta(e, e_{\mathbf{a}_j})$

counts how many times the alignment \mathbf{a} connects a word of the type f with one of type e

▣ $\delta(a, b) = 1$ if and only if $a = b$, otherwise 0

- We multiply with the probability of this alignment
- And sum over all alignments

Today

9

- Repetition:
 - ▣ Statistical machine translation:
 - The noisy channel model
 - IBM model 1
 - Training the intuitive way
- Training – the fast way
- Higher IBM-models

Training – the idea

10

1. From the translation probabilities, we may estimate alignment probabilities
 - ▣ (We do not choose only the best alignment)
 2. From alignment probabilities, we may maximize translation probabilities
- ▣ By alternating between (1) and (2), the numbers converge towards better results
 - ▣ For IBM Model 1 it may be proved that they converge towards a global optimum

Too many alignments

Words, $m=k$	2	3	4	5	6	7	8	9	10
Align.	9	64	625	2401	117 649	16807	43mill	282429536481	25 billions

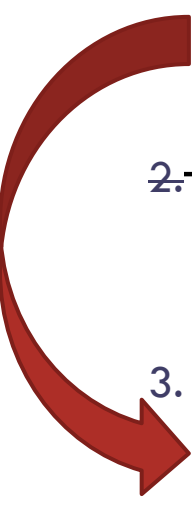
Training – the intuitive approach

12

1. Initialize the parameter values $t(f/e)$ for pairs of words f and e .
2. For each sentences pair f, e calculate the probabilities $P(a / f, e)$ of all alignments a .
3. Collect fractional counts, $tc(f/e)$:
 1. First, calculate this, $c(f/e ; f, e)$ for each sentence f, e ,
 2. Then add over all sentences
4. Calculate the new translation probabilities $t(f/e)$
5. Repeat from 2 as long as you like

Training – the efficient approach

13

1. Initialize the parameter values $t(f/e)$ for pairs of words f and e .
 - ~~2. For each sentences pair f, e calculate the probabilities $P(a / f, e)$ to all alignments a .~~
 3. Collect fractional counts, $tc(f/e)$:
 1. First, calculate this, $c(f/e ; f, e)$ for each sentence f, e ,
 2. Then add over all sentences
 4. Calculate the new translation probabilities $t(f/e)$
 5. Repeat from 2 as long as you like
- 

IBM Model 1

14

- Consider all possible alignments \mathbf{a} :

$$P(\mathbf{f} | \mathbf{e}) = \sum_{\mathbf{a}} P(\mathbf{f}, \mathbf{a} | \mathbf{e})$$

- For each alignment use the simplified generative model:

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = \frac{\varepsilon}{(k+1)^m} \prod_{j=1}^m t(f_j | e_{a_j})$$

- ε is a normalisation factor
- Formula 4.7 in the SMT book
 - (The book goes $f \rightarrow e$, not $e \rightarrow f$)

Necessary simplification

$$P(\mathbf{f} | \mathbf{e}) = \sum_{\mathbf{a}} P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = \sum_{\mathbf{a}} \frac{\varepsilon}{(k+1)^m} \prod_{j=1}^m t(f_j | e_{a_j})$$

$$P(\mathbf{f} | \mathbf{e}) = \sum_{a_1=0}^k \cdots \sum_{a_m=0}^k \frac{\varepsilon}{(k+1)^m} \prod_{j=1}^m t(f_j | e_{a_j})$$

□ This equals

$$P(\mathbf{f} | \mathbf{e}) = \frac{\varepsilon}{(k+1)^m} \prod_{j=1}^m \sum_{i=0}^k t(f_j | e_i)$$

□ Because

$$\prod_{j=1}^m \sum_{i=0}^k c_{j,i} = (c_{1,0} + c_{1,1} + \dots + c_{1,k})(c_{2,0} + \dots + c_{2,k}) \cdots (c_{m,0} + \dots + c_{m,k}) = \sum_{i=0}^k \cdots \sum_{i=0}^k \prod_{j=1}^m c_{j,i}$$

□ Reduces the problem from the order $(k+1)^n$ to roughly $k \times n$

Putting this together

□ So far

$$P(\mathbf{a} | \mathbf{f}, \mathbf{e}) = \frac{P(\mathbf{f}, \mathbf{a} | \mathbf{e})}{P(\mathbf{f} | \mathbf{e})}$$

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = \frac{\varepsilon}{(k+1)^m} \prod_{j=1}^m t(f_j | e_{a_j})$$

$$P(\mathbf{f} | \mathbf{e}) = \frac{\varepsilon}{(k+1)^m} \prod_{j=1}^m \sum_{i=0}^k t(f_j | e_i)$$

□ Hence

$$P(\mathbf{a} | \mathbf{f}, \mathbf{e}) = \frac{\frac{\varepsilon}{(k+1)^m} \prod_{j=1}^m t(f_j | e_{a_j})}{\frac{\varepsilon}{(k+1)^m} \prod_{j=1}^m \sum_{i=0}^k t(f_j | e_i)}$$

□ Formula 4.11

$$P(\mathbf{a} | \mathbf{f}, \mathbf{e}) = \frac{\prod_{j=1}^m t(f_j | e_{a_j})}{\prod_{j=1}^m \sum_{i=0}^k t(f_j | e_i)}$$

Fractional counts

Counting for one sentence

$$c(f | e; \mathbf{e}, \mathbf{f}) = \sum_{\mathbf{a}} (p(\mathbf{a} | \mathbf{e}, \mathbf{f}) \sum_{j=1}^m \delta(f, f_j) \times \delta(e, e_{\mathbf{a}_j}))$$

- The part $\sum_{j=1}^m \delta(f, f_j) \times \delta(e, e_{\mathbf{a}_j})$

counts how many times the alignment \mathbf{a} connects a word of the type f with one of type e

▣ $\delta(a, b) = 1$ if and only if $a = b$, otherwise 0

- We multiply with the probability of this alignment
- And sum over all alignments

Fractional counts

- Counting for one sentence

$$c(f | e; \mathbf{e}, \mathbf{f}) = \sum_{\mathbf{a}} (p(\mathbf{a} | \mathbf{e}, \mathbf{f}) \sum_{j=1}^m \delta(f, f_j) \times \delta(e, e_{\mathbf{a}_j}))$$

- Substituting in for $p(\mathbf{a} | \mathbf{e}, \mathbf{f})$

$$c(f | e; \mathbf{e}, \mathbf{f}) = \sum_{\mathbf{a}} \left(\frac{\prod_{j=1}^m t(f_j | e_{\mathbf{a}_j})}{\prod_{j=1}^m \sum_{i=0}^k t(f_j | e_i)} \sum_{j=1}^m \delta(f, f_j) \times \delta(e, e_{\mathbf{a}_j}) \right)$$

- and doing some non-trivial calculation:

$$c(f | e; \mathbf{e}, \mathbf{f}) = \frac{t(f | e)}{\sum_{i=0}^k t(f | e_i)} \sum_{j=1}^m \delta(f, f_j) \sum_{i=0}^k \delta(e, e_i)$$

Observe:
Directly from t to
 $c(f|e; \mathbf{e}, \mathbf{f})$ without
mentioning the \mathbf{a} -s

Fractional counts

- Counting over the whole corpus and normalize as before

$$t(f | e) = \frac{\sum_{(\mathbf{f}, \mathbf{e})} c(f | e; \mathbf{f}, \mathbf{e})}{\sum_{f'} \sum_{(\mathbf{f}, \mathbf{e})} c(f' | e; \mathbf{f}, \mathbf{e})}$$

Example – the efficient way

20

□ Corpus

e_1 : Dog barked
 f_1 : Hund bjeffet

e_2 : Dog bit dog
 f_2 : Hund bet hund

3 English words: dog bit barked
3 foreign words: hund bjeffet bet

Uniform initialization

$t(\text{hund} \text{dog}) = 1/3$	$t(\text{bet} \text{dog}) = 1/3$	$t(\text{bjeffet} \text{dog}) = 1/3$
$t(\text{hund} \text{bit}) = 1/3$	$t(\text{bet} \text{bit}) = 1/3$	$t(\text{bjeffet} \text{bit}) = 1/3$
$t(\text{hund} \text{barked}) = 1/3$	$t(\text{bet} \text{barked}) = 1/3$	$t(\text{bjeffet} \text{barked}) = 1/3$
$t(\text{hund} 0) = 1/3$	$t(\text{bet} 0) = 1/3$	$t(\text{bjeffet} 0) = 1/3$

$$c(f | e; \mathbf{e}, \mathbf{f}) = \frac{t(f | e)}{\sum_{i=0}^k t(f | e_i)} \sum_{j=1}^m \delta(f, f_j) \sum_{i=0}^k \delta(e, e_i)$$

e_1 : Dog barked
 f_1 : Hund bjeffet

$$c(\text{hund} | \text{barked}; \mathbf{e}_1, \mathbf{f}_1) = \frac{t(\text{hund} | \text{barked})}{\sum_{i=0}^2 t(f | e_i)} \sum_{j=1}^2 \delta(\text{hund}, f_j) \sum_{i=0}^2 \delta(\text{barked}, e_i) =$$

$$\frac{1/3}{\sum_{i=0}^2 (1/3)} (\delta(\text{hund}, \text{hund}) + \delta(\text{hund}, \text{bjeffet})) \times$$

$$(\delta(\text{barked}, 0) + \delta(\text{barked}, \text{dog}) + \delta(\text{barked}, \text{barked})) = 1/3$$

Uniform initialization

$t(\text{hund} \text{dog}) = 1/3$	$t(\text{bet} \text{dog}) = 1/3$	$t(\text{bjeffet} \text{dog}) = 1/3$
$t(\text{hund} \text{bit}) = 1/3$	$t(\text{bet} \text{bit}) = 1/3$	$t(\text{bjeffet} \text{bit}) = 1/3$
$t(\text{hund} \text{barked}) = 1/3$	$t(\text{bet} \text{barked}) = 1/3$	$t(\text{bjeffet} \text{barked}) = 1/3$
$t(\text{hund} 0) = 1/3$	$t(\text{bet} 0) = 1/3$	$t(\text{bjeffet} 0) = 1/3$

$$c(f | e; \mathbf{e}, \mathbf{f}) = \frac{t(f | e)}{\sum_{i=0}^k t(f | e_i)} \sum_{j=1}^m \delta(f, f_j) \sum_{i=0}^k \delta(e, e_i)$$

e_2 : Dog bit dog
 f_2 : Hund bet hund

$$c(\text{bet} | \text{bit}; \mathbf{e}_2, \mathbf{f}_2) = \frac{t(\text{bet} | \text{bit})}{\sum_{i=0}^3 t(\text{bet} | e_i)} \sum_{j=1}^3 \delta(\text{bet}, f_j) \sum_{i=0}^3 \delta(\text{bit}, e_i) = \frac{1/3}{\sum_{i=0}^3 (1/3)} \times 1 \times 1 = 1/4$$

$$c(\text{hund} | \text{dog}; \mathbf{e}_2, \mathbf{f}_2) = \frac{t(\text{hund} | \text{dog})}{\sum_{i=0}^3 t(\text{hund} | e_i)} \sum_{j=1}^3 \delta(\text{hund}, f_j) \sum_{i=0}^3 \delta(\text{dog}, e_i) = \frac{1/3}{\sum_{i=0}^3 (1/3)} \times 2 \times 2 = 1$$

Collect fractional counts

23

e_1 : Dog barked
 f_1 : Hund bjeffet

Results are the same as
the intuitive way

$$c(\text{hund} \mid \text{dog}; \mathbf{f}_1, \mathbf{e}_1) = 1/3$$

$$c(\text{bjeffet} \mid \text{dog}; \mathbf{f}_1, \mathbf{e}_1) = 1/3$$

$$c(\text{hund} \mid \text{barked}; \mathbf{f}_1, \mathbf{e}_1) = 1/3$$

$$c(\text{bjeffet} \mid \text{barked}; \mathbf{f}_1, \mathbf{e}_1) = 1/3$$

$$c(\text{hund} \mid 0; \mathbf{f}_1, \mathbf{e}_1) = 1/3$$

$$c(\text{bjeffet} \mid 0; \mathbf{f}_1, \mathbf{e}_1) = 1/3$$

e_2 : Dog bit dog
 f_2 : Hund bet hund

$$c(\text{hund} \mid \text{dog}; \mathbf{f}_2, \mathbf{e}_2) = 1$$

$$c(\text{bet} \mid \text{dog}; \mathbf{f}_2, \mathbf{e}_2) = 1/2$$

$$c(\text{hund} \mid \text{bit}; \mathbf{f}_2, \mathbf{e}_2) = 1/2$$

$$c(\text{bet} \mid \text{bit}; \mathbf{f}_2, \mathbf{e}_2) = 1/4$$

$$c(\text{hund} \mid 0; \mathbf{f}_2, \mathbf{e}_2) = 1/2$$

$$c(\text{bet} \mid 0; \mathbf{f}_2, \mathbf{e}_2) = 1/4$$

Step 3.2: Total counts (as before)

24

$$tc(f | e) = \sum_{(\mathbf{f}, \mathbf{e})} c(f | e; \mathbf{f}, \mathbf{e})$$

$tc(\text{hund} \text{dog}) = 1 + 1/3$	$tc(\text{bet} \text{dog}) = 1/2$	$tc(\text{bjeffet} \text{dog}) = 1/3$	$tc(* \text{dog}) = 4/3 + 1/2 + 1/3 = 13/6$
$tc(\text{hund} \text{bit}) = 1/2$	$tc(\text{bet} \text{bit}) = 1/4$	$tc(\text{bjeffet} \text{bit}) = 0$	$tc(* \text{bit}) = 3/4$
$tc(\text{hund} \text{barked}) = 1/3$	$tc(\text{bet} \text{barked}) = 0$	$tc(\text{bjeffet} \text{barked}) = 1/3$	$tc(* \text{barked}) = 2/3$
$tc(\text{hund} 0) = 1/2 + 1/3$	$tc(\text{bet} 0) = 1/4$	$tc(\text{bjeffet} 0) = 1/3$	$tc(* 0) = 17/12$

Step 4: new trans. probabilities

25

$$t(f|e) = \frac{tc(f|e)}{\sum_{f'} tc(f'|e)}$$

e	f	t(f e)	exact	decimal
0	hund	(5/6)/(17/12)	10/17	0.588235
0	bet	(1/4)/(17/12)	3/17	0.176471
0	bjeffet	(1/3)/(17/12)	4/17	0.235294
dog	hund	(4/3)/(13/6)	8/13	0.615385
dog	bet	(1/2)/(13/6)	3/13	0.230769
dog	bjeffet	(1/3)/(13/6)	2/13	0.153846
bit	hund	(1/2)/(3/4)	2/3	0.666667
bit	bet	(1/4)/(3/4)	1/3	0.333333
barked	hund	(1/3)/(2/3)	1/2	0.5
barked	bjeffet	(1/3)/(2/3)	1/2	0.5

Repeat: calculate fractional counts

26

□ Examples

$$c(\text{hund} \mid \text{barked}; \mathbf{e}_1, \mathbf{f}_1) = \frac{t(\text{hund} \mid \text{barked})}{\sum_{i=0}^2 t(f \mid e_i)} \sum_{j=1}^2 \delta(\text{hund}, f_j) \sum_{i=0}^2 \delta(\text{barked}, e_i) =$$
$$\frac{0.5}{0.588235 + 0.615385 + 0.5} = \frac{0.5}{1.70362} = 0.2934927$$

$$c(\text{hund} \mid \text{dog}; \mathbf{e}_2, \mathbf{f}_2) = \frac{t(\text{hund} \mid \text{dog})}{\sum_{i=0}^3 t(\text{hund} \mid e_i)} \sum_{j=1}^3 \delta(\text{hund}, f_j) \sum_{i=0}^3 \delta(\text{dog}, e_i) =$$
$$\frac{0.615385}{0.588235 + 0.615385 + 0.666667 + 0.615385} \times 2 \times 2 = ?$$

After some iterations

		1st iterat.	2nd iter.	5th iter.	25th iter	100th
0	hund	0.588235				
0	bet	0.176471				
0	bjeffet	0.235294				
dog	hund	0.615385				
dog	bet	0.230769				
dog	bjeffet	0.153846				
bit	hund	0.666667				
bit	bet	0.333333				
barked	hund	0.5				
barked	bjeffet	0.5				

After some iterations

		1st iterat.	2nd iter.	5th iter.	25th iter	100th
0	hund	0.588235	0.647158			
0	bet	0.176471	0.14363			
0	bjeffet	0.235294	0.209212			
dog	hund	0.615385	0.675859			
dog	bet	0.230769	0.237614			
dog	bjeffet	0.153846	0.086527			
bit	hund	0.666667	0.609848			
bit	bet	0.333333	0.390152			
barked	hund	0.5	0.342932			
barked	bjeffet	0.5	0.657068			

After some iterations

		1st iterat.	2nd iter.	5th iter.	25th iter	100th
0	hund	0.588235	0.647158	0.81929		
0	bet	0.176471	0.14363	0.067291		
0	bjeffet	0.235294	0.209212	0.113419		
dog	hund	0.615385	0.675859	0.773893		
dog	bet	0.230769	0.237614	0.214793		
dog	bjeffet	0.153846	0.086527	0.011313		
bit	hund	0.666667	0.609848	0.417491		
bit	bet	0.333333	0.390152	0.582509		
barked	hund	0.5	0.342932	0.097766		
barked	bjeffet	0.5	0.657068	0.902234		

After some iterations

		1st iterat.	2nd iter.	5th iter.	25th iter	100th
0	hund	0.588235	0.647158	0.81929	0.998457	
0	bet	0.176471	0.14363	0.067291	0.000122	
0	bjeffet	0.235294	0.209212	0.113419	0.001421	
dog	hund	0.615385	0.675859	0.773893	0.947458	
dog	bet	0.230769	0.237614	0.214793	0.052541	
dog	bjeffet	0.153846	0.086527	0.011313	0	
bit	hund	0.666667	0.609848	0.417491	0.005351	
bit	bet	0.333333	0.390152	0.582509	0.994648	
barked	hund	0.5	0.342932	0.097766	6.0e-07	
barked	bjeffet	0.5	0.657068	0.902234	0.999999	

After some iterations

		1st iterat.	2nd iter.	5th iter.	25th iter	100th
0	hund	0.588235	0.647158	0.81929	0.998457	1
0	bet	0.176471	0.14363	0.067291	0.000122	0
0	bjeffet	0.235294	0.209212	0.113419	0.001421	0
dog	hund	0.615385	0.675859	0.773893	0.947458	0.966031
dog	bet	0.230769	0.237614	0.214793	0.052541	0.033968
dog	bjeffet	0.153846	0.086527	0.011313	0	0
bit	hund	0.666667	0.609848	0.417491	0.005351	0
bit	bet	0.333333	0.390152	0.582509	0.994648	1
barked	hund	0.5	0.342932	0.097766	6.0e-07	0
barked	bjeffet	0.5	0.657068	0.902234	0.999999	1

Results (perplexity)

32

- Claim: «the numbers converge towards better results»
- Means: for each round

$$\prod_{(e,f)} P(\mathbf{f} | \mathbf{e})$$

does not decrease

- For IBM Model 1 it may be proved that they converge towards a global optimum

Today

33

- Repetition:
 - ▣ Statistical machine translation:
 - The noisy channel model
 - IBM model 1
 - Training the intuitive way
- Training – the fast way
- Higher IBM-models

IBM model 2

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = P(m | \mathbf{e}) \prod_{j=1}^m P(a_j | a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) P(f_j | a_1^j, f_1^{j-1}, m, \mathbf{e})$$

□ New

- $P(a_j | a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) = a(a_j | j, m, k)$

- For a probability distribution a

- i.e. it depends on the length of the string and the position

- (less likely to move far than to stay close)

□ As for Model 1

- $P(m | \mathbf{e})$ is a constant, independent of m and E

- the word translation probability only depends on source word

$$P(f_j | a_1^j, f_1^{j-1}, m, \mathbf{e}) = t(f_j | e_{a_j})$$

Model2

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = \varepsilon \prod_{j=1}^m a(a_j | j, m, k) t(f_j | e_{a_j})$$

- We can do similar steps as for Model1 for expressing $P(\mathbf{f} | \mathbf{e})$ and $P(\mathbf{a})$.
- We can do similar simplifications to bypass the exponential number of alignments, and
- Learn the alignment probabilities $a(a_j | j, m, k)$ at the same time as the translation probabilities
- You don't have to learn the details

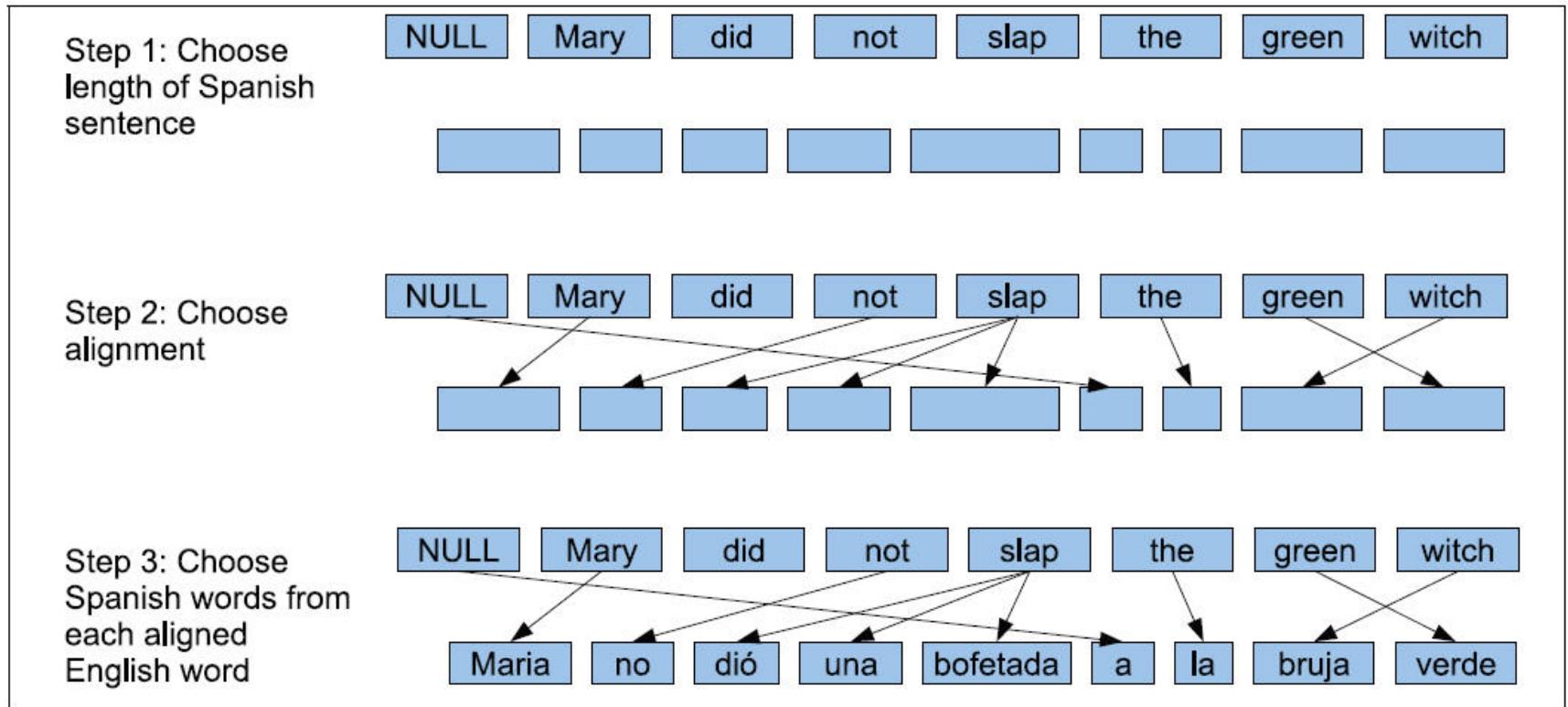
HMM Alignment

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = P(m | \mathbf{e}) \prod_{j=1}^m P(a_j | a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) P(f_j | a_1^j, f_1^{j-1}, m, \mathbf{e})$$

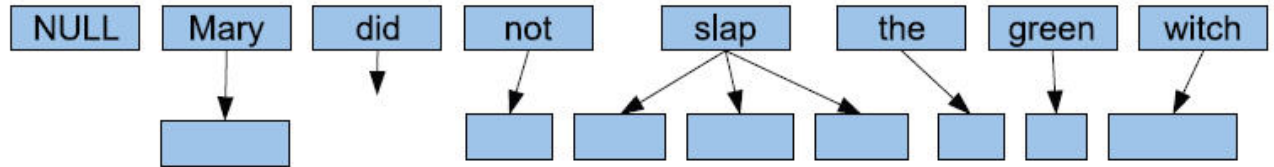
$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = P(m | k) \prod_{j=1}^m P(a_j | a_1^{j-1}, k) t(f_j | e_{a_j})$$

- $P(m | k)$ depends on the length k of \mathbf{e} .
- $P(a_j | a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) = P(a_j | a_{j-1}, k) = \lambda c(a_j - a_{j-1})$
 - Where word j should come from, depends on where word $j-1$ came from
 - This is again reduced to probabilities, c , of the distance between a_j and a_{j-1} independently of the actual j .

Model 1 & 2 and HMM alignment

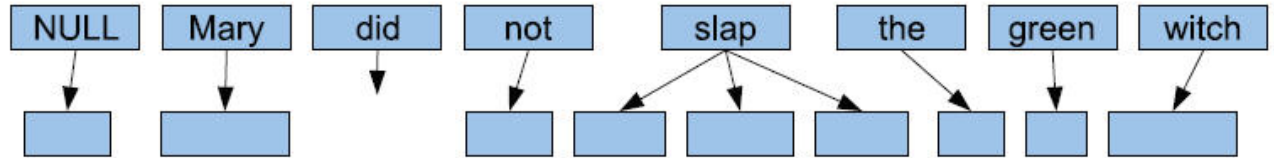


Step 1: Choose fertility for each English word

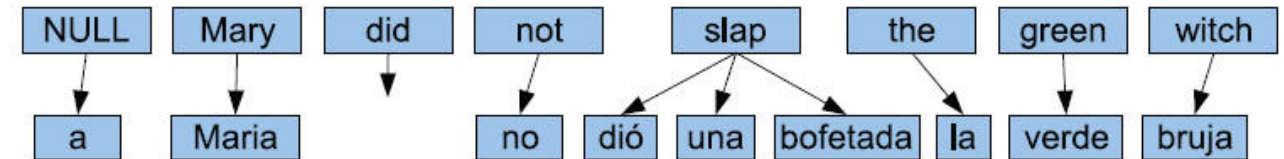


Model 3

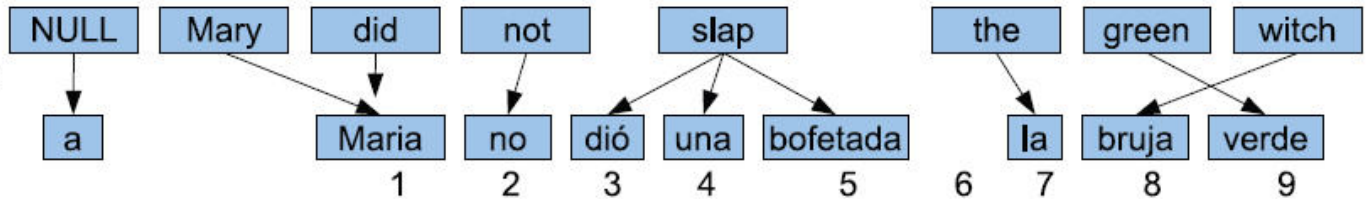
Step 2: Choose fertility for NULL



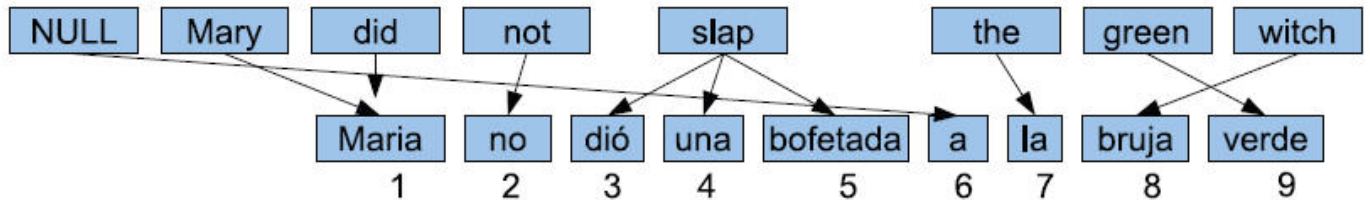
Step 3: Create Spanish words by translating aligned English word



Step 4: Move the Spanish words into final slots



Step 5: Move spurious Spanish words into unclaimed slots



IBM Model 3: Fertility

- Fertility: number of F words produced by an E word
- Modelled by a distribution $n(x|e)$

Example:

F = Norw.

$n(2 | \text{yesterday}) \approx 1$

$n(1 | \text{to}) \approx 0.8$

$n(2 | \text{to}) \approx 0.2$

$n(1 | \text{car}) \approx 1$

$n(0 | \text{the}) \approx 0.6$

$n(1 | \text{the}) \approx 0.4$

Example:

Norw. \rightarrow Eng.

$n(2 | \text{bilen}) \approx 0.7$

$n(1 | \text{bilen}) \approx 0.3$

$n(1 | \text{å}) \approx 0.8$

$n(0 | \text{å}) \approx 0.2$

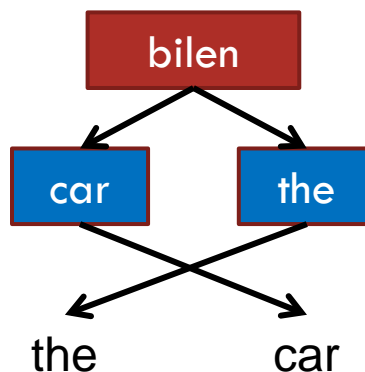
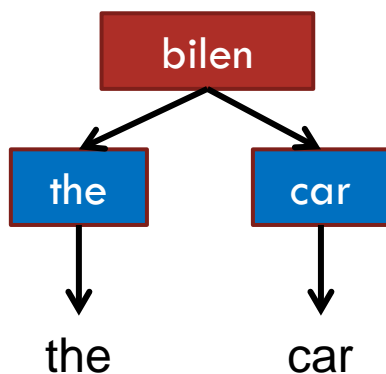
IBM Model 3: Null insertion

- Modelled by:
- There is a probability p_0 :
 - ▣ After each inserted word there is the probability p_0 of not inserting a null-word
 - ▣ And a probability $p_1 = (1-p_0)$ of inserting a null-word
- A rather complex expression for what this contributes into $P(\mathbf{a}, \mathbf{f} | \mathbf{e})$ which considers
 - ▣ Permutations
 - ▣ Length of \mathbf{f}

IBM Model 3: Distortion

$$d(j | a_j, m, k)$$

- A probability distribution which gives the probability of word a_i ending up in position j .
- Similar to alignment in model 2 but:
 - ▣ Oposite direction
 - ▣ Different choices of words + distortion may correpond to the same alignment



IBM model 3

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = \prod_{j=1}^m t(f_j | e_{a_j}) \prod_{j=1}^m d(j | a_j, k, m) \times \text{more}$$

- Where *more* is an expression which counts
 - ▣ $n(x | e_i)$ the right number of times
 - ▣ And uses p_0 to give the right probability to null-insertion.

Training Model 3

- In principle like Model 1, but
 - ▣ The trick to get rid of the alignments does not work
 - ▣ Too costly to calculate all alignments
- Strategy
 - ▣ Sample and use the most probable alignments
 - ▣ Start with alignments for Model 1 and Model 2
 - ▣ Use hill-climbing algorithm

Hill-climbing algorithm

- Assign some initial parameter values
- Consider several alternative sets of parameter values in the vicinity of where you are
- Compare the resulting values and choose the parameters which yield the best results
- Repeat

Training model 3

- Model 1: The optimum we find is global
- Model 3 (and model 2):
 - ▣ A local optimum does not have to be global
- First run some iterations of Model 1 and maybe some iterations of Model 2
- Use the results, in particular the alignment, as input to Model 3
- Hill-climb the space of alignments from here, doing minimal changes.

IBM Model 4

- Better reordering model
- Consider group of words (phrases)
- Distinguish between
 - ▣ the placement of the whole group
 - ▣ The placement within the group

The IBM-models

- IBM models 1-4 are not true probability models.
- Model 5 fixes this
 - ▣ Based on model 4
- We will not consider models 4 and 5
- Phrase Based translation makes use of Model 3