# Natural Language Processing
## CSE 325/425

Sihong Xie

Lecture 23:
- Statistical machine translation
- Alignment models (IBM model-1 and HMM)

# Statistical MT framework

- Goals of MT

  - Fluency and faithfulness

best-translation $\hat{T} = \text{argmax}_T\ P(T)\ P(S|T)$

English *(handwritten)* French *(handwritten)*

LM n-gram *(handwritten)*   faithfulness *(handwritten)*
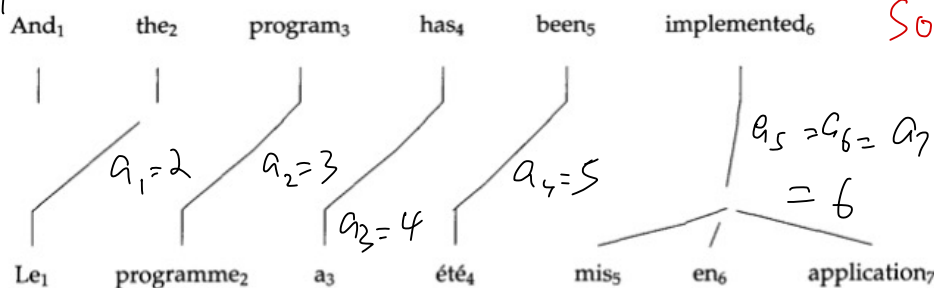
  - More formally,

    - $E=(e_1, \ldots, e_I)$: English (the target language).

    - $F=(f_1, \ldots, f_J)$: a foreign language (Spanish/French/…).

    - Language model  $P(E)$

    - Translation model:  $P(F|E)$

    - Using the Bayes theorem, *decoding* is defined as

$$E^* \quad F^* = \arg\max_{F\ E} P(E|F) = \arg\max_{F\ E} \frac{P(F|E)P(E)}{P(F)} = \arg\max_{F\ E} P(F|E)P(E)$$

# Alignment

- Word alignment: mapping words in *E* to words in *F*.

  $E^* = \arg\max_{E} P(E) \times P(F|E)$

  o Multiple target words can be mapped to one source word.

  o Example:

  Source *E*

  spurious

  | And$_1$ | the$_2$ | program$_3$ | has$_4$ | been$_5$ | implemented$_6$ |

  $a_1=2$   $a_2=3$   $a_3=4$   $a_4=5$   $a_5=a_6=a_7=6$

  Translate

  | Le$_1$ | programme$_2$ | a$_3$ | été$_4$ | mis$_5$ | en$_6$ | application$_7$ |

  Target *F*

  In ML:
  this is a so-called
  "instructured prediction"

  Applications
  CV, NLP

- Two representation of an alignment *A*

  o Recording the mapping *A* directly: $a_1=2$, $a_2=3$, ..., $a_7=6$.

  o Use an alignment matrix

  

  | *E* \ *F* | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
  |---|---|---|---|---|---|---|---|
  | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
  | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
  | 3 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
  | 4 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
  | 5 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
  | 6 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |

# Alignment

*target* *Source* *Alignment matrix* *§ 25*

- One source word mapped to multiple target words.

  *E (Source)* *E (Target)*

- Many-to-many mapping *(phase-based Alignment)*

*E*

*F*

*The_1   poor_2   don't_3   have_4   any_5   money_6*

*Les_1   pauvres_2   sont_3   demunis_4*

$The_1$
$balance_2$
$was_3$
$the_4$
$territory_5$
$of_6$
$the_7$
$aboriginal_8$
$people_9$

$Le_1$
$reste_2$
$appartenait_3$
$aux_4$
$autochtones_5$
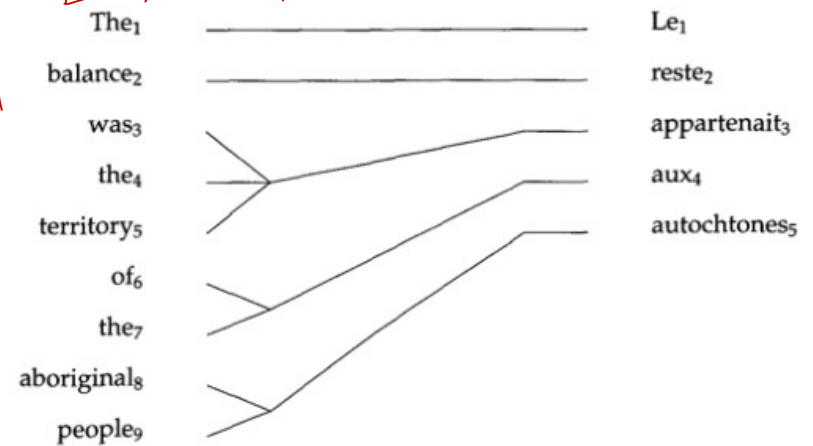
- Allow spurious word NULL for words that can't be mapped.

# Alignment (IBM model 1)

*Attention model in Neural Network ML*

- Published in *Brown, etc.1993*.

- Translation model $P(F|E) = \sum_A P(F,A|E)$

$A$: random matrix

$A \in \mathcal{A}$: Sample Space of all possible alignments

$Pr(A)$

- Generative story

  1. generate length of the source/foreign language *J;*

  2. generate an alignment *A;*

  3. generate words *E=(e₁, …, eₗ )* in the target language.



Step 1: Choose length of Spanish sentence

| NULL | Mary | did | not | slap | the | green | witch |

*F*

*E*

*J=5*

Step 2: Choose alignment

| NULL | Mary | did | not | slap | the | green | witch |

Step 3: Choose Spanish words from each aligned English word

| NULL | Mary | did | not | slap | the | green | witch |
| Maria | no | dió | una | bofetada | a | la | bruja | verde |

$P(F, A|E)?$

The Mathematics of Statistical Machine Translation. Brown, etc. 1993 Computational Linguistics

# Alignment (IBM model 1)

$j=2$

$J=9$

$S+1 = I+1$ rows

NULL

$A$

- Probability of generating a length $J$: a small positive number $\epsilon$

- Probability of an alignment between $I$ and $J$ words: $P(A(I,J)) = \dfrac{1}{(I+1)^J}$
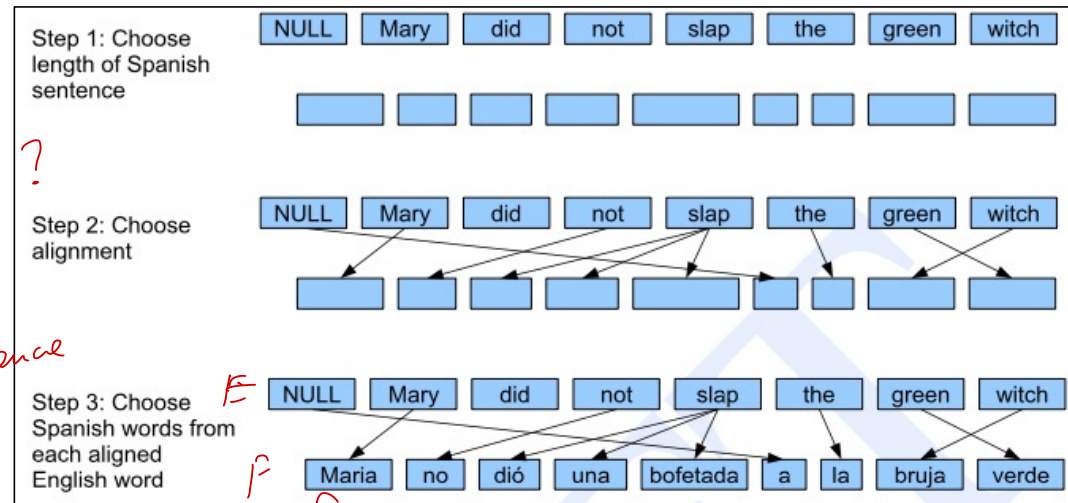
- Probability of the foreign word $f_j$

$j \in \{1, \cdots, J\}$ given the aligned target word $e_{a_j}$

$a_j \in \{0, \cdots, I\}$

$$P(f_j|e_{a_j})$$ $P(no|not)$ ?

- Probability of the foreign sentence

and an alignment $A$:

conditional

Independence

$$P(F,A|E) = \frac{\epsilon}{(I+1)^J} \prod_{j=1}^{J} P(f_j|e_{a_j})$$

$A$



Step 1: Choose length of Spanish sentence

| NULL | Mary | did | not | slap | the | green | witch |

Step 2: Choose alignment

Step 3: Choose Spanish words from each aligned English word

$E$

$F$

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

$f_1 \quad f_2 \quad f_3 \quad - \quad - \quad - \quad - \quad f_9$

$j=1 \cdots 9$

- Probability of $F$ given $E$:

$$P(F|E) = \sum_A P(F,A|E) = \sum_A \frac{\epsilon}{(I+1)^J} \prod_{j=1}^{J} P(f_j|e_{a_j})$$
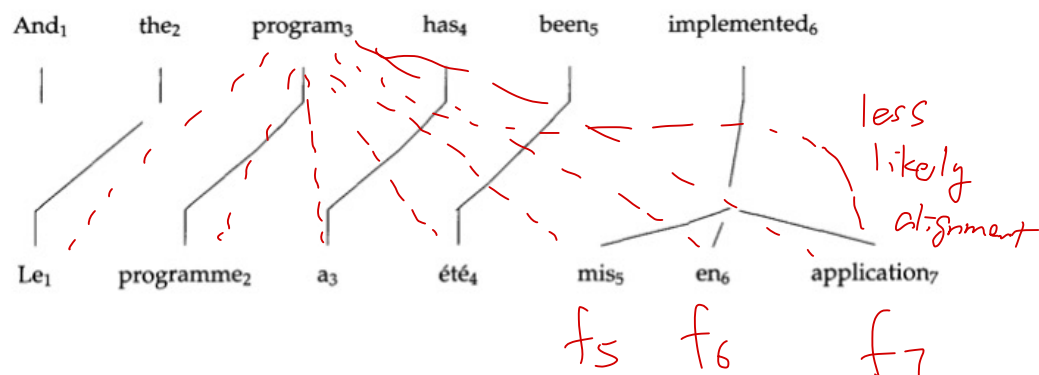
$j=1 \cdots 9$

# Alignment (HMM)

- IBM Model 1 makes several strong assumptions
  - two foreign words are mapped independently;
  - all alignments have the same probability.
- Both are not true
  - the last three French words are mapped jointly to the last English word.
    - need to use some joint probability to model such linguistic phenomena;
  - locality: two consecutive French words are likely mapped to consecutive English words.
    - some alignments are unlikely.
    - can you give one?



$$Pr(f_5, f_6, f_7 \mid A, E)$$
$$\neq Pr(f_5 \mid A, E) \cdots Pr(f_7 \mid AE)$$
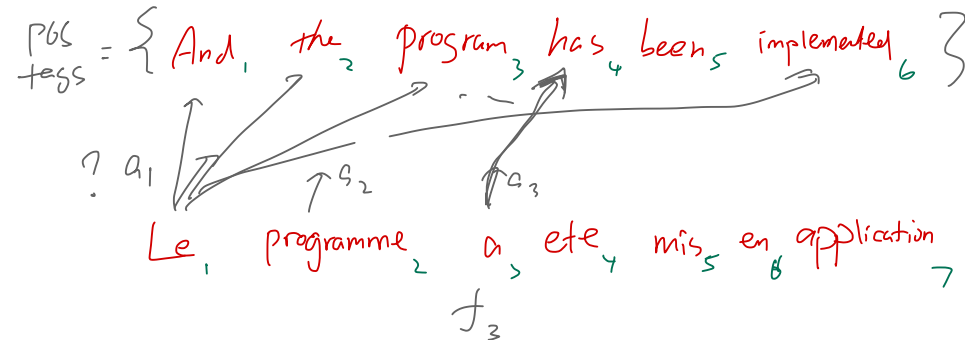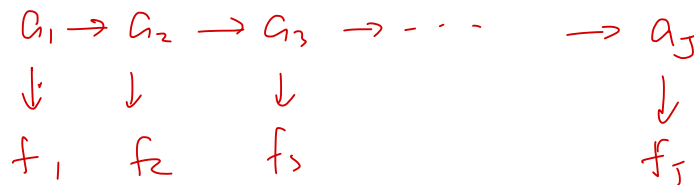
# Alignment (HMM)

- HMM alignment model

  - generates an alignment and the observed foreign sentence

    $a_j \in \{0, \cdots I\}$

    - alignment ⇔ POS-tags.  $a_j$  $a_1 \cdots a_{j-1}$

    - transition probabilities: align the next foreign word, given previous alignments.

    - foreign sentence ⇔ observed words in a sentence.

    - emission probabilities: emit a word given the up-to-date alignments.  $a_1 \cdots a_j$

$f_j$

$$P(f_1^J, a_1^J | e_1^I) = P(J|e_1^I) \times \prod_{j=1}^{J} P(f_j, a_j | f_1^{j-1}, a_1^{j-1}, e_1^I)$$

$$= P(J|e_1^I) \times \prod_{j-1}^{J} P(a_j | f_1^{j-1}, a_1^{j-1}, e_1^I) \times P(f_j | f_1^{j-1}, a_1^j, e_1^I)$$

$a_3$  $a_1^2$

POS tags $= \{$ And$_1$ the$_2$ program$_3$ has$_4$ been$_5$ implemented$_6 \}$

HMM  $a_1 \rightarrow a_2 \rightarrow a_3 \rightarrow \cdots \rightarrow a_J$  ? $a_1$  $a_2$  $a_3$

$\downarrow \quad \downarrow \quad \downarrow \qquad \qquad \downarrow$

$f_1 \quad f_2 \quad f_3 \qquad \qquad f_J$  Le$_1$ programme$_2$ a$_3$ ete$_4$ miss$_5$ en$_6$ application$_7$

$f_3$

# Alignment (HMM)

- Make some Markov assumptions to simplify the above joint probability.

disgard $(a_1 \cdots a_{j-2})$

$$P(a_j | f_1^{j-1}, a_1^{j-1}, e_1^I) = P(a_j | a_{j-1}, I)$$
$$P(f_j | f_1^{j-1}, a_1^j, e_1^I) = P(f_j | e_{a_j})$$

$Pr(a | has)$

$\uparrow$

$f_3 \qquad e_{a_3} = e_4 \ (b/c \ a_3 = 4)$

- The final joint distribution

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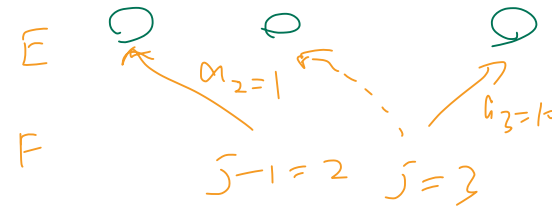
- The translation model

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

$$= \sum_A Pr(F, A | E)$$

$\nwarrow$ Sum - product

$\Updownarrow$
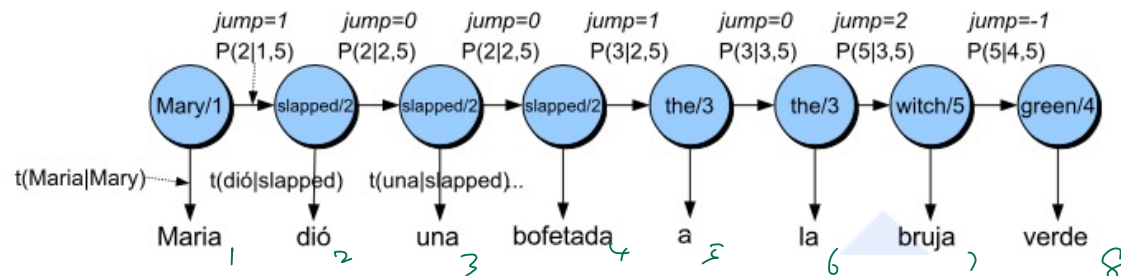
Product - sum

(forward Alg.)

# Alignment (HMM)

- The transition probability should encourage alignment locality.
  - *"the English words that generate neighboring Spanish words are likely to be nearby"* -- SLP
  - $P(a_j|a_{j-1}, I)$ should be large if $a_j$ is close to $a_{j-1}$
  - Locality is a relative concept and absolute positions are not relevant.
    - P(7|6, 15) = P(9|8, 15).   *place – invariance*
  - Model "jumps" of the alignment pointers
    - $P(a_j|a_{j-1}, I)$ is a decreasing function of the jump $|a_j - a_{j-1}|$

$E$ $\alpha_{2}=1$ $G_{3}=10$
$F$ $j-1=2$ $j=3$



| jump=1 | jump=0 | jump=0 | jump=1 | jump=0 | jump=2 | jump=-1 |
| P(2\|1,5) | P(2\|2,5) | P(2\|2,5) | P(3\|2,5) | P(3\|3,5) | P(5\|3,5) | P(5\|4,5) |

Mary/1 → slapped/2 → slapped/2 → slapped/2 → the/3 → the/3 → witch/5 → green/4

t(Maria|Mary) → t(dió|slapped) t(una|slapped)..

Maria₁   dió₂   una₃   bofetada₄   a₅   la₆   bruja₇   verde₈

$I = 8$
$J = 8$

Mary₁ slapped₂ the₃
Green₄ witch₅

$Pr(F, A | E) = P(J|I) \times Pr(2|1,5) \times Pr(Maria|Mary)$
$\qquad t_2 \quad e_{a_2} \quad f_1 \quad e_{a_1}$
$\times Pr(2|2,5) \times Pr(dió|slapped) \times ----$