

Quiz Question!

- Suppose I'm making a language model with a vocabulary size of 20,000 words
- In my training data, I saw the bigram *comes across* 10 times [these are authentic counts from a small corpus, BTW]
 - 5 times it was followed by *as*
 - 5 times it was followed by other words (*like, less, again, most, in*)
- The next time I see *comes across*:
 - According to Good-Turing smoothing, what is the probability of seeing a different, previously unseen word following it?
 - Using absolute discounting with $D = 0.75$, what is the probability of seeing *as* after it? [Writing the answer as a fraction is okay.]

Machine Translation: Word alignment models

Christopher Manning
CS224N

[Based on slides by Kevin Knight, Dan Klein,
Dan Jurafsky]

Centauri/Arcturan [Knight, 1997]: It's Really Spanish/English

Clients do not sell pharmaceuticals in Europe => Clientes no venden medicinas en Europa

1a. Garcia and associates . 1b. Garcia y asociados .	7a. the clients and the associates are enemies . 7b. los clientes y los asociados son enemigos .
2a. Carlos Garcia has three associates . 2b. Carlos Garcia tiene tres asociados .	8a. the company has three groups . 8b. la empresa tiene tres grupos .
3a. his associates are not strong . 3b. sus asociados no son fuertes .	9a. its groups are in Europe . 9b. sus grupos estan en Europa .
4a. Garcia has a company also . 4b. Garcia tambien tiene una empresa .	10a. the modern groups sell strong pharmaceuticals . 10b. los grupos modernos venden medicinas fuertes .
5a. its clients are angry . 5b. sus clientes estan enfadados .	11a. the groups do not sell zenzanine . 11b. los grupos no venden zanzanina .
6a. the associates are also angry . 6b. los asociados tambien estan enfadados .	12a. the small groups are not modern . 12b. los grupos pequenos no son modernos .

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok errok hihok yorok klok kantok ok-yurp

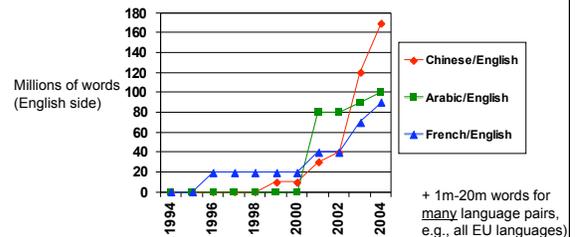
Your assignment, put these words in order: { jjat, arrat, mat, bat, oloat, at-yurp }

1a. ok-voon ororok sprok . 1b. at-voon bichat dat .	7a. lalok farok ororok lalok sprok izok enemok . 7b. wat jjat bichat wat dat vat eneak .
2a. ok-drubel ok-voon anak plok sprok . 2b. at-drubel at-voon pippat rrat dat .	8a. lalok brok anak plok nok . 8b. iat lat pippat rrat nmat .
3a. erok sprok izok hihok ghirok . 3b. totat dat arrat vat hilat .	9a. wiwok nok izok kantok ok-yurp . 9b. totat nmat quat oloat at-yurp .
4a. ok-voon anak drok brok jok . 4b. at-voon krat pippat sat lat .	10a. lalok mok nok yorok ghirok klok . 10b. wat nmat gat mat bat hilat .
5a. wiwok farok izok stok . 5b. totat jjat quat cat .	11a. lalok nok errok hihok yorok zanzanok . 11b. wat nmat arrat mat zanzañat .
6a. lalok sprok izok jok stok . 6b. wat dat krat quat cat .	12a. lalok rarok nok izok hihok mok . 12b. wat nmat forat arrat vat gat .

From No Data to Sentence Pairs

- Really hard way: pay \$\$\$
 - Suppose one billion words of parallel data were sufficient
 - At 20 cents/word, that's \$200 million
- Pretty hard way: Find it, and then earn it!
 - De-formatting
 - Remove strange characters
 - Character code conversion
 - Document alignment
 - **Sentence alignment**
 - **Tokenization (also called Segmentation)**
- Easy way: Linguistic Data Consortium (LDC)

Ready-to-Use Online Bilingual Data



(Data stripped of formatting, in sentence-pair format, available from the Linguistic Data Consortium at UPenn).

Tokenization (or Segmentation)

- English
 - Input (some character stream):
"There," said Bob.
 - Output (7 "tokens" or "words"):
" There , " said Bob .
- Chinese
 - Input (char stream): 美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯裔商拉登等发出的电子邮件。
 - Output: 美国 关岛 国际 机场 及其 办 公 室 均 接 获 一 名 自 称 沙 地 阿 拉 伯 裔 商 拉 登 等 发 出 的 电 子 邮 件 。

Sentence Alignment

The old man is happy. He has fished many times. His wife talks to him. The fish are jumping. The sharks await.	El viejo está feliz porque ha pescado muchos veces. Su mujer habla con él. Los tiburones esperan.
--	--

Sentence Alignment

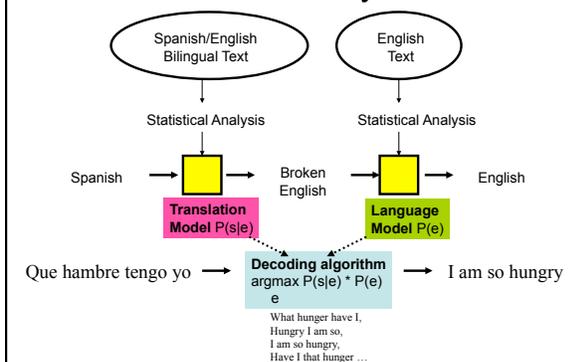
- | | |
|---------------------------------|---|
| 1. The old man is
happy. | 1. El viejo está feliz
porque ha
pescado muchos
veces. |
| 2. He has fished
many times. | 2. Su mujer habla
con él. |
| 3. His wife talks to
him. | 3. Los tiburones
esperan. |
| 4. The fish are
jumping. | |
| 5. The sharks await. | |

Sentence Alignment

- | | |
|---------------------------------|---|
| 1. The old man is
happy. | 1. El viejo está feliz
porque ha
pescado muchos
veces. |
| 2. He has fished
many times. | 2. Su mujer habla
con él. |
| 3. His wife talks to
him. | 3. Los tiburones
esperan. |
| 4. The fish are
jumping. | |
| 5. The sharks await. | |

Done by Dynamic Programming: see FSNLP ch. 13 for details

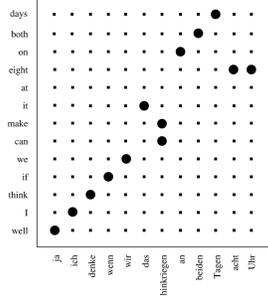
Statistical MT Systems



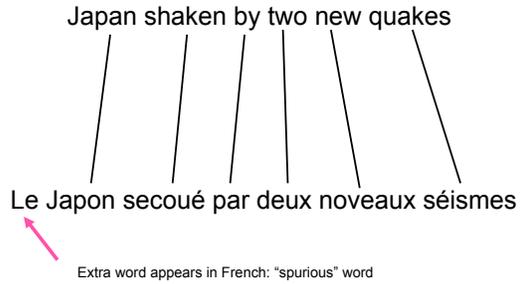
A division of labor

- Use of Bayes Rule ("the noisy channel model") allows a division of labor:
 - Job of the translation model $P(E|S)$ is just to model how various Spanish words typically get translated into English (perhaps in a certain context)
 - $P(E|S)$ doesn't have to worry about language-particular facts about English word order: that's the job of $P(E)$
 - The job of the language model is to choose felicitous bags of words and to correctly order them for English
 - $P(E)$ can do bag generation: putting a bag of words in order:
 - E.g., hungry I am so \rightarrow I am so hungry
- Both can be incomplete/sloppy

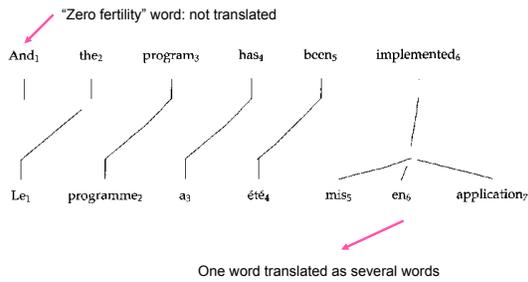
Word Alignment Examples: Grid



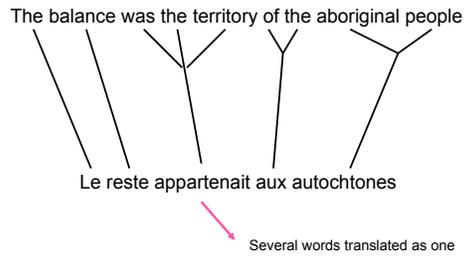
Word alignment examples: easy



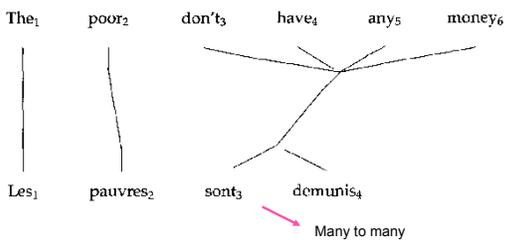
Alignments: harder



Alignments: harder

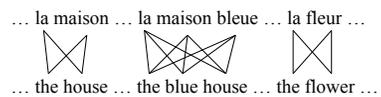


Alignments: hard



- A minimal aligned subset of words is called a 'cept' in the IBM work; often a 'bead' or '(aligned) statistical phrase' elsewhere.

Statistical Machine Translation



All word alignments equally likely

All $P(\text{french-word} \mid \text{english-word})$ equally likely

Statistical Machine Translation

... la maison ... la maison bleue ... la fleur ...

 ... the house ... the blue house ... the flower ...

"la" and "the" observed to co-occur frequently,
 so $P(\text{la} | \text{the})$ is increased.

Statistical Machine Translation

... la maison ... la maison bleue ... la fleur ...

 ... the house ... the blue house ... the flower ...

"house" co-occurs with both "la" and "maison", but
 $P(\text{maison} | \text{house})$ can be raised without limit, to 1.0,
 while $P(\text{la} | \text{house})$ is limited because of "the"

(pigeonhole principle)

Statistical Machine Translation

... la maison ... la maison bleue ... la fleur ...

 ... the house ... the blue house ... the flower ...

settling down after another iteration

Word alignment learning with EM

... la maison ... la maison bleue ... la fleur ...

 ... the house ... the blue house ... the flower ...

Hidden structure revealed by EM training!

That was IBM Model 1. For details, see later and:

- "A Statistical MT Tutorial Workbook" (Knight, 1999).
- "The Mathematics of Statistical Machine Translation" (Brown et al, 1993)
- Software: GIZA++

Statistical Machine Translation

... la maison ... la maison bleue ... la fleur ...

 ... the house ... the blue house ... the flower ...

$P(\text{juste} \text{fair}) = 0.411$
$P(\text{juste} \text{correct}) = 0.027$
$P(\text{juste} \text{right}) = 0.020$
...

NB! Confusing
 But true!

new French sentence → → Possible English translations, to be rescored by language model

IBM StatMT Translation Models

- IBM1 – lexical probabilities only
- IBM2 – lexicon plus absolute position
- HMM – lexicon plus relative position
- IBM3 – plus fertilities
- IBM4 – inverted relative position alignment
- IBM5 – non-deficient version of model 4

• All the models we discuss today handle 0:1, 1:0, 1:1, 1:n alignments *only*

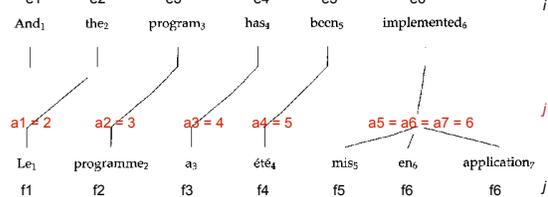
[Brown, et.al. 93, Vogel, et.al. 96]

IBM models 1,2,3,4,5

- Models for $P(F|E)$
- There is a set of English words and the extra English word NULL
- Each English word generates and places 0 or more French words
- Any remaining French words are deemed to have been produced by NULL

Model 1 parameters

- $P(F|E) = \prod_{(f,e)} P(f|e)$
- $P(f|e) = P(J|I) \sum_a P(f, a|e)$
- $P(f, a|e) = \prod_j P(a_j = i) P(f_j | e_j) = \prod_j [1/(i+1)] P(f_j | e_j)$



Model 1: Word alignment learning with Expectation-Maximization (EM)

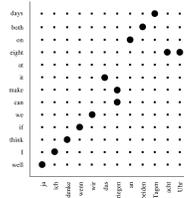
- Start with $P(f|e)$ uniform, including $P(f|NULL)$
- For each sentence
 - For each French position j
 - Calculate posterior over English positions $P(a_j|i)$

$$P(a_j = i | f, e) = \frac{P(f_j | e_i)}{\sum_{i'} P(f_j | e_{i'})}$$

- Increment count of word f_j with word e_{a_j}
 - $C(f_j e_{a_j}) += P(a_j = i | f, e)$
- Renormalize counts to give probs $P(f^p | e^q) = \frac{C(f^p | e^q)}{\sum_{j'} C(f^j | e^q)}$
- Iterate until convergence

IBM models 1,2,3,4,5

- In Model 2, the placement of a word in the French depends on where it was in the English



• Unlike Model 1, Model 2 captures the intuition that translations should usually "lie along the diagonal".

• The main focus of PA #2.

IBM models 1,2,3,4,5

- In model 3 we model how many French words an English word can produce, using a concept called fertility

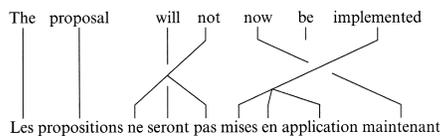
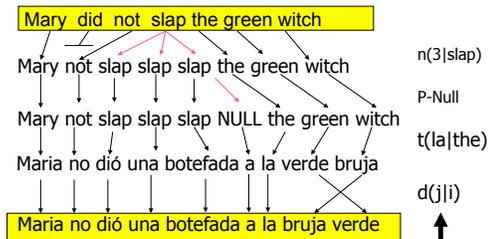


Figure 32.3 Alignment example.

IBM Model 3, Brown et al., 1993

Generative approach:



Probabilities can be learned from raw bilingual text.

IBM Model 3 (from Knight 1999)

- For each word e_i in English sentence, choose a **fertility** Φ_i . The choice of Φ_i depends only on e_i , not other words or Φ 's.
- For each word e_i , generate Φ_i Spanish words. Choice of French word depends only on English word e_i , not English context or any Spanish words.
- Permute all the Spanish words. Each Spanish word gets assigned absolute target position slot (1,2,3, etc). Choice of Spanish word position dependent only on absolute position of English word generating it.

Model 3: P(S|E) training parameters

- What are the parameters for this model?
- **Words:** $P(\text{casa}|\text{house})$
- **Spurious words:** $P(a|\text{null})$
- **Fertilities:** $n(1|\text{house})$: prob that "house" will produce 1 Spanish word whenever 'house' appears.
- **Distortions:** $d(5|2)$ prob. that English word in position 2 of English sentence generates French word in position 5 of French translation
 - Actually, distortions are $d(5|2,4,6)$ where 4 is length of English sentence, 6 is Spanish length

Spurious words

- We could have $n(3|\text{NULL})$ (probability of being exactly 3 spurious words in a Spanish translation)
- But instead, of $n(0|\text{NULL}), n(1|\text{NULL}) \dots n(25|\text{NULL})$, have a single parameter p_1
- After assign fertilities to non-NULL English words we want to generate (say) z Spanish words.
- As we generate each of z words, we optionally toss in spurious Spanish word with probability p_1
- Probability of not tossing in spurious word $p_0=1-p_1$

Distortion probabilities for spurious words

- Can't just have $d(5|0,4,6)$, i.e. chance that NULL word will end up in position 5.
- Why? These are spurious words! Could occur anywhere!! Too hard to predict
- Instead,
 - Use normal-word distortion parameters to choose positions for normally-generated Spanish words
 - Put Null-generated words into empty slots left over
 - If three NULL-generated words, and three empty slots, then there are 3!, or six, ways for slotting them all in
 - We'll assign a probability of 1/6 for each way

Real Model 3

- For each word e_i in English sentence, choose fertility Φ_i with prob $n(\Phi_i|e_i)$
- Choose number Φ_0 of spurious Spanish words to be generated from $e_0=\text{NULL}$ using p_1 and sum of fertilities from step 1
- Let m be sum of fertilities for all words including NULL
- For each $i=0,1,2,\dots,L$, $k=1,2,\dots,\Phi_i$:
 - choose Spanish word τ_{ik} with probability $t(\tau_{ik}|e_i)$
- For each $i=1,2,\dots,L$, $k=1,2,\dots,\Phi_i$:
 - choose target Spanish position π_{ik} with prob $d(\tau_{ik}|i,L,m)$
- For each $k=1,2,\dots,\Phi_0$ choose position π_{0k} from $\Phi_0 - k + 1$ remaining vacant positions in $1,2,\dots,m$ for total prob of $1/\Phi_0!$
- Output Spanish sentence with words τ_{ik} in positions π_{ik} ($0 \leq i \leq L, 1 \leq k \leq \Phi_i$)

Model 3 parameters

- n, t, p, d
- Again, if we had complete data of English strings and step-by-step rewritings into Spanish, we could:
 - Compute $n(0|\text{did})$ by locating every instance of "did", and seeing how many words it translates to
 - $t(\text{maison}|\text{house})$ how many of all French words generated by "house" were "maison"
 - $d(5|2,4,6)$ out of all times some word2 was translated, how many times did it become word5?

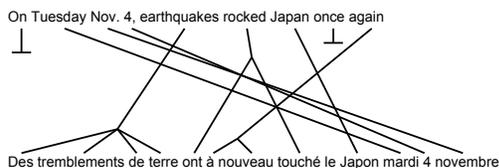
Since we don't have word-aligned data...

- We bootstrap alignments from incomplete data
- From a sentence-aligned bilingual corpus
 - 1) Assume some startup values for n, d, Φ , etc
 - 2) Use values for n, d, Φ , etc to use model 3 to work out chances of different possible alignments. Use these alignments to retrain n, d, Φ , etc
 - 3) Go to 2
- This is a more complicated case of the EM algorithm

IBM models 1,2,3,4,5

- In model 4 the placement of later French words produced by an English word depends on what happened to earlier French words generated by that same English word

Alignments: linguistics



IBM models 1,2,3,4,5

- In model 5 they do non-deficient alignment. That is, you can't put probability mass on impossible things.

Why all the models?

- We don't start with aligned text, so we have to get initial alignments from somewhere.
- Model 1 is words only, and is relatively easy and fast to train.
- We are working in a space with many local maxima, so output of model 1 can be a good place to start model 2. Etc.
- The sequence of models allows a better model to be found faster [the intuition is like deterministic annealing].

Alignments: linguistics



- There isn't enough linguistics to explain this in the translation model ... have to depend on the language model ... that may be unrealistic ... and may be harming our translation model