# Advanced Alignment Models

Philipp Koehn

11 February 2016

# IBM Model 1

- Generative model: break up translation process into smaller steps
  - IBM Model 1 only uses lexical translation

- Translation probability
  - for a foreign sentence $\mathbf{f} = (f_1, ..., f_{l_f})$ of length $l_f$
  - to an English sentence $\mathbf{e} = (e_1, ..., e_{l_e})$ of length $l_e$
  - with an alignment of each English word $e_j$ to a foreign word $f_i$ according to the alignment function $a : j \rightarrow i$

$$p(\mathbf{e}, a | \mathbf{f}) = \frac{\epsilon}{(l_f + 1)^{l_e}} \prod_{j=1}^{l_e} t(e_j | f_{a(j)})$$

  - parameter $\epsilon$ is a normalization constant
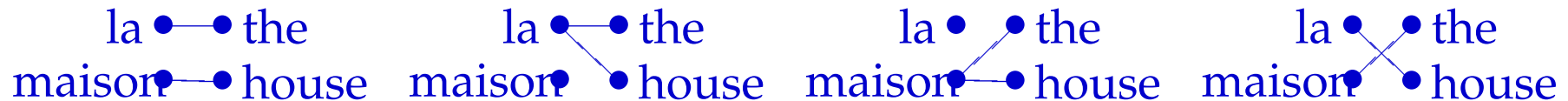
# IBM Model 1 and EM

- EM Algorithm consists of two steps

- Expectation-Step: Apply model to the data

  - parts of the model are hidden (here: alignments)
  - using the model, assign probabilities to possible values

- Maximization-Step: Estimate model from data

  - take assign values as fact
  - collect counts (weighted by probabilities)
  - estimate model from counts

- Iterate these steps until convergence

# IBM Model 1 and EM

- **Probabilities**

$$p(\text{the}|\text{la}) = 0.7 \qquad p(\text{house}|\text{la}) = 0.05$$
$$p(\text{the}|\text{maison}) = 0.1 \qquad p(\text{house}|\text{maison}) = 0.8$$

- **Alignments**



$$p(\mathbf{e}, a|\mathbf{f}) = 0.56 \qquad p(\mathbf{e}, a|\mathbf{f}) = 0.035 \qquad p(\mathbf{e}, a|\mathbf{f}) = 0.08 \qquad p(\mathbf{e}, a|\mathbf{f}) = 0.005$$

$$p(a|\mathbf{e}, \mathbf{f}) = 0.824 \qquad p(a|\mathbf{e}, \mathbf{f}) = 0.052 \qquad p(a|\mathbf{e}, \mathbf{f}) = 0.118 \qquad p(a|\mathbf{e}, \mathbf{f}) = 0.007$$

- **Counts**

$$c(\text{the}|\text{la}) = 0.824 + 0.052 \qquad c(\text{house}|\text{la}) = 0.052 + 0.007$$
$$c(\text{the}|\text{maison}) = 0.118 + 0.007 \qquad c(\text{house}|\text{maison}) = 0.824 + 0.118$$

- We need to be able to compute:

  – Expectation-Step: probability of alignments

  – Maximization-Step: count collection

- We need to compute $p(a|\mathbf{e}, \mathbf{f})$

- Applying the chain rule:

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

- We already have the formula for $p(\mathbf{e}, \mathbf{a}|\mathbf{f})$ (definition of Model 1)

- Now we have to collect counts

- Evidence from a sentence pair **e,f** that word $e$ is a translation of word $f$:

$$c(e|f; \mathbf{e}, \mathbf{f}) = \sum_a p(a|\mathbf{e}, \mathbf{f}) \sum_{j=1}^{l_e} \delta(e, e_j) \delta(f, f_{a(j)})$$
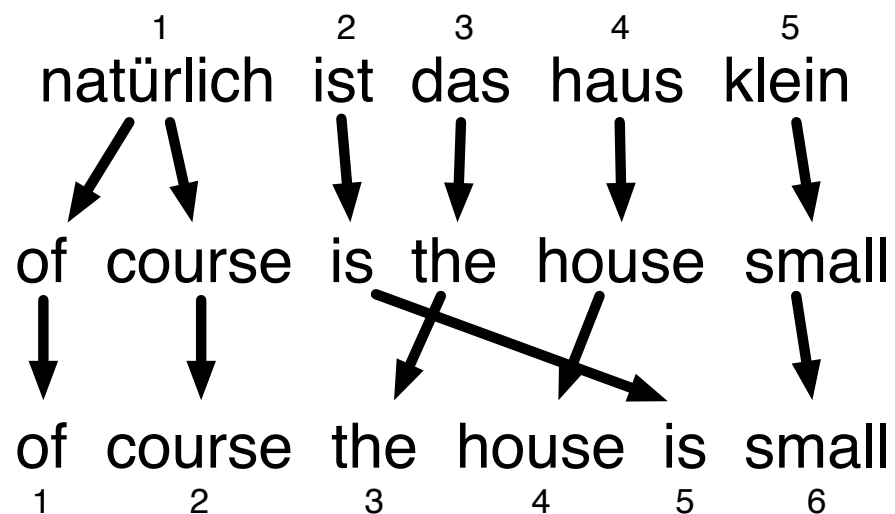
- With the same simplication as before:

$$c(e|f; \mathbf{e}, \mathbf{f}) = \frac{t(e|f)}{\sum_{i=0}^{l_f} t(e|f_i)} \sum_{j=1}^{l_e} \delta(e, e_j) \sum_{i=0}^{l_f} \delta(f, f_i)$$

# ibm model 2

Adding a model of alignment



lexical translation step

alignment step

- Modeling alignment with an alignment probability distribution

- Translating English word at position $j$ from foreign word at position $i = a(j)$:

$$a(i|j, l_e, l_f)$$

- Added to IBM Model 1

$$p(\mathbf{e}, a|\mathbf{f}) = \epsilon \prod_{j=1}^{l_e} t(e_j|f_{a(j)}) \, a(a(j)|j, l_e, l_f)$$

- Very similar to IBM Model 1 training

  - number of possible word alignments exponential with number of words
  - but: able to reduce complexity of computing $p(\mathbf{e}|\mathbf{f})$ to polynomial
  - same trick applies to IBM Model 2

$$p(\mathbf{e}|\mathbf{f}) = \sum_a p(\mathbf{e}, a|\mathbf{f})$$

$$= \epsilon \sum_{a(1)=0}^{l_f} ... \sum_{a(l_e)=0}^{l_f} \prod_{j=1}^{l_e} t(e_j|f_{a(j)}) \, a(a(j)|j, l_e, l_f)$$

$$= \epsilon \prod_{j=1}^{l_e} \sum_{i=0}^{l_f} t(e_j|f_{a(j)}) \, a(a(j)|j, l_e, l_f)$$

- Count collection for lexical translation

$$c(e|f; \mathbf{e}, \mathbf{f}) = \sum_{j=1}^{l_e} \sum_{i=0}^{l_f} \frac{t(e|f) \, a(a(j)|j, l_e, l_f) \, \delta(e, e_j) \, \delta(f, f_i)}{\sum_{i'=0}^{l_f} t(e|f_{i'}) \, a(i'|j, l_e, l_f))}$$

- Count collection for alignment

$$c(i|j, l_e, l_f; \mathbf{e}, \mathbf{f}) = \frac{t(e_j|f_i) \, a(a(j)|j, l_e, l_f)}{\sum_{i'=0}^{l_f} t(e_j|f_{i'}) \, a(i'|j, l_e, l_f))}$$

- Algorithm for training Model 2 is very similar to the one for IBM Model 1 (pseudo code in book)

- First run a few iterations of IBM Model 1 training

- Initialize probability distributions $t(e|f)$ and $a(i|j, l_e, l_f)$ from IBM Model 1

  – lexical translation probability distribution $t(e|f)$ can be taken verbatim
  – $a(i|j, l_e, l_f)$ initialized to $\frac{1}{l_f+1}$

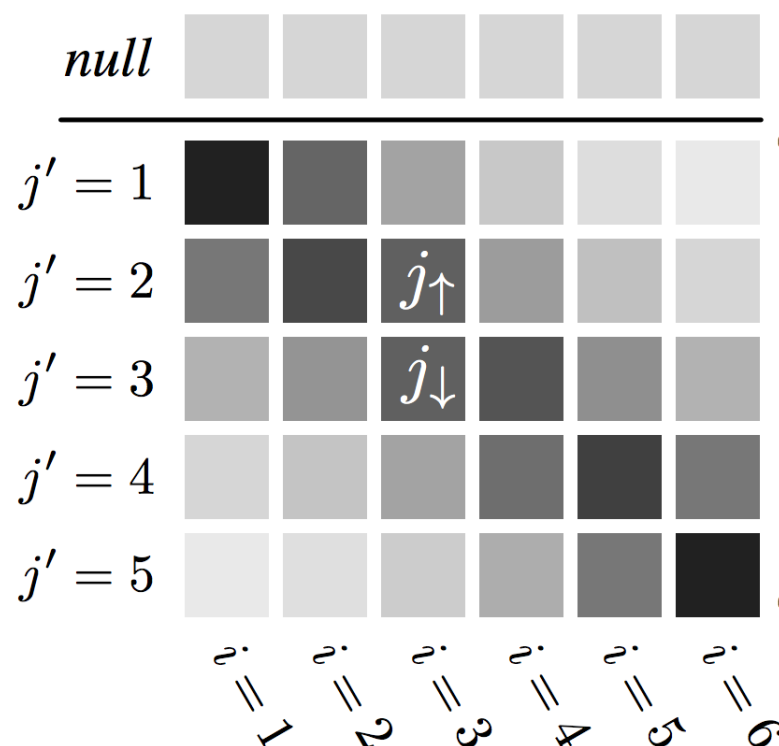# fast align:

# reparameterization of ibm model 2

# IBM Model 2: A Critique

- Alignment probability distribution has too many parameters $(l_e^2 l_f^2)$

$$a(i|j, l_e, l_f)$$

$\rightarrow$ too sparse data to estimate correctly

- Better: bias towards to diagonal

# Diagonal

- Distance from diagonal

$$h(i, j, l_e, l_f) = \left| \frac{i}{l_f} - \frac{j}{l_e} \right|$$

- Function that gives higher values to positions close to diagonal ($\lambda$ is a scaling factor)

$$e^{-\lambda h(i,j,l_e,l_f)}$$

- Special case: alignment to NULL token: $p_0$

- Alignment probability distribution

$$\delta(a(j) = i | j, l_e, l_f) = \begin{cases} p_0 & \text{if } i = 0 \\ (1 - p_0) \dfrac{e^{-\lambda h(i,j,l_e,l_f)}}{Z_\lambda(j,m,n)} & \text{if } 0 < i \le l_e \end{cases}$$

- This model was proposed by Dyer et al. (2013)

- It also changes the word translation probability distribution to include a prior

  – this was originally proposed by Mermer and Saraclar (2011)
  – an efficient estimation method (variational Bayes) was proposed by Riley and Gildea (2012)

- EM training is still simple

  – the probability to align an English word $e$ to a foreign word $f$ does not depend on the choices of other English words
  – the normalization function $Z_\lambda(j, m, n)$ can be computed in $O(1)$

# hmm model

# HMM Model

- Words do not move independently of each other

  - they often move in groups
  $\rightarrow$ condition word position on previous word's position

- HMM alignment model:

$$a(a(j)|a(j-1), l_e)$$

- EM algorithm application slightly harder, requires dynamic programming

- IBM Model 4 is similar, also conditions on word classes

# EM for the HMM Model

- Main objective: collect fractional counts to estimate

    - word translation probability distribution $t(e_j|f_{a(j)})$
    - alignment probability distribution $a(a(j)|a(j-1), l_e)$

- Consider all possible word alignments

- Collect evidence from each

- Exponentially many $\rightarrow$ need to do this efficiently

$j = 1$

$t(e_1|f_1) =$
$t(\text{of}|\text{natürlich})$

$a(a(1)|a(0)) =$
$a(1|0)$

$j = 2$

$t(e_2|f_1) =$
$t(\text{course}|\text{natürlich})$

$a(a(2)|a(1)) =$
$a(1|1)$

$j = 3$

$t(e_3|f_3) =$
$t(\text{the}|\text{das})$

$a(a(3)|a(2)) =$
$a(3|1)$

$j = 4$

$t(e_4|f_4) =$
$t(\text{house}|\text{hays})$

$a(a(4)|a(3)) =$
$a(4|3)$

$j = 5$

$t(e_5|f_2) =$
$t(\text{is}|\text{its})$

$a(a(5)|a(4)) =$
$a(2|4)$

$j = 6$

$t(e_6|f_5) =$
$t(\text{small}|\text{klein})$

$a(a(6)|a(5)) =$
$a(5|2)$

# First English Word

|  | **of** | **course** | **the** | ... |
|---|---|---|---|---|
|  | $j = 1$ | $j = 2$ | $j = 3$ | |

**natürlich**
$a(j) = 1$

$q_1(1) =$
$t(\text{of}|\text{natürlich})$
$\times a(1|0)$

**its**
$a(j) = 2$

$q_1(2) =$
$t(\text{of}|\text{ist})$
$\times a(2|0)$

**das**
$a(j) = 3$

$q_1(3) =$
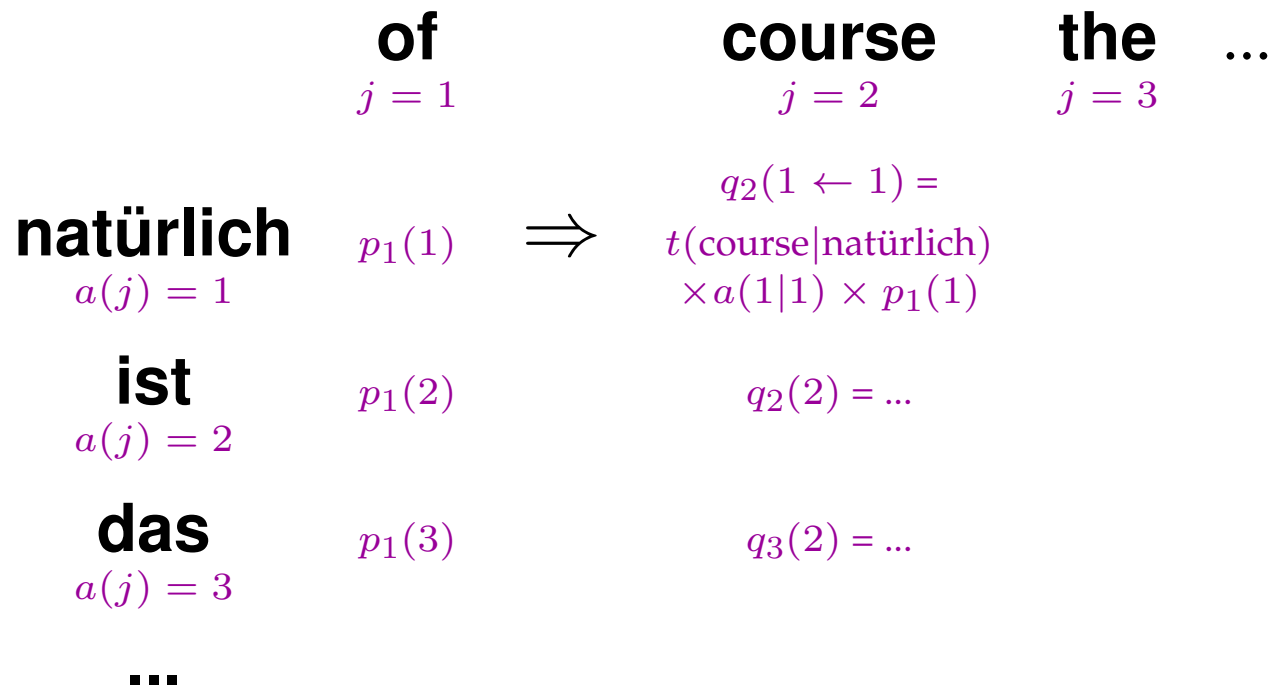$t(\text{of}|\text{das})$
$\times a(3|0)$

**...**

- Compute probabilities for each choice of $i = a(1)$ by normalizing $q_1(i)$

$$p_1(i) = \frac{q_1(i)}{\sum_{i'} q_1(i')}$$

- Use these probabilities for count collection for $t(\text{of}|\bullet)$ and $a(\bullet|0)$

# Next English Word

- One way to get there

$$
\begin{array}{ccc}
\textbf{of} & \textbf{course} & \textbf{the} \quad \text{...} \\
j = 1 & j = 2 & j = 3
\end{array}
$$

$$
q_2(1 \leftarrow 1) =
$$

$$
\textbf{natürlich} \quad p_1(1) \quad \Longrightarrow \quad t(\text{course}|\text{natürlich})
$$
$$
a(j) = 1 \qquad\qquad\qquad \times a(1|1) \times p_1(1)
$$

$$
\textbf{ist} \qquad p_1(2) \qquad\qquad q_2(2) = ...
$$
$$
a(j) = 2
$$

$$
\textbf{das} \qquad p_1(3) \qquad\qquad q_3(2) = ...
$$
$$
a(j) = 3
$$

$$
\textbf{...}
$$

- Another way to get there

$$\begin{array}{ccc}
\textbf{of} & \textbf{course} & \textbf{the} \quad ... \\
j = 1 & j = 2 & j = 3
\end{array}$$

$$
\begin{array}{cccc}
& & q_2(1 \leftarrow 2) = & \\
\textbf{natürlich} & p_1(1) & t(\text{course}|\text{natürlich}) & \\
a(j) = 1 & & \times a(1|2) \times p_1(2) & \\
& & & \\
\textbf{ist} & p_1(2) \quad \Longrightarrow & q_2(2) = ... & \\
a(j) = 2 & & & \\
& & & \\
\textbf{das} & p_1(3) & q_3(2) = ... & \\
a(j) = 3 & & & \\
& & & \\
\textbf{...} & & &
\end{array}
$$

- To compute the score of a state, we have to consider all of the paths

$$q_2(1) = t(e_2|f_1) \times \sum_i p_1(i) a(1|i)$$

# Summary of the Math

- Unnormalized score for a transition between two states

$$q_j(i \leftarrow i_{\text{previous}}) = t(e_j|f_i) \times a(i|i_{\text{previous}}) \times p_{j-1}(i_{\text{previous}})$$

- Normalization $\qquad p_j(i \leftarrow i_{\text{previous}}) = \dfrac{q_j(i \leftarrow i_{\text{previous}})}{\sum_{i,i_{\text{previous}}} q_j(i \leftarrow i_{\text{previous}})}$

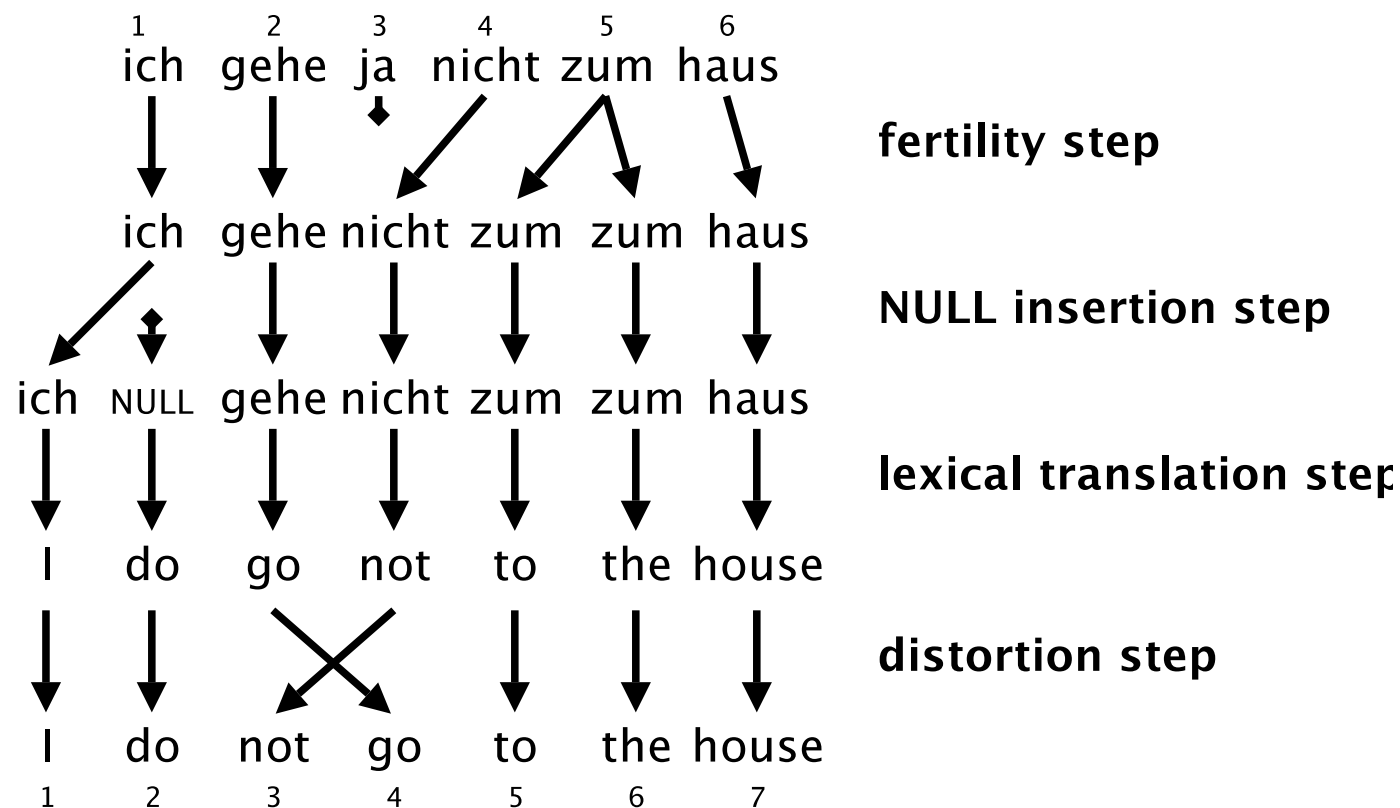- Probability of a state $\qquad p_j(i) = \sum_{i_{\text{previous}}} p_j(i \leftarrow i_{\text{previous}})$

- Count collection $\qquad c(e_j|f_i) = \sum_{i,j} p_j(i)$

$$c(i|i_{\text{previous}}) = \sum_{i,j,i_{\text{previous}}} p_j(i \leftarrow i_{\text{previous}})$$

# ibm model 3

## Adding a model of fertilty



fertility step

NULL insertion step

lexical translation step

distortion step

# IBM Model 3: Fertility

- Fertility: number of English words generated by a foreign word

- Modelled by distribution $n(\phi|f)$

- Example:

$$n(1|\text{haus}) \simeq 1$$
$$n(2|\text{zum}) \simeq 1$$
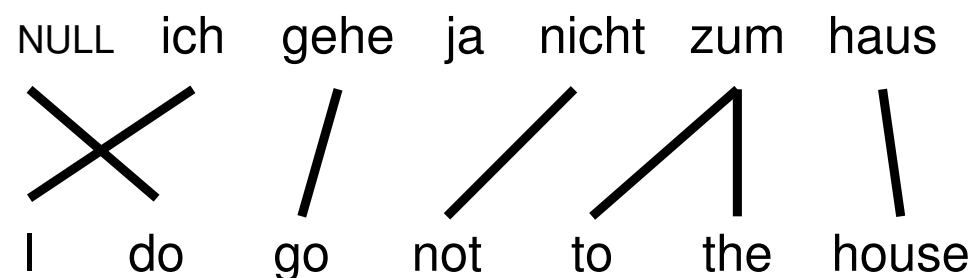$$n(0|\text{ja}) \simeq 1$$

- Training IBM Model 3 with the EM algorithm

  - The trick that reduces exponential complexity does not work anymore
  $\rightarrow$ Not possible to exhaustively consider all alignments

- Finding the most probable alignment by hillclimbing

  - start with initial alignment
  - change alignments for individual words
  - keep change if it has higher probability
  - continue until convergence

- Sampling: collecting variations to collect statistics

  - all alignments found during hillclimbing
  - neighboring alignments that differ by a move or a swap

# IBM Model 4

- Better reordering model

- Reordering in IBM Model 2 and 3

  - recall: $d(j|i, l_e, l_f)$
  - for large sentences (large $l_f$ and $l_e$), sparse and unreliable statistics
  - phrases tend to move together

- Relative reordering model: relative to previously translated words (cepts)

Foreign words with non-zero fertility forms cepts
(here 5 cepts)

NULL    ich    gehe    ja    nicht    zum    haus

I    do    go    not    to    the    house

| cept $\pi_i$ | $\pi_1$ | $\pi_2$ | $\pi_3$ | $\pi_4$ | $\pi_5$ |
|---|---|---|---|---|---|
| foreign position $[i]$ | 1 | 2 | 4 | 5 | 6 |
| foreign word $f_{[i]}$ | ich | gehe | nicht | zum | haus |
| English words $\{e_j\}$ | I | go | not | to,the | house |
| English positions $\{j\}$ | 1 | 4 | 3 | 5,6 | 7 |
| center of cept $\odot_i$ | 1 | 4 | 3 | 6 | 7 |

# IBM Model 4: Relative Distortion

| $j$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| $e_j$ | I | do | not | go | to | the | house |
| in cept $\pi_{i,k}$ | $\pi_{1,0}$ | $\pi_{0,0}$ | $\pi_{3,0}$ | $\pi_{2,0}$ | $\pi_{4,0}$ | $\pi_{4,1}$ | $\pi_{5,0}$ |
| $\odot_{i-1}$ | 0 | - | 4 | 1 | 3 | - | 6 |
| $j - \odot_{i-1}$ | +1 | - | $-1$ | +3 | +2 | - | +1 |
| distortion | $d_1(+1)$ | 1 | $d_1(-1)$ | $d_1(+3)$ | $d_1(+2)$ | $d_{>1}(+1)$ | $d_1(+1)$ |

- Center $\odot_i$ of a cept $\pi_i$ is $\mathrm{ceiling}(\mathrm{avg}(j))$

- Three cases:

  - uniform for NULL generated words
  - first word of a cept: $d_1$
  - next words of a cept: $d_{>1}$

- Some words may trigger reordering $\rightarrow$ condition reordering on words

$$\text{for initial word in cept:} \quad d_1(j - \odot_{[i-1]}|f_{[i-1]}, e_j)$$

$$\text{for additional words:} \quad d_{>1}(j - \Pi_{i,k-1}|e_j)$$

- Sparse data concerns $\rightarrow$ cluster words into classes

$$\text{for initial word in cept:} \quad d_1(j - \odot_{[i-1]}|\mathcal{A}(f_{[i-1]}), \mathcal{B}(e_j))$$

$$\text{for additional words:} \quad d_{>1}(j - \Pi_{i,k-1}|\mathcal{B}(e_j))$$

- IBM Models 1–4 are *deficient*

  - some impossible translations have positive probability
  - multiple output words may be placed in the same position
  $\rightarrow$ probability mass is wasted

- IBM Model 5 fixes deficiency by keeping track of vacancies (available positions)

- IBM Models were the pioneering models in statistical machine translation

- Introduced important concepts

  - generative model
  - EM training
  - reordering models

- Only used for niche applications as translation model

- ... but still in common use for word alignment (e.g., GIZA++ toolkit)

# word alignment

# Word Alignment

Given a sentence pair, which words correspond to each other?

# Word Alignment?



Is the English word does aligned to
the German wohnt (verb) or nicht (negation) or neither?

# Word Alignment?

How do the idioms kicked the bucket and biss ins grass match up?
Outside this exceptional context, bucket is never a good translation for grass

# Measuring Word Alignment Quality

- Manually align corpus with *sure* ($S$) and *possible* ($P$) alignment points ($S \subseteq P$)

- Common metric for evaluation word alignments: Alignment Error Rate (AER)

$$\text{AER}(S, P; A) = \frac{|A \cap S| + |A \cap P|}{|A| + |S|}$$

- AER = 0: alignment $A$ matches all sure, any possible alignment points

- However: different applications require different precision/recall trade-offs

# symmetrization

- IBM Models create a **many-to-one** mapping

  – words are aligned using an alignment function

  – a function may return the same value for different input
    (one-to-many mapping)

  – a function can not return multiple values for one input
    (no many-to-one mapping)

- Real word alignments have **many-to-many** mappings

# Symmetrization

- Run IBM Model training in both directions

$\rightarrow$ two sets of word alignment points

- Intersection: high precision alignment points

- Union: high recall alignment points

- Refinement methods explore the sets between intersection and union

# Example



english to spanish

spanish to english

intersection

**black**: intersection        **grey**: additional points in union

- Add alignment points from union based on heuristics:
  - directly/diagonally neighboring points
  - finally, add alignments that connect unaligned words in source and/or target
- Popular method: grow-diag-final-and

**grow-diag-final**(e2f,f2e)

1: neighboring = {(-1,0),(0,-1),(1,0),(0,1),(-1,-1),(-1,1),(1,-1),(1,1)}
2: alignment $A$ = intersect(e2f,f2e); grow-diag(); final(e2f); final(f2e);

**grow-diag**()

1: **while** new points added **do**
2:     **for all** English word $e \in [1...e_n]$, foreign word $f \in [1...f_n]$, $(e, f) \in A$ **do**
3:         **for all** neighboring alignment points $(e_{\text{new}}, f_{\text{new}})$ **do**
4:             **if** ($e_{\text{new}}$ unaligned OR $f_{\text{new}}$ unaligned) AND $(e_{\text{new}}, f_{\text{new}}) \in$ union(e2f,f2e) **then**
5:                 add $(e_{\text{new}}, f_{\text{new}})$ to $A$
6:             **end if**
7:         **end for**
8:     **end for**
9: **end while**

**final**()

1: **for all** English word $e_{\text{new}} \in [1...e_n]$, foreign word $f_{\text{new}} \in [1...f_n]$ **do**
2:     **if** ($e_{\text{new}}$ unaligned OR $f_{\text{new}}$ unaligned) AND $(e_{\text{new}}, f_{\text{new}}) \in$ union(e2f,f2e) **then**
3:         add $(e_{\text{new}}, f_{\text{new}})$ to $A$
4:     **end if**
5: **end for**

# More Work on Symmetrization

- Symmetrize after each iteration of IBM Models [Matusov et al., 2004]

  - run one iteration of E-step for each direction
  - symmetrize the two directions
  - count collection (M-step)

- Use of posterior probabilities in symmetrization

  - generate n-best alignments for each direction
  - calculate how often an alignment point occurs in these alignments
  - use this posterior probability during symmetrization

# Link Deletion / Addition Models

- Link deletion [Fossum et al., 2008]

  – start with union of IBM Model alignment points
  – delete one alignment point at a time
  – uses a neural network classifiers that also considers aspects such as how useful the alignment is for learning translation rules


- Link addition [Ren et al., 2007] [Ma et al., 2008]

  – possibly start with a skeleton of highly likely alignment points
  – add one alignment point at a time

- Given some annotated training data, supervised learning methods are possible

- Structured prediction

  – not just a classification problem
  – solution structure has to be constructed in steps

- Many approaches: maximum entropy, neural networks, support vector machines, conditional random fields, MIRA, ...

- Small labeled corpus may be used for parameter tuning of unsupervised aligner [Fraser and Marcu, 2007]

# Better Generative Models

- Aligning phrases

  – joint model [Marcu and Wong, 2002]
  – problem: EM algorithm likes really long phrases

- Fraser's LEAF

  – decomposes word alignment into many steps
  – similar in spirit to IBM Models
  – includes step for grouping into phrase

- Riesa's NILE

  – use syntactic parse trees to guide word alignment
  – build up words bottom up following the parse tree

- Research on word alignment has recently picked up again

  – speed matters
  – incremental ("online") training


- Unclear link betwwn

  – word alignment quality measured against manual gold standard
  – impact on machine translation quality


- Advice: if you develop method, make easy-to-use toolkit available