

# CMPE 409 Machine Translation

## IBM Models and Alignments

Murat ORHUN

(most of them adapted from the original slides of Prof. Philipp Koehn)

Istanbul Bilgi University

April 13, 2022

## 1 IBM Models

- IBM Model 1
- IBM Model 2
- IBM Model 3
- IBM Model 4
- IBM Model 5

## 2 Alignment

## 3 Further Reading

## 4 References

# 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  $\mathbf{f} = (f_1, \dots, f_{l_f})$  of length  $l_f$
  - to an English sentence  $\mathbf{e} = (e_1, \dots, e_{l_e})$  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

# Translation Probability

das	
$e$	$t(e f)$
the	0.7
that	0.15
which	0.075
who	0.05
this	0.025

Haus	
$e$	$t(e f)$
house	0.8
building	0.16
home	0.02
household	0.015
shell	0.005

ist	
$e$	$t(e f)$
is	0.8
's	0.16
exists	0.02
has	0.015
are	0.005

klein	
$e$	$t(e f)$
small	
little	
short	
minor	
petty	

$$\begin{aligned}
 p(e, a|f) &= \frac{\epsilon}{4^3} \times t(\text{the}|\text{das}) \times t(\text{house}|\text{Haus}) \times t(\text{is}|\text{ist}) \times t(\text{small}|\text{klein}) \\
 &= \frac{\epsilon}{4^3} \times 0.7 \times 0.8 \times 0.8 \times 0.4 \\
 &= 0.0028\epsilon
 \end{aligned}$$

# 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

- We need to be able to compute:
  - Expectation-Step: probability of alignments
  - Maximization-Step: count collection

# IBM Model 1 and EM: Pseudocode

**Input:** set of sentence pairs  $(\mathbf{e}, \mathbf{f})$

**Output:** translation prob.  $t(e|f)$

```

1: initialize  $t(e|f)$  uniformly
2: while not converged do
3:   // initialize
4:    $\text{count}(e|f) = 0$  for all  $e, f$ 
5:    $\text{total}(f) = 0$  for all  $f$ 
6:   for all sentence pairs  $(\mathbf{e}, \mathbf{f})$  do
7:     // compute normalization
8:     for all words  $e$  in  $\mathbf{e}$  do
9:        $\text{s-total}(e) = 0$ 
10:      for all words  $f$  in  $\mathbf{f}$  do
11:         $\text{s-total}(e) += t(e|f)$ 
12:      end for
13:    end for
```

```

14:   // collect counts
15:   for all words  $e$  in  $\mathbf{e}$  do
16:     for all words  $f$  in  $\mathbf{f}$  do
17:        $\text{count}(e|f) += \frac{t(e|f)}{\text{s-total}(e)}$ 
18:        $\text{total}(f) += \frac{t(e|f)}{\text{s-total}(e)}$ 
19:     end for
20:   end for
21:   // estimate probabilities
22:   for all foreign words  $f$  do
23:     for all English words  $e$  do
24:        $t(e|f) = \frac{\text{count}(e|f)}{\text{total}(f)}$ 
25:     end for
26:   end for
27:   end for
28: end while
```

# Conversion

das Haus  
the house

das Buch  
the book

ein Buch  
a book

<i>e</i>	<i>f</i>	initial	1st it.	2nd it.	3rd it.	...	final
the	das	0.25	0.5	0.6364	0.7479	...	1
book	das	0.25	0.25	0.1818	0.1208	...	0
house	das	0.25	0.25	0.1818	0.1313	...	0
the	buch	0.25	0.25	0.1818	0.1208	...	0
book	buch	0.25	0.5	0.6364	0.7479	...	1
a	buch	0.25	0.25	0.1818	0.1313	...	0
book	ein	0.25	0.5	0.4286	0.3466	...	0
a	ein	0.25	0.5	0.5714	0.6534	...	1
the	haus	0.25	0.5	0.4286	0.3466	...	0
house	haus	0.25	0.5	0.5714	0.6534	...	1



# Perplexity

- Perplexity is a measurement of how well a probability distribution or probability model predicts a sample.
- A low perplexity indicates the probability distribution is good at predicting the sample.

$$\log_2 PP = - \sum_s \log_2 p(\mathbf{e}_s | \mathbf{f}_s)$$

# Example

das

$e$	$t(e f)$
the	0.7
that	0.15
which	0.075
who	0.05
this	0.025

Haus

$e$	$t(e f)$
house	0.8
building	0.16
home	0.02
household	0.015
shell	0.005

ist

$e$	$t(e f)$
is	0.8
's	0.16
exists	0.02
has	0.015
are	0.005

klein

$e$	$t(e f)$
small	0.4
little	0.4
short	0.1
minor	0.06
petty	0.04

$$\begin{aligned}
 p(e, a|f) &= \frac{\epsilon}{4^3} \times t(\text{the}|\text{das}) \times t(\text{house}|\text{Haus}) \times t(\text{is}|\text{ist}) \times t(\text{small}|\text{klein}) \\
 &= \frac{\epsilon}{4^3} \times 0.7 \times 0.8 \times 0.8 \times 0.4 \\
 &= 0.0028\epsilon
 \end{aligned}$$

# Perplexity

	Initial	1st it.	2nd it.	3rd it.	...	Final
$p(\text{the haus} \text{das haus})$	0.0625	0.1875	0.1905	0.1913	...	0.1875
$p(\text{the book} \text{das buch})$	0.0625	0.1406	0.1790	0.2075	...	0.25
$p(\text{a book} \text{ein buch})$	0.0625	0.1875	0.1907	0.1913	...	0.1875
perplexity ( $\epsilon=1$ )	4095	202.3	153.6	131.6	...	113.8

# Fluent Output

- small step, little step
- which one is better?
- show google result
- language model

# Language Model

- relying on Google counts is a neat trick
- longer and longer sentences
- generative modeling
- The most common method for language modeling is the use of **n-gram** language models.

# Trigram language models

$$\begin{aligned}
 p(\mathbf{e}) &= p(e_1, e_2, \dots, e_n) \\
 &= p(e_1)p(e_2|e_1) \cdots p(e_n|e_1, e_2, \dots, e_{n-1}) \\
 &\simeq p(e_1)p(e_2|e_1) \cdots p(e_n|e_{n-2}, e_{n-1})
 \end{aligned}$$

# Noisy-Channel Model

Combining a language model and translation model in this way is called the noisy-channel model.

$$\begin{aligned}\operatorname{argmax}_{\mathbf{e}} p(\mathbf{e}|\mathbf{f}) &= \operatorname{argmax}_{\mathbf{e}} \frac{p(\mathbf{f}|\mathbf{e})p(\mathbf{e})}{p(\mathbf{f})} \\ &= \operatorname{argmax}_{\mathbf{e}} p(\mathbf{f}|\mathbf{e})p(\mathbf{e})\end{aligned}$$

## Higher IBM Models

IBM Model 1	lexical translation
IBM Model 2	adds absolute reordering model
IBM Model 3	adds fertility model
IBM Model 4	relative reordering model
IBM Model 5	fixes deficiency

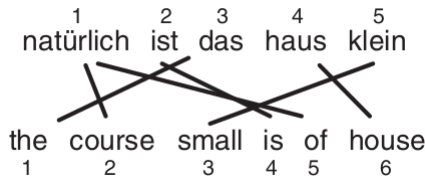
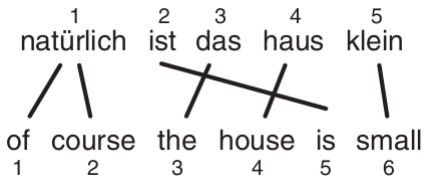


# Fertility

- input words produce a specific number of output words in the output language.
- But some words produce multiple words or get dropped (producing zero words).
- A model for the fertility of words addresses this aspect of translation.

## IBM Model 2

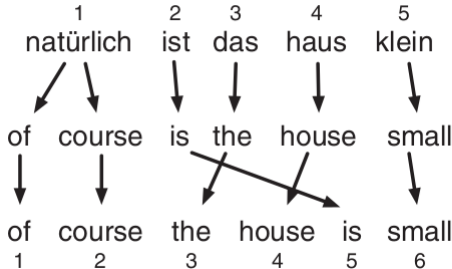
In IBM Model 1 the translation probabilities for the following two alternative translations are the same



## IBM Model 2

- IBM Model 2 addresses the issue of alignment
- The translation of a foreign input word in position  $i$  to an English word in position  $j$  is modeled by an alignment probability distribution
- We can view translation under IBM Model 2 as a two-step process with a lexical translation step and an alignment step:

# IBM Model 2



**lexical translation step**

**alignment step**

# IBM Model 2

- Modeling alignment with an alignment probability distribution
- Translating foreign word at position  $i$  to English word at position  $j$ :

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

- Putting everything together

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

- EM training of this model works the same way as IBM Model 1

## IBM Model 3

- how many words are generated from each input word.
- some German words like **zum** typically translate to two English words, i.e., **to the**.
- Others, such as the flavoring particle **ja**, get dropped.

# IBM Model 3

Model the fertility of input words directly with a probability distribution:

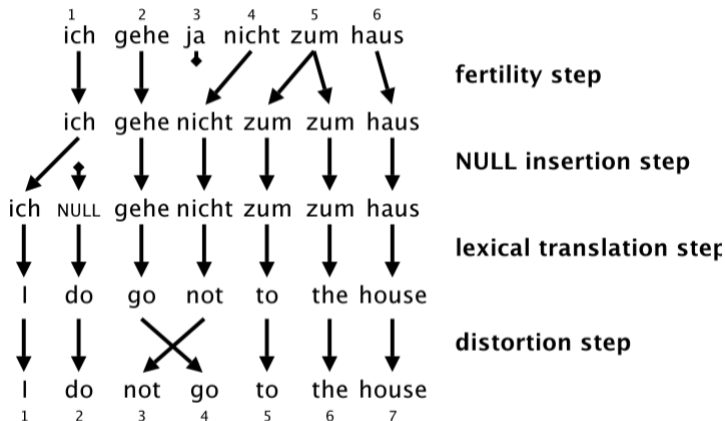
$$n(\phi|f)$$

$$n(1|\text{haus}) \simeq 1$$

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

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

# Adding a model of fertility





# Sampling the Alignment Space

- training IBM Model 3 with the EM algorithm
- Finding the most probable alignment by hillclimbing
- Sampling: collecting variations to collect statistics

## Model 3 is already a powerful model

- Translation of words
- Reordering (distortion)
- Insertion of words ( NULL insertion)
- Dropping of words (words with fertility 0)
- One-to-many translation (fertility)
- **IBM Model 4 further improves on Model 3**

## Problem with IBM Model 3

- Distortion probability
- Large input and outputs
- Sparse data
- Relations sentences

## IBM Model 4:Relative Distortion

In this model, the placement of the translation of an input word is typically based on the placement of the translation of the proceeding input word.

# IBM Model 4

- Relative Distortion
- Word classes
- Not all words need to be reordered

# IBM Model 5

- IBM Models 1–4 are *deficient*
  - some impossible translations have positive probability
  - multiple output words may be placed in the same position
  - probability mass is wasted
- IBM Model 5 fixes deficiency by keeping track of vacancies (available positions)

# Conclusion

- 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)

# Alignment Matrix

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

	michael	geht	davon	aus	,	dass	er	im	haus	bleibt
michael										
assumes										
that										
he										
will										
stay										
in										
the										
house										



















# Word Alignment

	john	wohnt	hier	nicht
john				
does		?		?
not				
live				
here				

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

# Word Alignment

	john	biss	ins	grass
john				
kicked				
the				
bucket				

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)

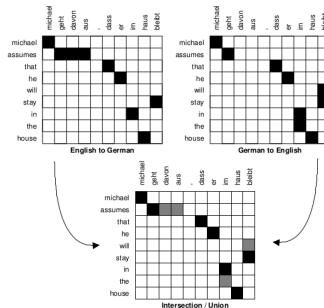
$$AER(S, P; A) = \frac{|A \cap S| + |A \cap P|}{|A| + |S|}$$

- AER = 0: alignment  $A$  matches all sure, any possible alignment points
- However: different applications require different precision/recall trade-offs

# Word Alignment with IBM Models

- 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

# Symmetrizing Word Alignments



- Intersection of GIZA++ bidirectional alignments
- Grow additional alignment points [Och and Ney, CompLing2003]

## Further Reading

- Chapter 4 of Koehn. P. (200p)

# Recourse

- Jurafsky, D. and J. H. Martin. Speech and language processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition, Second Edition, Upper Saddle River, NJ: Prentice-Hall, 2008.
- Koehn, P. (2009). Statistical Machine Translation. Cambridge: Cambridge University Press.  
doi:10.1017/CBO9780511815829