# Natural Language Processing

### Anoop Sarkar
anoopsarkar.github.io/nlp-class

Simon Fraser University

October 25, 2014

# Natural Language Processing

Anoop Sarkar

anoopsarkar.github.io/nlp-class

Simon Fraser University

Part 1: Generative Models for Word Alignment

# Statistical Machine Translation
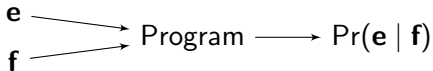
# Statistical Machine Translation

Noisy Channel Model

$$\mathbf{e}^* \; = \; \arg\max_{\mathbf{e}} \; \underbrace{\Pr(\mathbf{e})}_{\textbf{\textcolor{red}{Language Model}}} \; \cdot \; \underbrace{\Pr(\mathbf{f} \mid \mathbf{e})}_{\textbf{\textcolor{blue}{Alignment Model}}}$$

## Alignment Task

$$\mathbf{e} \searrow$$
$$\text{Program} \longrightarrow \Pr(\mathbf{e} \mid \mathbf{f})$$
$$\mathbf{f} \nearrow$$

## Training Data

- **Alignment Model**: learn a mapping between **f** and **e**.
  Training data: lots of translation pairs between **f** and **e**.

# Statistical Machine Translation

## The IBM Models

- The first statistical machine translation models were developed at IBM Research (Yorktown Heights, NY) in the 1980s

- The models were published in 1993:
  Brown et. al. The Mathematics of Statistical Machine Translation. *Computational Linguistics*. 1993.
  http://aclweb.org/anthology/J/J93/J93-2003.pdf

- These models are the basic SMT models, called:
  IBM Model 1, IBM Model 2, IBM Model 3, IBM Model 4, IBM Model 5
  as they were called in the 1993 paper.

- We use **e** and **f** in the equations in honor of their system which translated from French to English.
  Trained on the Canadian Hansards (Parliament Proceedings)

Generative Model of Word Alignment
    Word Alignments: IBM Model 3
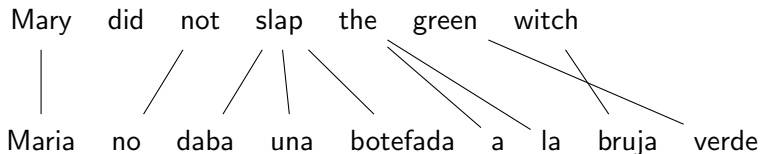    Word Alignments: IBM Model 1
    Finding the best alignment: IBM Model 1
    Learning Parameters: IBM Model 1

# Generative Model of Word Alignment

- English **e**: Mary did not slap the green witch
- "French" **f**: Maria no daba una botefada a la bruja verde
- Alignment **a**: $\{1, 3, 4, 4, 4, 5, 5, 7, 6\}$
  e.g. $(f_8, e_{a_8}) = (f_8, e_7) = (\text{bruja}, \text{witch})$

## Visualizing alignment **a**

# Generative Model of Word Alignment

## Data Set

- Data set $\mathcal{D}$ of $N$ sentences:
  $$\mathcal{D} = \{(\mathbf{f}^{(1)}, \mathbf{e}^{(1)}), \ldots, (\mathbf{f}^{(N)}, \mathbf{e}^{(N)})\}$$

- French $\mathbf{f}$: $(f_1, f_2, \ldots, f_I)$
- English $\mathbf{e}$: $(e_1, e_2, \ldots, e_J)$
- Alignment $\mathbf{a}$: $(a_1, a_2, \ldots, a_I)$

# Generative Model of Word Alignment

Find the best alignment for each translation pair

$$\mathbf{a}^* = \arg\max_{\mathbf{a}} \Pr(\mathbf{a} \mid \mathbf{f}, \mathbf{e})$$

Alignment probability

$$
\begin{aligned}
\Pr(\mathbf{a} \mid \mathbf{f}, \mathbf{e}) &= \frac{\Pr(\mathbf{f}, \mathbf{a}, \mathbf{e})}{\Pr(\mathbf{f}, \mathbf{e})} \\
&= \frac{\Pr(\mathbf{e}) \Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})}{\Pr(\mathbf{e}) \Pr(\mathbf{f} \mid \mathbf{e})} \\
&= \frac{\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})}{\Pr(\mathbf{f} \mid \mathbf{e})} \\
&= \frac{\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})}{\sum_{\mathbf{a}} \Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})}
\end{aligned}
$$

Generative Model of Word Alignment
    Word Alignments: IBM Model 3
    Word Alignments: IBM Model 1
    Finding the best alignment: IBM Model 1
    Learning Parameters: IBM Model 1

# Word Alignments: IBM Model 3

## Generative "story" for $P(\mathbf{f}, \mathbf{a} \mid \mathbf{e})$

# Word Alignments: IBM Model 3

Fertility parameter

$$n(\phi_j \mid e_j) : n(3 \mid \textit{slap}); n(0 \mid \textit{did})$$

Translation parameter

$$t(f_i \mid e_{a_i}) : t(\textit{bruja} \mid \textit{witch})$$

Distortion parameter

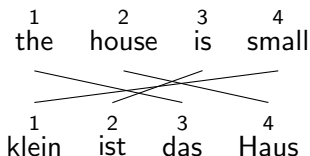$$d(f_{pos} = i \mid e_{pos} = j, I, J) : d(8 \mid 7, 8, 6)$$

# Word Alignments: IBM Model 3

Generative model for $P(\mathbf{f}, \mathbf{a} \mid \mathbf{e})$

$$
\begin{aligned}
P(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) &= \prod_{j=1}^{J} n(\phi_i \mid e_i) \\
&\times \prod_{i=1}^{I} t(f_i \mid e_{a_j}) \\
&\times \prod_{i=1}^{I} d(i \mid j, I, J)
\end{aligned}
$$

# Word Alignments: IBM Model 3

## Sentence pair with alignment $\mathbf{a} = (4, 3, 1, 2)$



If we know the parameter values we can easily compute the probability of this aligned sentence pair.

$\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) =$

$$
\begin{array}{llll}
n(1 \mid \text{the}) & \times & t(\text{das} \mid \text{the}) & \times & d(3 \mid 1, 4, 4) \times \\
n(1 \mid \text{house}) & \times & t(\text{Haus} \mid \text{house}) & \times & d(4 \mid 2, 4, 4) \times \\
n(1 \mid \text{is}) & \times & t(\text{ist} \mid \text{is}) & \times & d(2 \mid 3, 4, 4) \times \\
n(1 \mid \text{small}) & \times & t(\text{klein} \mid \text{small}) & \times & d(1 \mid 4, 4, 4)
\end{array}
$$

# Word Alignments: IBM Model 3



Parameter Estimation

- What is $n(1 \mid \text{very}) = ?$ and $n(0 \mid \text{very}) = ?$
- What is $t(\text{Haus} \mid \text{house}) = ?$ and $t(\text{klein} \mid \text{small}) = ?$
- What is $d(1 \mid 4, 4, 4) = ?$ and $d(1 \mid 1, 4, 4) = ?$

# Word Alignments: IBM Model 3



Parameter Estimation: Sum over all alignments

$$\sum_{\mathbf{a}} \Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) = \sum_{\mathbf{a}} \prod_{i=1}^{l} n(\phi_{a_i} \mid e_{a_i}) \times t(f_i \mid e_{a_i}) \times d(i \mid a_i, \mathbf{f}_{\text{len}}, \mathbf{e}_{\text{len}})$$

# Word Alignments: IBM Model 3

## Summary

- If we know the parameter values we can easily compute the probability $\Pr(\mathbf{a} \mid \mathbf{f}, \mathbf{e})$ given an aligned sentence pair
- If we are given a corpus of sentence pairs with alignments we can easily learn the parameter values by using relative frequencies.
- If we do not know the alignments then perhaps we can produce all possible alignments each with a certain probability?

IBM Model 3 is too hard: Let us try learning only $t(f_i \mid e_{a_i})$

$$\sum_{\mathbf{a}} \Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) = \sum_{\mathbf{a}} \prod_{i=1}^{I} n(\phi_{a_i} \mid e_{a_i}) \times t(f_i \mid e_{a_i}) \times d(i \mid a_i, \mathbf{f}_{\text{len}}, \mathbf{e}_{\text{len}})$$

Generative Model of Word Alignment
    Word Alignments: IBM Model 3
    Word Alignments: IBM Model 1
    Finding the best alignment: IBM Model 1
    Learning Parameters: IBM Model 1

# Word Alignments: IBM Model 1

### Alignment probability

$$\Pr(\mathbf{a} \mid \mathbf{f}, \mathbf{e}) = \frac{\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})}{\sum_{\mathbf{a}} \Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})}$$

### Example alignment

```
1      2       3    4
the   house    is   small
|      /       |    |
1      2       3    4
das   Haus    ist  klein
```

$$\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) = \prod_{i=1}^{l} t(f_i \mid e_{a_i})$$

$\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) =$
  $t(\text{das} \mid \text{the}) \times$
  $t(\text{Haus} \mid \text{house}) \times$
  $t(\text{ist} \mid \text{is}) \times$
  $t(\text{klein} \mid \text{small})$

# Word Alignments: IBM Model 1

### Generative "story" for Model 1

the    house   is   small

↓        ↓      ↓    ↓

das   Haus   ist   klein   (translate)

$$\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) = \prod_{i=1}^{I} t(f_i \mid e_{a_i})$$

Generative Model of Word Alignment
    Word Alignments: IBM Model 3
    Word Alignments: IBM Model 1
    Finding the best alignment: IBM Model 1
    Learning Parameters: IBM Model 1
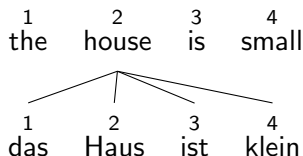
# Finding the best word alignment: IBM Model 1

Compute the arg max word alignment

$$\hat{\mathbf{a}} = \arg \max_{\mathbf{a}} \Pr(\mathbf{a} \mid \mathbf{e}, \mathbf{f})$$
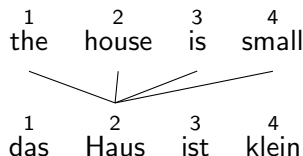
▶ For each $f_i$ in $(f_1, \ldots, f_I)$ build $\mathbf{a} = (\hat{a_1}, \ldots, \hat{a_I})$

$$\hat{a_i} = \arg \max_{a_i} t(f_i \mid e_{a_i})$$

Many to one alignment ✓

```
1      2      3     4
the  house    is  small


1      2      3     4
das   Haus   ist  klein
```

One to many alignment ✗

```
1      2      3     4
the  house    is  small


1      2      3     4
das   Haus   ist  klein
```

Statistical Machine Translation

Generative Model of Word Alignment
Word Alignments: IBM Model 3
Word Alignments: IBM Model 1
Finding the best alignment: IBM Model 1
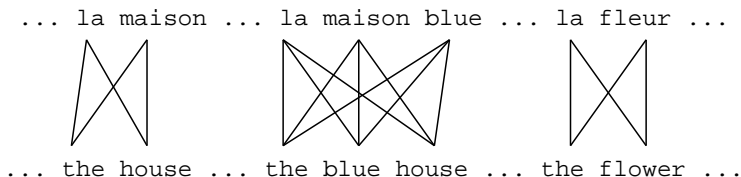Learning Parameters: IBM Model 1

# Learning parameters [from P.Koehn SMT book slides]

- ▶ We would like to estimate the lexical translation probabilities $t(e|f)$ from a parallel corpus
- ▶ ... but we do not have the alignments
- ▶ Chicken and egg problem
    - ▶ if we had the *alignments*,
      $\rightarrow$ we could estimate the *parameters* of our generative model
    - ▶ if we had the *parameters*,
      $\rightarrow$ we could estimate the *alignments*

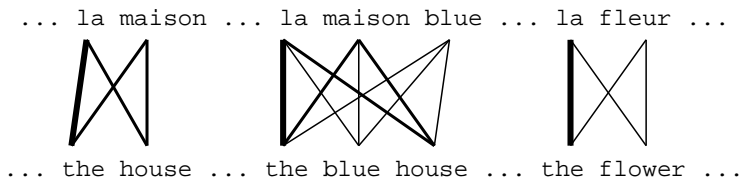# EM Algorithm<sub>[from P.Koehn SMT book slides]</sub>

- ▶ Incomplete data
  - ▶ if we had *complete data*, we could estimate *model*
  - ▶ if we had *model*, we could fill in the *gaps in the data*
- ▶ Expectation Maximization (EM) in a nutshell
  1. initialize model parameters (e.g. uniform)
  2. assign probabilities to the missing data
  3. estimate model parameters from completed data
  4. iterate steps 2–3 until convergence

... la maison ... la maison blue ... la fleur ...

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

- Initial step: all alignments equally likely
- Model learns that, e.g., *la* is often aligned with *the*

... la maison ... la maison blue ... la fleur ...

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

- ▶ After one iteration
- ▶ Alignments, e.g., between *la* and *the* are more likely

```
... la maison ... la maison bleu ... la fleur ...
```

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

- ▶ After another iteration
- ▶ It becomes apparent that alignments, e.g., between *fleur* and *flower* are more likely (pigeon hole principle)

... la maison ... la maison bleu ... la fleur ...

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

- Convergence
- Inherent hidden structure revealed by EM

# EM Algorithm[from P.Koehn SMT book slides]



```
... la maison ... la maison bleu ... la fleur ...
... the house ... the blue house ... the flower ...
```

$$p(la|the) = 0.453$$
$$p(le|the) = 0.334$$
$$p(maison|house) = 0.876$$
$$p(bleu|blue) = 0.563$$
$$...$$

▶ Parameter estimation from the aligned corpus

# IBM Model 1 and the EM Algorithm[from P.Koehn SMT book slides]

- ▶ 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 the EM Algorithm[from P.Koehn SMT book slides]

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

# Word Alignments: IBM Model 1

## Alignment probability

$$
\begin{aligned}
\Pr(\mathbf{a} \mid \mathbf{f}, \mathbf{e}) &= \frac{\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})}{\Pr(\mathbf{f} \mid \mathbf{e})} \\[2mm]
&= \frac{\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})}{\sum_{\mathbf{a}} \Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})} \\[2mm]
&= \frac{\prod_{i=1}^{I} t(f_i \mid e_{a_i})}{\sum_{\mathbf{a}} \prod_{i=1}^{I} t(f_i \mid e_{a_i})}
\end{aligned}
$$

## Computing the denominator

- The denominator above is summing over $J^I$ alignments
- An interlude on how compute the denominator faster ...

# Word Alignments: IBM Model 1

Sum over all alignments

$$\sum_{\mathbf{a}} \Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) \;=\; \sum_{a_1=1}^{J} \sum_{a_2=1}^{J} \cdots \sum_{a_I=1}^{J} \prod_{i=1}^{I} t(f_i \mid e_{a_i})$$

Assume $(f_1, f_2, f_3)$ and $(e_1, e_2)$

$$\sum_{a_1=1}^{2} \sum_{a_2=1}^{2} \sum_{a_3=1}^{2} t(f_1 \mid e_{a_1}) \times t(f_2 \mid e_{a_2}) \times t(f_3 \mid e_{a_3})$$

# Word Alignments: IBM Model 1

Assume $(f_1, f_2, f_3)$ and $(e_1, e_2)$: $I = 3$ and $J = 2$

$$\sum_{a_1=1}^{2} \sum_{a_2=1}^{2} \sum_{a_3=1}^{2} t(f_1 \mid e_{a_1}) \times t(f_2 \mid e_{a_2}) \times t(f_3 \mid e_{a_3})$$

$J^I = 2^3$ terms to be added:

$$
\begin{array}{ccccc}
t(f_1 \mid e_1) & \times & t(f_2 \mid e_1) & \times & t(f_3 \mid e_1) & + \\
t(f_1 \mid e_1) & \times & t(f_2 \mid e_1) & \times & t(f_3 \mid e_2) & + \\
t(f_1 \mid e_1) & \times & t(f_2 \mid e_2) & \times & t(f_3 \mid e_1) & + \\
t(f_1 \mid e_1) & \times & t(f_2 \mid e_2) & \times & t(f_3 \mid e_2) & + \\
t(f_1 \mid e_2) & \times & t(f_2 \mid e_1) & \times & t(f_3 \mid e_1) & + \\
t(f_1 \mid e_2) & \times & t(f_2 \mid e_1) & \times & t(f_3 \mid e_2) & + \\
t(f_1 \mid e_2) & \times & t(f_2 \mid e_2) & \times & t(f_3 \mid e_1) & + \\
t(f_1 \mid e_2) & \times & t(f_2 \mid e_2) & \times & t(f_3 \mid e_2) &
\end{array}
$$

# Word Alignments: IBM Model 1

Factor the terms:

$$
\begin{array}{llll}
(t(f_1 \mid e_1) \times t(f_2 \mid e_1)) & \times & (t(f_3 \mid e_1) + t(f_3 \mid e_2)) & + \\
(t(f_1 \mid e_1) \times t(f_2 \mid e_2)) & \times & (t(f_3 \mid e_1) + t(f_3 \mid e_2)) & + \\
(t(f_1 \mid e_2) \times t(f_2 \mid e_1)) & \times & (t(f_3 \mid e_1) + t(f_3 \mid e_2)) & + \\
(t(f_1 \mid e_2) \times t(f_2 \mid e_2)) & \times & (t(f_3 \mid e_1) + t(f_3 \mid e_2)) &
\end{array}
$$

$$
(t(f_3 \mid e_1) + t(f_3 \mid e_2)) \left(
\begin{array}{llll}
t(f_1 \mid e_1) & \times & t(f_2 \mid e_1) & + \\
t(f_1 \mid e_1) & \times & t(f_2 \mid e_2) & + \\
t(f_1 \mid e_2) & \times & t(f_2 \mid e_1) & + \\
t(f_1 \mid e_2) & \times & t(f_2 \mid e_2) &
\end{array}
\right)
$$

$$
(t(f_3 \mid e_1) + t(f_3 \mid e_2)) \left(
\begin{array}{llll}
t(f_1 \mid e_1) & \times & (t(f_2 \mid e_1) + t(f_2 \mid e_2)) & + \\
t(f_1 \mid e_2) & \times & (t(f_2 \mid e_1) + t(f_2 \mid e_2)) &
\end{array}
\right)
$$

# Word Alignments: IBM Model 1

Assume $(f_1, f_2, f_3)$ and $(e_1, e_2)$: $I = 3$ and $J = 2$

$$\prod_{i=1}^{3} \sum_{a_i=1}^{2} t(f_i \mid e_{a_i})$$

$I \times J = 2 \times 3$ terms to be added:

$$(t(f_1 \mid e_1) \ + \ t(f_1 \mid e_2)) \ \times$$
$$(t(f_2 \mid e_1) \ + \ t(f_2 \mid e_2)) \ \times$$
$$(t(f_3 \mid e_1) \ + \ t(f_3 \mid e_2))$$

# Word Alignments: IBM Model 1

Alignment probability

$$\begin{aligned}
\Pr(\mathbf{a} \mid \mathbf{f}, \mathbf{e}) &= \frac{\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})}{\Pr(\mathbf{f} \mid \mathbf{e})} \\
&= \frac{\prod_{i=1}^{I} t(f_i \mid e_{a_i})}{\sum_{\mathbf{a}} \prod_{i=1}^{I} t(f_i \mid e_{a_i})} \\
&= \frac{\prod_{i=1}^{I} t(f_i \mid e_{a_i})}{\prod_{i=1}^{I} \sum_{j=1}^{J} t(f_i \mid e_j)}
\end{aligned}$$

# Learning Parameters: IBM Model 1



Learning parameters $t(f|e)$ when alignments are known

$$t(das \mid the) = \frac{c(das,the)}{\sum_f c(f,the)} \quad t(house \mid Haus) = \frac{c(Haus,house)}{\sum_f c(f,house)}$$
$$t(ein \mid a) = \frac{c(ein,a)}{\sum_f c(f,a)} \quad t(Buch \mid book) = \frac{c(Buch,book)}{\sum_f c(f,book)}$$

$$t(f|e) = \sum_{s=1}^{N} \sum_{f \to e \in \mathbf{f}^{(s)}, \mathbf{e}^{(s)}} \frac{c(f,e)}{\sum_f c(f,e)}$$

# Learning Parameters: IBM Model 1



Learning parameters $t(f|e)$ when alignments are *unknown*



Also list alignments for *(the book, das Buch)* and *(a book, ein Buch)*

# Learning Parameters: IBM Model 1

Initialize $t^0(f|e)$

| | | | | | | |
|---|---|---|---|---|---|---|
| $t(Haus \mid the)$ | = | 0.25 | | $t(das \mid house)$ | = | 0.5 |
| $t(das \mid the)$ | = | 0.5 | | $t(Haus \mid house)$ | = | 0.5 |
| $t(Buch \mid the)$ | = | 0.25 | | $t(Buch \mid house)$ | = | 0.0 |

Compute posterior for each alignment



$$\Pr(\mathbf{a} \mid \mathbf{f}, \mathbf{e}) = \frac{\Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})}{\Pr(\mathbf{f} \mid \mathbf{e})} = \frac{\prod_{i=1}^{I} t(f_i \mid e_{a_i})}{\prod_{i=1}^{I} \sum_{j=1}^{J} t(f_i \mid e_j)}$$

## Learning Parameters: IBM Model 1

Initialize $t^0(f|e)$

| | | | | | |
|---|---|---|---|---|---|
| $t(Haus \mid the)$ | = | 0.25 | $t(das \mid house)$ | = | 0.5 |
| $t(das \mid the)$ | = | 0.5 | $t(Haus \mid house)$ | = | 0.5 |
| $t(Buch \mid the)$ | = | 0.25 | $t(Buch \mid house)$ | = | 0.0 |

Compute $\Pr(\mathbf{a}, \mathbf{f} \mid \mathbf{e})$ for each alignment



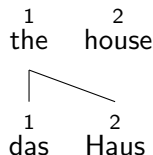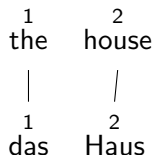| | | | |
|---|---|---|---|
| $0.5 \times 0.25$ | $0.5 \times 0.5$ | $0.25 \times 0.5$ | $0.5 \times 0.5$ |
| 0.125 | 0.25 | 0.125 | 0.25 |

## Learning Parameters: IBM Model 1

Compute $\Pr(\mathbf{a} \mid \mathbf{f}, \mathbf{e}) = \frac{\Pr(\mathbf{a}, \mathbf{f} \mid \mathbf{e})}{\Pr(\mathbf{f} \mid \mathbf{e})}$

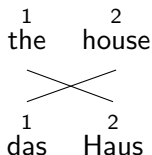$\Pr(\mathbf{f} \mid \mathbf{e}) = 0.125 + 0.25 + 0.125 + 0.25 = 0.75$
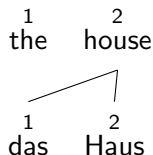


$$\frac{0.125}{0.75} \qquad \frac{0.25}{0.75} \qquad \frac{0.125}{0.75} \qquad \frac{0.25}{0.75}$$
$$0.167 \qquad\qquad 0.334 \qquad\qquad 0.167 \qquad\qquad 0.334$$

Compute fractional counts $c(f, e)$

| | | | | | |
|---|---|---|---|---|---|
| $c(\textit{Haus}, \textit{the})$ | $=$ | $0.125 + 0.125$ | $c(\textit{das}, \textit{house})$ | $=$ | $0.125 + 0.25$ |
| $c(\textit{das}, \textit{the})$ | $=$ | $0.125 + 0.25$ | $c(\textit{Haus}, \textit{house})$ | $=$ | $0.25 + 0.25$ |
| $c(\textit{Buch}, \textit{the})$ | $=$ | $0.0$ | $c(\textit{Buch}, \textit{house})$ | $=$ | $0.0$ |

# Learning Parameters: IBM Model 1



$$\Pr(\mathbf{f} \mid \mathbf{e}) = 0.125 + 0.25 + 0.125 + 0.25 = 0.75$$

Expectation step: expected counts $g(f, e)$

| | | | | | |
|---|---|---|---|---|---|
| $g(das, the)$ | = | $\frac{0.125+0.25}{0.75}$ | $g(das, house)$ | = | $\frac{0.125+0.25}{0.75}$ |
| $g(Haus, the)$ | = | $\frac{0.125+0.125}{0.75}$ | $g(Haus, house)$ | = | $\frac{0.25+0.25}{0.75}$ |
| $g(Buch, the)$ | = | $0.0$ | $g(Buch, house)$ | = | $0.0$ |

Maximization step: get new $t^{(1)}(f \mid e) = \frac{g(f,e)}{\sum_f g(f,e)}$

# Learning Parameters: IBM Model 1

Expectation step: expected counts $g(f, e)$

| | | | | | |
|---|---|---|---|---|---|
| $g(das, the)$ | = | 0.5 | $g(das, house)$ | = | 0.5 |
| $g(Haus, the)$ | = | 0.334 | $g(Haus, house)$ | = | 0.667 |
| $g(Buch, the)$ | = | 0.0 | $g(Buch, house)$ | = | 0.0 |
| **total** | = | 0.834 | **total** | = | 1.167 |

Maximization step: get new $t^{(1)}(f \mid e) = \frac{g(f,e)}{\sum_f g(f,e)}$

| | | | | | |
|---|---|---|---|---|---|
| $t(Haus \mid the)$ | = | 0.4 | $t(das \mid house)$ | = | 0.43 |
| $t(das, \mid the)$ | = | 0.6 | $t(Haus \mid house)$ | = | 0.57 |
| $t(Buch \mid the)$ | = | 0.0 | $t(Buch \mid house)$ | = | 0.0 |

Keep iterating: Compute $t^{(0)}, t^{(1)}, t^{(2)}, \ldots$ until convergence

# Parameter Estimation: IBM Model 1

EM learns the parameters $t(\cdot \mid \cdot)$ that maximizes the log-likelihood of the training data:

$$\arg\max_t L(t) = \arg\max_t \sum_s \log \Pr(\mathbf{f}^{(s)} \mid \mathbf{e}^{(s)}, t)$$

- Start with an initial estimate $t_0$
- Modify it iteratively to get $t_1, t_2, \ldots$
- Re-estimate $t$ from parameters at previous time step $t_{-1}$
- The convergence proof of EM guarantees that $L(t) \geq L(t_{-1})$
- EM converges when $L(t) - L(t_{-1})$ is zero (or almost zero).

## Acknowledgements

Many slides borrowed or inspired from lecture notes by Michael
Collins, Chris Dyer, Kevin Knight, Philipp Koehn, Adam Lopez,
and Luke Zettlemoyer from their NLP course materials.

All mistakes are my own.