# Lexical alignment: IBM models 1 and 2
## MLE via EM for categorical distributions

Miguel Rios

April 3, 2019

# Translation data

Let's assume we are confronted with a new language
and luckily we managed to obtain some sentence-aligned data

| | |
|---|---|
| the black dog | □ ⊛ |
| the nice dog | □ ∪ |
| the black cat | ⊡ ⊛ |
| a dog chasing a cat | ⊡ ◁ □ |

# Translation data

Let's assume we are confronted with a new language
and luckily we managed to obtain some sentence-aligned data

| | |
|---|---|
| the black dog | □ ⊛ |
| the nice dog | □ ∪ |
| the black cat | ⊡ ⊛ |
| a dog chasing a cat | ⊡ ◁ □ |

Is there anything we could say about this language?

# Translation by analogy

| | |
|---:|:---|
| the black dog | $\square$ ⊛ |
| the nice dog | $\square$ ∪ |
| the black cat | $\boxdot$ ⊛ |
| a dog chasing a cat | $\boxdot$ ◁ $\square$ |

A few hypotheses:

# Translation by analogy

| | |
|---:|:---|
| the black dog | $\square$ ⊛ |
| the nice dog | $\square$ ∪ |
| the black cat | $\boxdot$ ⊛ |
| a dog chasing a cat | $\boxdot$ ◁ $\square$ |

A few hypotheses:

- $\square$ ⟺ dog

# Translation by analogy

| | |
|---:|:---|
| the black dog | □ ⊛ |
| the nice dog | □ ∪ |
| the black cat | ⊡ ⊛ |
| a dog chasing a cat | ⊡ ◁ □ |

A few hypotheses:

- □ ⟺ dog
- ⊡ ⟺ cat

## Translation by analogy

| | |
|---:|:---|
| the black dog | □ ⊛ |
| the nice dog | □ ∪ |
| the black cat | ⊡ ⊛ |
| a dog chasing a cat | ⊡ ◁ □ |

A few hypotheses:

- ► □ ⟺ dog
- ► ⊡ ⟺ cat
- ► ⊛ ⟺ black

# Translation by analogy

| | |
|---:|:---|
| the black dog | □ ⊛ |
| the nice dog | □ ∪ |
| the black cat | ⊡ ⊛ |
| a dog chasing a cat | ⊡ ◁ □ |

A few hypotheses:

- ▶ □ ⟺ dog
- ▶ ⊡ ⟺ cat
- ▶ ⊛ ⟺ black
- ▶ nouns seem to preceed adjectives

# Translation by analogy

| | |
|---:|:---|
| the black dog | □ ⊛ |
| the nice dog | □ ∪ |
| the black cat | ⊡ ⊛ |
| a dog chasing a cat | ⊡ ◁ □ |

A few hypotheses:

- ▶ □ ⟺ dog
- ▶ ⊡ ⟺ cat
- ▶ ⊛ ⟺ black
- ▶ nouns seem to preceed adjectives
- ▶ determines are probably not expressed

# Translation by analogy

| | |
|---:|:---|
| the black dog | □ ⊛ |
| the nice dog | □ ∪ |
| the black cat | ⊡ ⊛ |
| a dog chasing a cat | ⊡ ◁ □ |

A few hypotheses:

- ▶ □ ⟺ dog
- ▶ ⊡ ⟺ cat
- ▶ ⊛ ⟺ black
- ▶ nouns seem to preceed adjectives
- ▶ determines are probably not expressed
- ▶ *chasing* may be expressed by ◁
  and perhaps this language is OVS

# Translation by analogy

| | |
|---|---|
| the black dog | □ ⊛ |
| the nice dog | □ ∪ |
| the black cat | ⊡ ⊛ |
| a dog chasing a cat | ⊡ ◁ □ |

A few hypotheses:

- ▶ □ ⟺ dog
- ▶ ⊡ ⟺ cat
- ▶ ⊛ ⟺ black
- ▶ nouns seem to preceed adjectives
- ▶ determines are probably not expressed
- ▶ *chasing* may be expressed by ◁
  and perhaps this language is OVS
- ▶ or perhaps *chasing* is realised by a verb with swapped
  arguments

# Probabilistic lexical alignment models

This lecture is about operationalising this intuition

- through a probabilistic learning algorithm
- for a non-probabilistic approach see for example [**?**]

# Content

# Word-to-word alignments

Imagine you are given a text

| the black dog | el perro negro |
| the nice dog | el perro bonito |
| the black cat | el gato negro |
| a dog chasing a cat | un perro presiguiendo a un gato |

# Word-to-word alignments

Now imagine the French words were replaced by placeholders

| | |
|---|---|
| the black dog | $F_1$ $F_2$ $F_3$ |
| the nice dog | $F_1$ $F_2$ $F_3$ |
| the black cat | $F_1$ $F_2$ $F_3$ |
| a dog chasing a cat | $F_1$ $F_2$ $F_3$ $F_4$ $F_5$ |

# Word-to-word alignments

Now imagine the French words were replaced by placeholders

| | |
|---|---|
| the black dog | $F_1\ F_2\ F_3$ |
| the nice dog | $F_1\ F_2\ F_3$ |
| the black cat | $F_1\ F_2\ F_3$ |
| a dog chasing a cat | $F_1\ F_2\ F_3\ F_4\ F_5$ |

and suppose our task is to have a model explain the original data

# Word-to-word alignments

Now imagine the French words were replaced by placeholders

| | |
|---|---|
| the black dog | $F_1$ $F_2$ $F_3$ |
| the nice dog | $F_1$ $F_2$ $F_3$ |
| the black cat | $F_1$ $F_2$ $F_3$ |
| a dog chasing a cat | $F_1$ $F_2$ $F_3$ $F_4$ $F_5$ |

and suppose our task is to have a model explain the original data
*by generating each French word from exactly one English word*

# Generative story

For each sentence pair independently,

1. observe an English sentence $e_1, \cdots, e_m$
   and a French sentence length $n$
2. for each French word position $j$ from 1 to $n$
   2.1 select an English position $a_j$
   2.2 conditioned on the English word $e_{a_j}$, generate $f_j$

# Generative story

For each sentence pair independently,

1. observe an English sentence $e_1, \cdots, e_m$
   and a French sentence length $n$
2. for each French word position $j$ from 1 to $n$
   2.1 select an English position $a_j$
   2.2 conditioned on the English word $e_{a_j}$, generate $f_j$

We have introduced an alignment
which is not directly visible in the data

# Data augmentation

Observations:

$$\text{the black dog} \mid \text{el perro negro}$$

Imagine data is made of pairs: $(a_j, f_j)$ and $e_{a_j} \rightarrow f_j$

# Data augmentation

Observations:

$$\text{the black dog} \mid \text{el perro negro}$$

Imagine data is made of pairs: $(a_j, f_j)$ and $e_{a_j} \to f_j$

$$\text{the black dog} \mid (A_1, E_{A_1} \to F_1) \ (A_2, E_{A_2} \to F_2) \ (A_3, E_{A_3} \to F_3)$$

# Data augmentation

Observations:

$$\text{the black dog} \mid \text{el perro negro}$$

Imagine data is made of pairs: $(a_j, f_j)$ and $e_{a_j} \to f_j$

$$\text{the black dog} \mid (1, E_{A_1} \to F_1) \ (A_2, E_{A_2} \to F_2) \ (A_3, E_{A_3} \to F_3)$$

# Data augmentation

Observations:

$$\text{the black dog} \mid \text{el perro negro}$$

Imagine data is made of pairs: $(a_j, f_j)$ and $e_{a_j} \to f_j$

$\text{the black dog} \mid (1, \text{the} \to \text{el}) \ (A_2, E_{A_2} \to F_2) \ (A_3, E_{A_3} \to F_3)$

# Data augmentation

Observations:

$$\text{the black dog} \mid \text{el perro negro}$$

Imagine data is made of pairs: $(a_j, f_j)$ and $e_{a_j} \to f_j$

the black dog $\mid$ $(1, \text{the} \to \text{el})$ $(3, E_{A_2} \to F_2)$ $(A_3, E_{A_3} \to F_3)$

# Data augmentation

Observations:

$$\text{the black dog} \mid \text{el perro negro}$$

Imagine data is made of pairs: $(a_j, f_j)$ and $e_{a_j} \rightarrow f_j$

the black dog $\mid$ $(1, \text{the} \rightarrow \text{el})$ $(3, \text{dog} \rightarrow \text{perro})$ $(A_3, E_{A_3} \rightarrow F_3)$

# Data augmentation

Observations:

$$\text{the black dog} \mid \text{el perro negro}$$

Imagine data is made of pairs: $(a_j, f_j)$ and $e_{a_j} \to f_j$

the black dog $\mid$ $(1, \text{the} \to \text{el})$ $(3, \text{dog} \to \text{perro})$ $(2, E_{A_3} \to F_3)$

# Data augmentation

Observations:

$$\text{the black dog} \mid \text{el perro negro}$$

Imagine data is made of pairs: $(a_j, f_j)$ and $e_{a_j} \to f_j$

the black dog $\mid$ $(1, \text{the} \to \text{el})$ $(3, \text{dog} \to \text{perro})$ $(2, \text{black} \to \text{negro})$

# Data augmentation

Observations:

the black dog │ el perro negro

Imagine data is made of pairs: $(a_j, f_j)$ and $e_{a_j} \to f_j$

the black dog │ $(1, \text{the} \to \text{el})$ $(3, \text{dog} \to \text{perro})$ $(2, \text{black} \to \text{negro})$

the black dog │ $(A_1, \text{the} \to \text{el})$ $(A_1, \text{the} \to \text{perro})$ $(A_1, \text{the} \to \text{negro})$

# Data augmentation

Observations:

$$\text{the black dog} \mid \text{el perro negro}$$

Imagine data is made of pairs: $(a_j, f_j)$ and $e_{a_j} \to f_j$

$$\text{the black dog} \mid (1, \text{the} \to \text{el})\ (3, \text{dog} \to \text{perro})\ (2, \text{black} \to \text{negro})$$

$$\text{the black dog} \mid (A_1, \text{the} \to \text{el})\ (A_1, \text{the} \to \text{perro})\ (A_1, \text{the} \to \text{negro})$$

$$\text{the black dog} \mid (a_1, e_{a_1} \to f_1)\ (a_2, e_{a_2} \to f_2)\ (a_3, e_{a_3} \to f_3)$$

# Content

# Mixture models: generative story



- $c$ mixture components
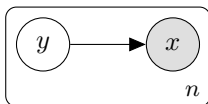- each defines a distribution over the same data space $\mathcal{X}$
- plus a distribution over components themselves
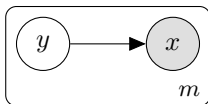
# Mixture models: generative story



- ▶ $c$ mixture components
- ▶ each defines a distribution over the same data space $\mathcal{X}$
- ▶ plus a distribution over components themselves

Generative story

1. select a mixture component $y \sim p(y)$
2. generate an observation from it $x \sim p(x|y)$

# Mixture models: likelihood



Incomplete-data likelihood

$$p(x_1^m) = \prod_{i=1}^{m} p(x_i) \tag{1}$$

$$= \prod_{i=1}^{m} \sum_{y=1}^{c} \underbrace{p(x_i, y)}_{\text{complete-data likelihood}} \tag{2}$$

$$= \prod_{i=1}^{m} \sum_{y=1}^{c} p(z) p(x_i | y) \tag{3}$$

# Interpretation

Missing data

- Let $y$ take one of $c$ mixture components
- Assume data consists of pairs $(x, y)$
- $x$ is always observed
- $y$ is always missing

## Interpretation

Missing data

- Let $y$ take one of $c$ mixture components
- Assume data consists of pairs $(x, y)$
- $x$ is always observed
- $y$ is always missing

Inference: posterior distribution over possible $y$ for each $x$

$$p(y|x) = \frac{p(y, x)}{\sum_{y'=1}^{c} p(y', x)} \tag{4}$$

$$= \frac{p(y)p(x|y)}{\sum_{y'=1}^{c} p(y')p(x|y')} \tag{5}$$

# Non-identifiability

Different parameter settings, same distribution

Suppose $\mathcal{X} = \{a, b\}$ and $c = 2$
and let $p(y = 1) = p(y = 2) = 0.5$

| $y$ | $x = a$ | $x = b$ |
|---|---|---|
| 1 | 0.2 | 0.8 |
| 2 | 0.7 | 0.3 |
| $p(x)$ | 0.45 | 0.55 |

| $y$ | $x = a$ | $x = b$ |
|---|---|---|
| 1 | 0.7 | 0.3 |
| 2 | 0.2 | 0.8 |
| $p(x)$ | 0.45 | 0.55 |

# Non-identifiability

Different parameter settings, same distribution

Suppose $\mathcal{X} = \{a, b\}$ and $c = 2$
and let $p(y = 1) = p(y = 2) = 0.5$

| $y$ | $x = a$ | $x = b$ |
|---|---|---|
| 1 | 0.2 | 0.8 |
| 2 | 0.7 | 0.3 |
| $p(x)$ | 0.45 | 0.55 |

| $y$ | $x = a$ | $x = b$ |
|---|---|---|
| 1 | 0.7 | 0.3 |
| 2 | 0.2 | 0.8 |
| $p(x)$ | 0.45 | 0.55 |

Problem for parameter estimation by hillclimbing

# Maximum likelihood estimation

Suppose a dataset $\mathcal{D} = \{x^{(1)}, x^{(2)}, \cdots, x^{(m)}\}$

# Maximum likelihood estimation

Suppose a dataset $\mathcal{D} = \{x^{(1)}, x^{(2)}, \cdots, x^{(m)}\}$

Suppose $p(x)$ is one of a parametric family with parameters $\theta$

# Maximum likelihood estimation

Suppose a dataset $\mathcal{D} = \{x^{(1)}, x^{(2)}, \cdots, x^{(m)}\}$

Suppose $p(x)$ is one of a parametric family with parameters $\theta$

Likelihood of iid observations

$$p(\mathcal{D}) = \prod_{i=1}^{m} p_\theta(x^{(i)})$$

# Maximum likelihood estimation

Suppose a dataset $\mathcal{D} = \{x^{(1)}, x^{(2)}, \cdots, x^{(m)}\}$

Suppose $p(x)$ is one of a parametric family with parameters $\theta$

Likelihood of iid observations

$$p(\mathcal{D}) = \prod_{i=1}^{m} p_\theta(x^{(i)})$$

the score function is

$$l(\theta) = \sum_{i=1}^{m} \log p_\theta(x^{(i)})$$

# Maximum likelihood estimation

Suppose a dataset $\mathcal{D} = \{x^{(1)}, x^{(2)}, \cdots, x^{(m)}\}$
Suppose $p(x)$ is one of a parametric family with parameters $\theta$
Likelihood of iid observations

$$p(\mathcal{D}) = \prod_{i=1}^{m} p_\theta(x^{(i)})$$

the score function is

$$l(\theta) = \sum_{i=1}^{m} \log p_\theta(x^{(i)})$$

then we choose

$$\theta^\star = \arg\max_\theta l(\theta)$$

# MLE for categorical: estimation from fully observed data

Suppose we have **complete data**

- $\mathcal{D}_{\mathsf{complete}} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

# MLE for categorical: estimation from fully observed data

Suppose we have **complete data**

- $\mathcal{D}_{\text{complete}} = \{(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})\}$

Then, for a **categorical distribution**

$$p(x|y) = \theta_{y,x}$$

and $n(y, x|\mathcal{D}_{\text{complete}}) =$ *count of* $(y, x)$ *in* $\mathcal{D}_{\text{complete}}$

MLE solution:

$$\theta_{y,x} = \frac{n(y, x|\mathcal{D}_{\text{complete}})}{\sum_{x'} n(y, x'|\mathcal{D}_{\text{complete}})}$$

# MLE for categorical: estimation from incomplete data

**Expectation-Maximisation algorithm** [**?**]

E-step:

- for every observation $x$, imagine that every possible latent assignment $y$ happened with probability $p_\theta(y|x)$

$$\mathcal{D}_{\text{completed}} = \{(x, y = 1), \dots, (x, y = c) : x \in \mathcal{D}\}$$

# MLE for categorical: estimation from incomplete data

**Expectation-Maximisation algorithm**                                   [?]

M-step:

- reestimate $\theta$ as to climb the likelihood surface
- for categorical distributions $p(x|y) = \theta_{y,x}$
  $y$ and $x$ are categorical
  $0 \le \theta_{y,x} \le 1$   and   $\sum_{x \in X} \theta_{y,x} = 1$

$$
\begin{align}
\theta_{y,x} &= \frac{\mathbb{E}[n(y \to x|\mathcal{D}_{\text{completed}})]}{\sum_{x'} \mathbb{E}[n(y \to x'|\mathcal{D}_{\text{completed}})]} \tag{6} \\
&= \frac{\sum_{i=1}^{m} \sum_{y'} p(y'|x^{(i)}) \mathbb{1}_y(y') \mathbb{1}_x(x^{(i)})}{\sum_{i=1}^{m} \sum_{x'} \sum_{y'} p(y'|x^{(i)}) \mathbb{1}_y(y') \mathbb{1}_{x'}(x^{(i)})} \tag{7} \\
&= \frac{\sum_{i=1}^{m} p(y|x^{(i)}) \mathbb{1}_x(x^{(i)})}{\sum_{i=1}^{m} \sum_{x'} p(y|x^{(i)}) \mathbb{1}_{x'}(x^{(i)})} \tag{8}
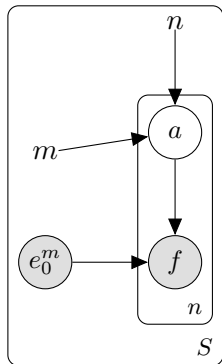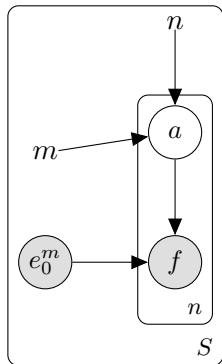\end{align}
$$

# Content

# IBM1: a constrained mixture model



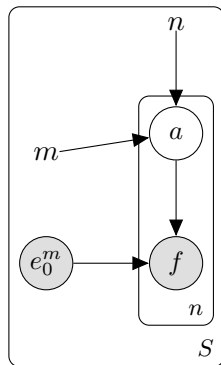Constrained mixture model

# IBM1: a constrained mixture model



Constrained mixture model

- mixture components are English words

# IBM1: a constrained mixture model



Constrained mixture model

- mixture components are English words
- but only English words that appear in the English sentence can be assigned
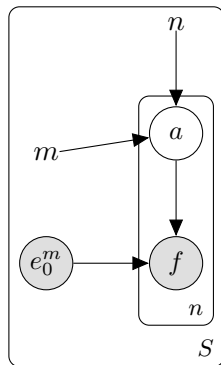
# IBM1: a constrained mixture model



Constrained mixture model

- ▶ mixture components are English words
- ▶ but only English words that appear in the English sentence can be assigned
- ▶ $a_j$ acts as an indicator for the mixture component that generates French word $f_j$
- ▶ $e_0$ is occupied by a special NULL component

## Parameterisation

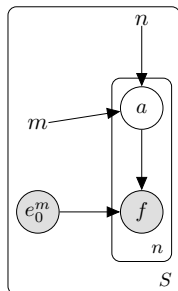Alignment distribution: uniform

$$p(a|m, n) = \frac{1}{m+1} \tag{9}$$

Lexical distribution: categorical

$$p(f|e) = \mathrm{Cat}(f|\theta_e) \tag{10}$$

- ▶ where $\theta_e \in \mathbb{R}^{v_F}$
- ▶ $0 \leq \theta_{e,f} \leq 1$
- ▶ $\sum_f \theta_{e,f} = 1$

# IBM1: incomplete-data likelihood

Incomplete-data likelihood



$$p(f_1^n | e_0^m) = \sum_{a_1=0}^{m} \cdots \sum_{a_n=0}^{m} p(f_1^n, a_1^n | e_{a_j}) \tag{11}$$

$$= \sum_{a_1=0}^{m} \cdots \sum_{a_n=0}^{m} \prod_{j=1}^{n} p(a_j | m, n) p(f_j | e_{a_j}) \tag{12}$$

$$= \prod_{j=1}^{n} \sum_{a_j=0}^{m} p(a_j | m, n) p(f_j | e_{a_j}) \tag{13}$$

# IBM1: posterior

Posterior

$$p(a_1^n | f_1^n, e_0^m) = \frac{p(f_1^n, a_1^n | e_0^m)}{p(f_1^n | e_0^m)} \qquad (14)$$

Factorised

$$p(a_j | f_1^n, e_0^m) = \frac{p(a_j | m, n) p(f_j | e_{a_j})}{\sum_{i=0}^{m} p(i | m, n) p(f_j | e_i)} \qquad (15)$$

## MLE via EM

E-step:

$$\mathbb{E}[n(\mathsf{e} \to \mathsf{f}|a_1^n)] = \sum_{a_1=0}^{m} \cdots \sum_{a_n=0}^{m} p(a_1^n|f_1^n, e_0^m) n(\mathsf{e} \to \mathsf{f}|A_1^n) \tag{16}$$

$$= \sum_{a_1=0}^{m} \cdots \sum_{a_n=0}^{m} \prod_{j=1}^{n} p(a_j|f_1^n, e_0^m) \mathbb{1}_{\mathsf{e}}(e_{a_j}) \mathbb{1}_{\mathsf{f}}(f_j) \tag{17}$$

$$= \prod_{j=1}^{n} \sum_{i=0}^{m} p(a_j = i|f_1^n, e_0^m) \mathbb{1}_{\mathsf{e}}(e_i) \mathbb{1}_{\mathsf{f}}(f_j) \tag{18}$$

M-step:

$$\theta_{e,f} = \frac{\mathbb{E}[n(e \to f|a_1^n)]}{\sum_{f'} \mathbb{E}[n(e \to f'|a_1^n)]} \tag{19}$$

# EM algorithm

Repeat until convergence to a local optimum
1. For each sentence pair
    1.1 compute posterior per alignment link
    1.2 accumulate fractional counts
2. Normalise counts for each English word

# Content

# Alignment distribution

Positional distribution

$p(a_j|m, n) = \mathrm{Cat}(a|\lambda_{j,m,n})$

- one distribution for each tuple $(j, m, n)$
- support must include length of longest English sentence
- extremely over-parameterised!

# Alignment distribution

Positional distribution

$p(a_j|m,n) = \mathrm{Cat}(a|\lambda_{j,m,n})$

- one distribution for each tuple $(j, m, n)$
- support must include length of longest English sentence
- extremely over-parameterised!

Jump distribution                                                    [**?**]

- define a jump function $\delta(a_j, j, m, n) = a_j - \lfloor j\frac{m}{n} \rfloor$
- $p(a_j|m,n) = \mathrm{Cat}(\Delta|\lambda)$
- $\Delta$ takes values from $-$longest to $+$longest

# Content

# Note on terminology: source/target vs French/English

From an alignment model perspective all that matters is

- ▶ we condition on one language and generate the other
- ▶ in IBM models terminology, we condition on *English* and generate *French*

From a noisy channel perspective, where we want to translate a *source* sentence $f_1^n$ into some *target* sentence $e_1^m$

- ▶ Bayes rule decomposes $p(e_1^m|f_1^n) \propto p(f_1^n|e_1^m)p(e_1^m)$
- ▶ train $p(e_1^m)$ and $p(f_1^n|e_1^m)$ independently
- ▶ **language model:** $p(e_1^m)$
- ▶ **alignment model:** $p(f_1^n|e_1^m)$
- ▶ note that the alignment model conditions on the target sentence (English) and generates the source sentence (French)

# Limitations of IBM1-2

- ► too strong independence assumptions
- ► categorical parameterisation suffers from data sparsity
- ► EM suffers from local optima

# Extensions

Fertility, distortion, and concepts [**?**]

Dirichlet priors and posterior inference [**?**]

- ► + no NULL words [**?**]
- ► + HMM and efficient sampler [**?**]

Log-linear distortion parameters and variational Bayes
[**?**]

First-order dependency (HMM) [**?**]

- ► E-step requires dynamic programming
  [**?**]

# References I