

6.864 (Fall 2007)
Machine Translation Part II

1

Roadmap for the Next Few Lectures

- Lecture 1 (today): IBM Models 1 and 2
- Lecture 2: *phrase-based* models
- Lecture 3: Syntax in statistical machine translation

3

Recap: The Noisy Channel Model

- Goal: translation system from French to English
- Have a model $P(\mathbf{e} | \mathbf{f})$ which estimates conditional probability of any English sentence \mathbf{e} given the French sentence \mathbf{f} . Use the training corpus to set the parameters.
- A Noisy Channel Model has two components:

$P(\mathbf{e})$ **the language model**

$P(\mathbf{f} | \mathbf{e})$ **the translation model**

- Giving:

$$P(\mathbf{e} | \mathbf{f}) = \frac{P(\mathbf{e}, \mathbf{f})}{P(\mathbf{f})} = \frac{P(\mathbf{e})P(\mathbf{f} | \mathbf{e})}{\sum_{\mathbf{e}} P(\mathbf{e})P(\mathbf{f} | \mathbf{e})}$$

and

$$\operatorname{argmax}_{\mathbf{e}} P(\mathbf{e} | \mathbf{f}) = \operatorname{argmax}_{\mathbf{e}} P(\mathbf{e})P(\mathbf{f} | \mathbf{e})$$

2

Overview

- IBM Model 1
- IBM Model 2
- EM Training of Models 1 and 2
- Some examples of training Models 1 and 2
- Decoding

4

IBM Model 1: Alignments

- How do we model $P(\mathbf{f} | \mathbf{e})$?
- English sentence \mathbf{e} has l words $e_1 \dots e_l$,
French sentence \mathbf{f} has m words $f_1 \dots f_m$.
- An **alignment** \mathbf{a} identifies which English word each French word originated from
- Formally, an **alignment** \mathbf{a} is $\{a_1, \dots, a_m\}$, where each $a_j \in \{0 \dots l\}$.
- There are $(l + 1)^m$ possible alignments.

5

Alignments in the IBM Models

- We'll define models for $P(\mathbf{a} | \mathbf{e})$ and $P(\mathbf{f} | \mathbf{a}, \mathbf{e})$, giving

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

- Also,

$$P(\mathbf{f} | \mathbf{e}) = \sum_{\mathbf{a} \in \mathcal{A}} P(\mathbf{a} | \mathbf{e})P(\mathbf{f} | \mathbf{a}, \mathbf{e})$$

where \mathcal{A} is the set of all possible alignments

7

IBM Model 1: Alignments

- e.g., $l = 6, m = 7$

\mathbf{e} = And the program has been implemented

\mathbf{f} = Le programme a ete mis en application

- One alignment is

$$\{2, 3, 4, 5, 6, 6, 6\}$$

- Another (bad!) alignment is

$$\{1, 1, 1, 1, 1, 1, 1\}$$

6

A By-Product: Most Likely Alignments

- Once we have a model $P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = P(\mathbf{a} | \mathbf{e})P(\mathbf{f} | \mathbf{a}, \mathbf{e})$ we can also calculate

$$P(\mathbf{a} | \mathbf{f}, \mathbf{e}) = \frac{P(\mathbf{f}, \mathbf{a} | \mathbf{e})}{\sum_{\mathbf{a} \in \mathcal{A}} P(\mathbf{f}, \mathbf{a} | \mathbf{e})}$$

for any alignment \mathbf{a}

- For a given \mathbf{f}, \mathbf{e} pair, we can also compute the most likely alignment,

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

- Nowadays, the original IBM models are rarely (if ever) used for translation, but they **are** used for recovering alignments

8

An Example Alignment

French:

le conseil a rendu son avis , et nous devons à présent adopter un nouvel avis sur la base de la première position .

English:

the council has stated its position , and now , on the basis of the first position , we again have to give our opinion .

Alignment:

the/**le** council/**conseil** has/**à** stated/**rendu** its/**son** position/**avis** ./,
and/**et** now/**présent** ./,NULL on/**sur** the/**le** basis/**base** of/**de** the/**la**
first/**première** position/**position** ./,NULL we/**nous** again/**NULL**
have/**devons** to/**a** give/**adopter** our/**nouvel** opinion/**avis** ./.

9

IBM Model 1: Translation Probabilities

- Next step: come up with an estimate for

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

- In model 1, this is:

$$P(\mathbf{f} \mid \mathbf{a}, \mathbf{e}) = \prod_{j=1}^m P(f_j \mid e_{a_j})$$

11

IBM Model 1: Alignments

- In IBM model 1 all alignments \mathbf{a} are equally likely:

$$P(\mathbf{a} \mid \mathbf{e}) = C \times \frac{1}{(l+1)^m}$$

where $C = \text{prob}(\text{length}(\mathbf{f}) = m)$ is a constant.

- This is a **major** simplifying assumption, but it gets things started...

10

- e.g., $l = 6, m = 7$

\mathbf{e} = And the program has been implemented

\mathbf{f} = Le programme a ete mis en application

- $\mathbf{a} = \{2, 3, 4, 5, 6, 6, 6\}$

$$\begin{aligned} P(\mathbf{f} \mid \mathbf{a}, \mathbf{e}) &= P(\text{Le} \mid \text{the}) \times \\ &P(\text{programme} \mid \text{program}) \times \\ &P(\text{a} \mid \text{has}) \times \\ &P(\text{ete} \mid \text{been}) \times \\ &P(\text{mis} \mid \text{implemented}) \times \\ &P(\text{en} \mid \text{implemented}) \times \\ &P(\text{application} \mid \text{implemented}) \end{aligned}$$

12

IBM Model 1: The Generative Process

To generate a French string f from an English string e :

- **Step 1:** Pick the length of f (all lengths equally probable, probability C)
- **Step 2:** Pick an alignment a with probability $\frac{1}{(l+1)^m}$
- **Step 3:** Pick the French words with probability

$$P(f | a, e) = \prod_{j=1}^m P(f_j | e_{a_j})$$

The final result:

$$P(f, a | e) = P(a | e) \times P(f | a, e) = \frac{C}{(l+1)^m} \prod_{j=1}^m P(f_j | e_{a_j})$$

13

A Hidden Variable Problem

- Training data is a set of (f_i, e_i) pairs, likelihood is

$$\sum_i \log P(f_i | e_i) = \sum_i \log \sum_{a \in \mathcal{A}} P(a | e_i) P(f_i | a, e_i)$$

where \mathcal{A} is the set of all possible alignments.

- We need to maximize this function w.r.t. the translation parameters $P(f_j | e_{a_j})$.
- EM can be used for this problem: initialize translation parameters randomly, and at each iteration choose

$$\Theta_t = \operatorname{argmax}_{\Theta} \sum_i \sum_{a \in \mathcal{A}} P(a | e_i, f_i, \Theta^{t-1}) \log P(f_i | a, e_i, \Theta)$$

where Θ^t are the parameter values at the t 'th iteration.

15

A Hidden Variable Problem

- We have:

$$P(f, a | e) = \frac{C}{(l+1)^m} \prod_{j=1}^m P(f_j | e_{a_j})$$

- And:

$$P(f | e) = \sum_{a \in \mathcal{A}} \frac{C}{(l+1)^m} \prod_{j=1}^m P(f_j | e_{a_j})$$

where \mathcal{A} is the set of all possible alignments.

14

An Example

- I have the following training examples

the dog \Rightarrow le chien
the cat \Rightarrow le chat

- Need to find estimates for:

$$\begin{array}{lll} P(le | the) & P(chien | the) & P(chat | the) \\ P(le | dog) & P(chien | dog) & P(chat | dog) \\ P(le | cat) & P(chien | cat) & P(chat | cat) \end{array}$$

- As a result, each (e_i, f_i) pair will have a most likely alignment.

16

An Example Lexical Entry

English	French	Probability
position	position	0.756715
position	situation	0.0547918
position	mesure	0.0281663
position	vue	0.0169303
position	point	0.0124795
position	attitude	0.0108907

... de la **situation** au niveau des négociations de l'OMPI ...
... of the current **position** in the wipo negotiations ...

nous ne sommes pas en **mesure** de décider, ...
we are not in a **position** to decide, ...

... le **point de vue** de la commission face à ce problème complexe .
... the commission's **position** on this complex problem .

... cette **attitude** laxiste et irresponsable .
... this irresponsibly lax **position** .

17

IBM Model 2

- Only difference: we now introduce **alignment** or **distortion** parameters

$D(i | j, l, m)$ = Probability that j 'th French word is connected to i 'th English word, given sentence lengths of e and f are l and m respectively

- Define

$$P(\mathbf{a} | \mathbf{e}, l, m) = \prod_{j=1}^m D(a_j | j, l, m)$$

where $\mathbf{a} = \{a_1, \dots, a_m\}$

- Gives

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}, l, m) = \prod_{j=1}^m D(a_j | j, l, m) T(f_j | e_{a_j})$$

19

Overview

- IBM Model 1
- **IBM Model 2**
- EM Training of Models 1 and 2
- Some examples of training Models 1 and 2
- Decoding

18

- Note: Model 1 is a special case of Model 2, where $D(i | j, l, m) = \frac{1}{l+1}$ for all i, j .

20

An Example

$l = 6$
 $m = 7$
 $e =$ And the program has been implemented
 $f =$ Le programme a ete mis en application
 $a = \{2, 3, 4, 5, 6, 6, 6\}$

$$P(\mathbf{a} | \mathbf{e}, l, m) = \mathbf{D}(2 | 1, 6, 7) \times \\ \mathbf{D}(3 | 2, 6, 7) \times \\ \mathbf{D}(4 | 3, 6, 7) \times \\ \mathbf{D}(5 | 4, 6, 7) \times \\ \mathbf{D}(6 | 5, 6, 7) \times \\ \mathbf{D}(6 | 6, 6, 7) \times \\ \mathbf{D}(6 | 7, 6, 7)$$

21

$$P(\mathbf{f} | \mathbf{a}, \mathbf{e}) = \mathbf{T}(Le | the) \times \\ \mathbf{T}(programme | program) \times \\ \mathbf{T}(a | has) \times \\ \mathbf{T}(ete | been) \times \\ \mathbf{T}(mis | implemented) \times \\ \mathbf{T}(en | implemented) \times \\ \mathbf{T}(application | implemented)$$

22

IBM Model 2: The Generative Process

To generate a French string \mathbf{f} from an English string \mathbf{e} :

- **Step 1:** Pick the length of \mathbf{f} (all lengths equally probable, probability C)
- **Step 2:** Pick an alignment $\mathbf{a} = \{a_1, a_2 \dots a_m\}$ with probability

$$\prod_{j=1}^m \mathbf{D}(a_j | j, l, m)$$

- **Step 3:** Pick the French words with probability

$$P(\mathbf{f} | \mathbf{a}, \mathbf{e}) = \prod_{j=1}^m \mathbf{T}(f_j | e_{a_j})$$

The final result:

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = P(\mathbf{a} | \mathbf{e})P(\mathbf{f} | \mathbf{a}, \mathbf{e}) = C \prod_{j=1}^m \mathbf{D}(a_j | j, l, m) \mathbf{T}(f_j | e_{a_j})$$

23

Overview

- IBM Model 1
- IBM Model 2
- **EM Training of Models 1 and 2**
- Some examples of training Models 1 and 2
- Decoding

24

A Hidden Variable Problem

- We have:

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = C \prod_{j=1}^m \mathbf{D}(a_j | j, l, m) \mathbf{T}(f_j | e_{a_j})$$

- And:

$$P(\mathbf{f} | \mathbf{e}) = \sum_{\mathbf{a} \in \mathcal{A}} C \prod_{j=1}^m \mathbf{D}(a_j | j, l, m) \mathbf{T}(f_j | e_{a_j})$$

where \mathcal{A} is the set of all possible alignments.

25

Model 2 as a Product of Multinomials

- The model can be written in the form

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = \prod_r \Theta_r^{Count(\mathbf{f}, \mathbf{a}, \mathbf{e}, r)}$$

where the parameters Θ_r correspond to the $\mathbf{T}(f|e)$ and $\mathbf{D}(i|j, l, m)$ parameters

- To apply EM, we need to calculate expected counts

$$\overline{Count}(r) = \sum_k \sum_{\mathbf{a}} P(\mathbf{a} | \mathbf{e}_k, \mathbf{f}_k, \bar{\Theta}) Count(\mathbf{f}_k, \mathbf{a}, \mathbf{e}_k, r)$$

27

A Hidden Variable Problem

- Training data is a set of $(\mathbf{f}_k, \mathbf{e}_k)$ pairs, likelihood is

$$\sum_k \log P(\mathbf{f}_k | \mathbf{e}_k) = \sum_k \log \sum_{\mathbf{a} \in \mathcal{A}} P(\mathbf{a} | \mathbf{e}_k) P(\mathbf{f}_k | \mathbf{a}, \mathbf{e}_k)$$

where \mathcal{A} is the set of all possible alignments.

- We need to maximize this function w.r.t. the translation parameters, and the alignment probabilities
- EM can be used for this problem: initialize parameters randomly, and at each iteration choose

$$\Theta_t = \operatorname{argmax}_{\Theta} \sum_k \sum_{\mathbf{a} \in \mathcal{A}} P(\mathbf{a} | \mathbf{e}_k, \mathbf{f}_k, \Theta^{t-1}) \log P(\mathbf{f}_k, \mathbf{a} | \mathbf{e}_k, \Theta)$$

where Θ^t are the parameter values at the t 'th iteration.

26

A Crucial Step in the EM Algorithm

- Say we have the following (\mathbf{e}, \mathbf{f}) pair:

\mathbf{e} = And the program has been implemented

\mathbf{f} = Le programme a ete mis en application

- Given that \mathbf{f} was generated according to Model 2, what is the probability that $a_1 = 2$? **Formally:**

$$Prob(a_1 = 2 | \mathbf{f}, \mathbf{e}) = \sum_{\mathbf{a}: a_1=2} P(\mathbf{a} | \mathbf{f}, \mathbf{e}, \bar{\Theta})$$

28

Calculating Expected Translation Counts

- One example:

$$\overline{\text{Count}}(\mathbf{T}(le|the)) = \sum_{(i,j,k) \in \mathcal{S}} P(a_j = i | \mathbf{e}_k, \mathbf{f}_k, \bar{\Theta})$$

where \mathcal{S} is the set of all (i, j, k) triples such that $e_{k,i} = the$ and $f_{k,j} = le$

29

Models 1 and 2 Have a Simple Structure

- We have $\mathbf{f} = \{f_1 \dots f_m\}$, $\mathbf{a} = \{a_1 \dots a_m\}$, and

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}, l, m) = \prod_{j=1}^m P(a_j, f_j | \mathbf{e}, l, m)$$

where

$$P(a_j, f_j | \mathbf{e}, l, m) = \mathbf{D}(a_j | j, l, m) \mathbf{T}(f_j | e_{a_j})$$

- **We can think of the m (f_j, a_j) pairs as being generated independently**

31

Calculating Expected Distortion Counts

- One example:

$$\overline{\text{Count}}(\mathbf{D}(i = 5 | j = 6, l = 10, m = 11)) = \sum_{k \in \mathcal{S}} P(a_6 = 5 | \mathbf{e}_k, \mathbf{f}_k, \bar{\Theta})$$

where \mathcal{S} is the set of all values of k such that $\text{length}(\mathbf{e}_k) = 10$ and $\text{length}(\mathbf{f}_k) = 11$

30

The Answer

$$\begin{aligned} \text{Prob}(a_1 = 2 | \mathbf{f}, \mathbf{e}) &= \sum_{\mathbf{a}: a_1=2} P(\mathbf{a} | \mathbf{f}, \mathbf{e}, l, m) \\ &= \frac{\mathbf{D}(a_1 = 2 | j = 1, l = 6, m = 7) \mathbf{T}(le | the)}{\sum_{i=0}^l \mathbf{D}(a_1 = i | j = 1, l = 6, m = 7) \mathbf{T}(le | e_i)} \end{aligned}$$

Follows directly because the (a_j, f_j) pairs are independent:

$$P(a_1 = 2 | \mathbf{f}, \mathbf{e}, l, m) = \frac{P(a_1 = 2, f_1 = le | f_2 \dots f_m, \mathbf{e}, l, m)}{P(f_1 = le | f_2 \dots f_m, \mathbf{e}, l, m)} \quad (1)$$

$$= \frac{P(a_1 = 2, f_1 = le | \mathbf{e}, l, m)}{P(f_1 = le | \mathbf{e}, l, m)} \quad (2)$$

$$= \frac{P(a_1 = 2, f_1 = le | \mathbf{e}, l, m)}{\sum_i P(a_1 = i, f_1 = le | \mathbf{e}, l, m)}$$

where (2) follows from (1) because $P(\mathbf{f}, \mathbf{a} | \mathbf{e}, l, m) = \prod_{j=1}^m P(a_j, f_j | \mathbf{e}, l, m)$

32

A General Result

$$\begin{aligned} \text{Prob}(a_j = i \mid \mathbf{f}, \mathbf{e}) &= \sum_{\mathbf{a}: a_j = i} P(\mathbf{a} \mid \mathbf{f}, \mathbf{e}, l, m) \\ &= \frac{\mathbf{D}(a_j = i \mid j, l, m) \mathbf{T}(f_j \mid e_i)}{\sum_{i'=0}^l \mathbf{D}(a_j = i' \mid j, l, m) \mathbf{T}(f_j \mid e_{i'})} \end{aligned}$$

33

Alignment Probabilities have a Simple Solution!

- e.g., Say we have $l = 6, m = 7$,

\mathbf{e} = And the program has been implemented

\mathbf{f} = Le programme a ete mis en application

- Probability of “mis” being connected to “the”:

$$P(a_5 = 2 \mid \mathbf{f}, \mathbf{e}) = \frac{\mathbf{D}(a_5 = 2 \mid j = 5, l = 6, m = 7) \mathbf{T}(\text{mis} \mid \text{the})}{Z}$$

where

$$\begin{aligned} Z = & \mathbf{D}(a_5 = 0 \mid j = 5, l = 6, m = 7) \mathbf{T}(\text{mis} \mid \text{NULL}) \\ & + \mathbf{D}(a_5 = 1 \mid j = 5, l = 6, m = 7) \mathbf{T}(\text{mis} \mid \text{And}) \\ & + \mathbf{D}(a_5 = 2 \mid j = 5, l = 6, m = 7) \mathbf{T}(\text{mis} \mid \text{the}) \\ & + \mathbf{D}(a_5 = 3 \mid j = 5, l = 6, m = 7) \mathbf{T}(\text{mis} \mid \text{program}) \\ & + \dots \end{aligned}$$

34

The EM Algorithm for Model 2

- Define

$\mathbf{e}[k]$ for $k = 1 \dots n$ is the k 'th English sentence

$\mathbf{f}[k]$ for $k = 1 \dots n$ is the k 'th French sentence

$l[k]$ is the length of $\mathbf{e}[k]$

$m[k]$ is the length of $\mathbf{f}[k]$

$\mathbf{e}[k, i]$ is the i 'th word in $\mathbf{e}[k]$

$\mathbf{f}[k, j]$ is the j 'th word in $\mathbf{f}[k]$

- Current parameters Θ^{t-1} are

$\mathbf{T}(f \mid e)$ for all $f \in \mathcal{F}, e \in \mathcal{E}$

$\mathbf{D}(i \mid j, l, m)$

- We'll see how the EM algorithm re-estimates the \mathbf{T} and \mathbf{D} parameters

35

Step 1: Calculate the Alignment Probabilities

- Calculate an array of alignment probabilities (for $(k = 1 \dots n), (j = 1 \dots m[k]), (i = 0 \dots l[k])$):

$$\begin{aligned} a[i, j, k] &= P(a_j = i \mid \mathbf{e}[k], \mathbf{f}[k], \Theta^{t-1}) \\ &= \frac{\mathbf{D}(a_j = i \mid j, l, m) \mathbf{T}(f_j \mid e_i)}{\sum_{i'=0}^l \mathbf{D}(a_j = i' \mid j, l, m) \mathbf{T}(f_j \mid e_{i'})} \end{aligned}$$

where $e_i = \mathbf{e}[k, i]$, $f_j = \mathbf{f}[k, j]$, and $l = l[k], m = m[k]$

i.e., the probability of $\mathbf{f}[k, j]$ being aligned to $\mathbf{e}[k, i]$.

36

Step 2: Calculating the Expected Counts

- Calculate the translation counts

$$tcount(e, f) = \sum_{\substack{i,j,k: \\ e[k,i]=e, \\ f[k,j]=f}} a[i, j, k]$$

- $tcount(e, f)$ is expected number of times that e is aligned with f in the corpus

37

Step 3: Re-estimating the Parameters

- New translation probabilities are then defined as

$$\mathbf{T}(f | e) = \frac{tcount(e, f)}{\sum_f tcount(e, f)}$$

- New alignment probabilities are defined as

$$\mathbf{D}(i | j, l, m) = \frac{acount(i, j, l, m)}{\sum_i acount(i, j, l, m)}$$

This defines the mapping from Θ^{t-1} to Θ^t

39

Step 2: Calculating the Expected Counts

- Calculate the alignment counts

$$acount(i, j, l, m) = \sum_{\substack{k: \\ l[k]=l, m[k]=m}} a[i, j, k]$$

- Here, $acount(i, j, l, m)$ is expected number of times that e_i is aligned to f_j in English/French sentences of lengths l and m respectively

38

A Summary of the EM Procedure

- Start with parameters Θ^{t-1} as

$$\begin{array}{l} \mathbf{T}(f | e) \quad \text{for all } f \in \mathcal{F}, e \in \mathcal{E} \\ \mathbf{D}(i | j, l, m) \end{array}$$

- Calculate **alignment probabilities** under current parameters

$$a[i, j, k] = \frac{\mathbf{D}(a_j = i | j, l, m) \mathbf{T}(f_j | e_i)}{\sum_{i'=0}^l \mathbf{D}(a_j = i' | j, l, m) \mathbf{T}(f_j | e_{i'})}$$

- Calculate **expected counts** $tcount(e, f)$, $acount(i, j, l, m)$ from the alignment probabilities

- Re-estimate $\mathbf{T}(f | e)$ and $\mathbf{D}(i | j, l, m)$ from the expected counts

40

The Special Case of Model 1

- Start with parameters Θ^{t-1} as

$$\mathbf{T}(f | e) \quad \text{for all } f \in \mathcal{F}, e \in \mathcal{E}$$

(no alignment parameters)

- Calculate **alignment probabilities** under current parameters

$$a[i, j, k] = \frac{\mathbf{T}(f_j | e_i)}{\sum_{i'=0}^l \mathbf{T}(f_j | e_{i'})}$$

(because $\mathbf{D}(a_j = i | j, l, m) = 1/(l+1)^m$ for all i, j, l, m).

- Calculate **expected counts** $tcount(e, f)$
- Re-estimate $\mathbf{T}(f | e)$ from the expected counts

41

An Example of Training Models 1 and 2

Example will use following translations:

e[1] = the dog
f[1] = le chien

e[2] = the cat
f[2] = le chat

e[3] = the bus
f[3] = l' autobus

NB: I won't use a NULL word e_0

43

Overview

- IBM Model 1
- IBM Model 2
- EM Training of Models 1 and 2
- Some examples of training Models 1 and 2**
- Decoding

42

Initial (random) parameters:

e	f	$\mathbf{T}(f e)$
the	le	0.23
the	chien	0.2
the	chat	0.11
the	l'	0.25
the	autobus	0.21
dog	le	0.2
dog	chien	0.16
dog	chat	0.33
dog	l'	0.12
dog	autobus	0.18
cat	le	0.26
cat	chien	0.28
cat	chat	0.19
cat	l'	0.24
cat	autobus	0.03
bus	le	0.22
bus	chien	0.05
bus	chat	0.26
bus	l'	0.19
bus	autobus	0.27

44

Alignment probabilities:

i	j	k	a(i,j,k)
1	1	0	0.526423237959726
2	1	0	0.473576762040274
1	2	0	0.552517995605817
2	2	0	0.447482004394183
1	1	1	0.466532602066533
2	1	1	0.533467397933467
1	2	1	0.356364544422507
2	2	1	0.643635455577493
1	1	2	0.571950438336247
2	1	2	0.428049561663753
1	2	2	0.439081311724508
2	2	2	0.560918688275492

Expected counts:

<i>e</i>	<i>f</i>	<i>tcount(e, f)</i>
the	le	0.99295584002626
the	chien	0.552517995605817
the	chat	0.356364544422507
the	l'	0.571950438336247
the	autobus	0.439081311724508
dog	le	0.473576762040274
dog	chien	0.447482004394183
dog	chat	0
dog	l'	0
dog	autobus	0
cat	le	0.533467397933467
cat	chien	0
cat	chat	0.643635455577493
cat	l'	0
cat	autobus	0
bus	le	0
bus	chien	0
bus	chat	0
bus	l'	0.428049561663753
bus	autobus	0.560918688275492

Old and new parameters:

<i>e</i>	<i>f</i>	old	new
the	le	0.23	0.34
the	chien	0.2	0.19
the	chat	0.11	0.12
the	l'	0.25	0.2
the	autobus	0.21	0.15
dog	le	0.2	0.51
dog	chien	0.16	0.49
dog	chat	0.33	0
dog	l'	0.12	0
dog	autobus	0.18	0
cat	le	0.26	0.45
cat	chien	0.28	0
cat	chat	0.19	0.55
cat	l'	0.24	0
cat	autobus	0.03	0
bus	le	0.22	0
bus	chien	0.05	0
bus	chat	0.26	0
bus	l'	0.19	0.43
bus	autobus	0.27	0.57

<i>e</i>	<i>f</i>						
the	le	0.23	0.34	0.46	0.56	0.64	0.71
the	chien	0.2	0.19	0.15	0.12	0.09	0.06
the	chat	0.11	0.12	0.1	0.08	0.06	0.04
the	l'	0.25	0.2	0.17	0.15	0.13	0.11
the	autobus	0.21	0.15	0.12	0.1	0.08	0.07
dog	le	0.2	0.51	0.46	0.39	0.33	0.28
dog	chien	0.16	0.49	0.54	0.61	0.67	0.72
dog	chat	0.33	0	0	0	0	0
dog	l'	0.12	0	0	0	0	0
dog	autobus	0.18	0	0	0	0	0
cat	le	0.26	0.45	0.41	0.36	0.3	0.26
cat	chien	0.28	0	0	0	0	0
cat	chat	0.19	0.55	0.59	0.64	0.7	0.74
cat	l'	0.24	0	0	0	0	0
cat	autobus	0.03	0	0	0	0	0
bus	le	0.22	0	0	0	0	0
bus	chien	0.05	0	0	0	0	0
bus	chat	0.26	0	0	0	0	0
bus	l'	0.19	0.43	0.47	0.47	0.47	0.48
bus	autobus	0.27	0.57	0.53	0.53	0.53	0.52

<i>e</i>	<i>f</i>	
the	le	0.94
the	chien	0
the	chat	0
the	l'	0.03
the	autobus	0.02
dog	le	0.06
dog	chien	0.94
dog	chat	0
dog	l'	0
dog	autobus	0
cat	le	0.06
cat	chien	0
cat	chat	0.94
cat	l'	0
cat	autobus	0
bus	le	0
bus	chien	0
bus	chat	0
bus	l'	0.49
bus	autobus	0.51

After 20 iterations:

<i>e</i>	<i>f</i>	$T(f e)$
the	le	0
the	chien	0.4
the	chat	0.3
the	l'	0
the	autobus	0.3
dog	le	0.5
dog	chien	0.5
dog	chat	0
dog	l'	0
dog	autobus	0
cat	le	0.5
cat	chien	0
cat	chat	0.5
cat	l'	0
cat	autobus	0
bus	le	0
bus	chien	0
bus	chat	0
bus	l'	0.5
bus	autobus	0.5

Model 2 has several local maxima – bad one:

<i>e</i>	<i>f</i>	$T(f e)$
the	le	0.67
the	chien	0
the	chat	0
the	l'	0.33
the	autobus	0
dog	le	0
dog	chien	1
dog	chat	0
dog	l'	0
dog	autobus	0
cat	le	0
cat	chien	0
cat	chat	1
cat	l'	0
cat	autobus	0
bus	le	0
bus	chien	0
bus	chat	0
bus	l'	0
bus	autobus	1

Model 2 has several local maxima – good one:

<i>e</i>	<i>f</i>	$T(f e)$
the	le	0
the	chien	0.33
the	chat	0.33
the	l'	0
the	autobus	0.33
dog	le	1
dog	chien	0
dog	chat	0
dog	l'	0
dog	autobus	0
cat	le	1
cat	chien	0
cat	chat	0
cat	l'	0
cat	autobus	0
bus	le	0
bus	chien	0
bus	chat	0
bus	l'	1
bus	autobus	0

another bad one:

Improving the Convergence Properties of Model 2

- Alignment parameters for good solution:

$$\mathbf{T}(i = 1 \mid j = 1, l = 2, m = 2) = 1$$

$$\mathbf{T}(i = 2 \mid j = 1, l = 2, m = 2) = 0$$

$$\mathbf{T}(i = 1 \mid j = 2, l = 2, m = 2) = 0$$

$$\mathbf{T}(i = 2 \mid j = 2, l = 2, m = 2) = 1$$

log probability = -1.91

- Alignment parameters for first bad solution:

$$\mathbf{T}(i = 1 \mid j = 1, l = 2, m = 2) = 0$$

$$\mathbf{T}(i = 2 \mid j = 1, l = 2, m = 2) = 1$$

$$\mathbf{T}(i = 1 \mid j = 2, l = 2, m = 2) = 0$$

$$\mathbf{T}(i = 2 \mid j = 2, l = 2, m = 2) = 1$$

log probability = -4.16

53

- **Out of 100 random starts, only 60 converged to the best local maxima**
- Model 1 converges to the same, global maximum every time (see the Brown et. al paper)
- Method in IBM paper: run Model 1 to estimate \mathbf{T} parameters, then use these as the initial parameters for Model 2
- In 100 tests using this method, Model 2 converged to the correct point every time.

55

Overview

- Alignment parameters for second bad solution:

$$\mathbf{T}(i = 1 \mid j = 1, l = 2, m = 2) = 0$$

$$\mathbf{T}(i = 2 \mid j = 1, l = 2, m = 2) = 1$$

$$\mathbf{T}(i = 1 \mid j = 2, l = 2, m = 2) = 1$$

$$\mathbf{T}(i = 2 \mid j = 2, l = 2, m = 2) = 0$$

log probability = -3.30

- IBM Model 1
- IBM Model 2
- EM Training of Models 1 and 2
- Some examples of training Models 1 and 2
- **Decoding**

54

56

Decoding

- Problem: for a given French sentence f , find

$$\operatorname{argmax}_e P(e)P(f | e)$$

or the “Viterbi approximation”

$$\operatorname{argmax}_{e,a} P(e)P(f, a | e)$$

57

First Stage of the Greedy Method

- For each French word f_j , pick the English word e which maximizes

$$\mathbf{T}(e | f_j)$$

(an inverse translation table $\mathbf{T}(e | f)$ is required for this step)

- This gives us an initial alignment, e.g.,

Bien entendu , il parle de une belle victoire
Well heard , it talking NULL a beautiful victory

(Correct translation: *quite naturally, he talks about a great victory*)

59

Decoding

- Decoding is NP-complete (see (Knight, 1999))
- IBM papers describe a *stack-decoding* or *A* search* method
- A recent paper on decoding:
Fast Decoding and Optimal Decoding for Machine Translation.
Germann, Jahr, Knight, Marcu, Yamada. In ACL 2001.
- Introduces a *greedy* search method
- Compares the two methods to exact (integer-programming) solution

58

Next Stage: Greedy Search

- First stage gives us an initial (e^0, a^0) pair
- Basic idea: define a set of local transformations that map an (e, a) pair to a new (e', a') pair
- Say $\Pi(e, a)$ is the set of all (e', a') reachable from (e, a) by some transformation, then at each iteration take

$$(e^t, a^t) = \operatorname{argmax}_{(e,a) \in \Pi(e^{t-1}, a^{t-1})} P(e)P(f, a | e)$$

i.e., take the highest probability output from results of all transformations

- Basic idea: iterate this process until convergence

60

The Space of Transforms

- CHANGE(j, e):
Changes translation of f_j from e_{a_j} into e
- Two possible cases (take $e_{old} = e_{a_j}$):
 - e_{old} is aligned to more than 1 word, or $e_{old} = NULL$
Place e at position in string that maximizes the alignment probability
 - e_{old} is aligned to exactly one word
In this case, simply replace e_{old} with e
- Typically consider only (e, f) pairs such that e is in top 10 ranked translations for f under $\mathbf{T}(e | f)$
(an inverse table of probabilities $\mathbf{T}(e | f)$ is required – this is described in Germann 2003)

61

The Space of Transforms

- TranslateAndInsert($j, e1, e2$):
Implements CHANGE($j, e1$),
(i.e. Changes translation of f_j from e_{a_j} into $e1$)
and inserts $e2$ at most likely point in the string
- Typically, $e2$ is chosen from the English words which have high probability of being aligned to 0 French words

63

The Space of Transforms

- CHANGE2($j1, e1, j2, e2$):
Changes translation of f_{j1} from $e_{a_{j1}}$ into $e1$,
and changes translation of f_{j2} from $e_{a_{j2}}$ into $e2$
- Just like performing CHANGE($j1, e1$) and CHANGE($j2, e2$)
in sequence

62

The Space of Transforms

- RemoveFertilityZero(i):
Removes e_i , providing that e_i is aligned to nothing in the alignment

64

The Space of Transforms

- $\text{SwapSegments}(i1, i2, j1, j2)$:
Swaps words $e_{i1} \dots e_{i2}$ with words e_{j1} and e_{j2}
- Note: the two segments cannot overlap

65

An Example from Germann et. al 2001

Bien entendu , il parle de une belle victoire
Well heard , it **talking** NULL a **beautiful** victory

⇓

Bien entendu , il parle de une belle victoire
Well heard , it **talks** NULL a **great** victory

$\text{CHANGE2}(5, \text{talks}, 8, \text{great})$

67

The Space of Transforms

- $\text{JoinWords}(i1, i2)$:
Deletes English word at position $i1$, and links all French words that were linked to e_{i1} to e_{i2}

66

An Example from Germann et. al 2001

Bien entendu , il parle de une belle victoire
Well **heard** , it talks **NULL** a great victory

⇓

Bien entendu , il parle de une belle victoire
Well **understood** , it talks **about** a great victory

$\text{CHANGE2}(2, \text{understood}, 6, \text{about})$

68

An Example from Germann et. al 2001

Bien entendu , il parle de une belle victoire

Well understood , it talks about a great victory

↓

Bien entendu , il parle de une belle victoire

Well understood , he talks about a great victory

CHANGE(4, *he*)

69

An Example from Germann et. al 2001

Bien entendu , il parle de une belle victoire

Well understood , he talks about a great victory

↓

Bien entendu , il parle de une belle victoire

quite naturally , he talks about a great victory

CHANGE2(1, *quite*, 2, *naturally*)

70

An Exact Method Based on Integer Programming

Method from Germann et. al 2001:

- Integer programming problems

$$3.2x_1 + 4.7x_2 - 2.1x_3 \quad \text{Minimize objective function}$$

$$\begin{aligned} x_1 - 2.6x_3 &> 5 && \text{Subject to linear constraints} \\ 7.3x_2 &> 7 \end{aligned}$$

- Generalization of travelling salesman problem:
Each town has a number of hotels; some hotels can be in multiple towns. Find the lowest cost tour of hotels such that each town is visited exactly once.

71

- In the MT problem:

- Each city is a French word (all cities visited \Rightarrow all French words must be accounted for)
- Each hotel is an English word matched with one or more French words
- The “cost” of moving from hotel i to hotel j is a sum of a number of terms. E.g., the cost of choosing “not” after “what”, and aligning it with “ne” and “pas” is

$$\begin{aligned} &\log(\text{bigram}(\text{not} \mid \text{what})) + \\ &\log(\mathbf{T}(\text{ne} \mid \text{not})) + \log(\mathbf{T}(\text{pas} \mid \text{not})) \end{aligned}$$

...

72

An Exact Method Based on Integer Programming

- Say distance between hotels i and j is d_{ij} ;
Introduce x_{ij} variables where $x_{ij} = 1$ if path from hotel i to hotel j is taken, zero otherwise

- Objective function: maximize

$$\sum_{i,j} x_{ij} d_{ij}$$

- All cities must be visited once \Rightarrow constraints

$$\forall c \in \text{cities} \quad \sum_{\substack{i \\ \text{i located in } c}} \sum_j x_{ij} = 1$$

73

- Every hotel must have equal number of incoming and outgoing edges \Rightarrow

$$\forall i \quad \sum_j x_{ij} = \sum_j x_{ji}$$

- Another constraint is added to ensure that the tour is fully connected

74