

89688: Statistical Machine Translation

Lecture 3: Statistical Methods, IBM models and the EM algorithm

Roee Aharoni Computer Science Department Bar Ilan University

Based in part on slides from Edinburgh University's MT class and by Kevin Knight

April 2020

 ok-voon ororok sprok . 1b. at-voon bichat dat . 2a. ok-drubel ok-voon anok plok sprok **`** 2b. at-drubel at-voon pippat rrat dat . 3a. erok sprok izok hihok ghirok . \sim 3b. totat dat arrat vat hilat . -----4a. ok-voon anok drok brok jok . 4b. at-voon krat pippat sat lat . 5a. wiwok farok izok stok . 5b. totat jjat quat cat . ______

VS.

6a.	lalok sprok izok jok stok .
6b. 1	wat dat krat quat cat .
7a.	lalok farok ororok lalok sprok izok enemok .
7b. 1	wat jjat bichat wat dat vat eneat .
8a.	lalok brok anok plok nok .
8b.	iat lat pippat rrat nnat .
9a. 1	wiwok nok izok kantok ok-yurp .
9b. ·	totat nnat quat oloat at-yurp .
10a.	lalok mok nok yorok ghirok clok .
10b. v	wat nnat gat mat bat hilat .





• What can we learn from?



- What can we learn from?
 - Parallel Corpora (human translations)



- What can we learn from?
 - Parallel Corpora (human translations)
 - Monolingual corpora (books, wikipedia, the web...)



- What can we learn from?
 - Parallel Corpora (human translations)
 - Monolingual corpora (books, wikipedia, the web...)
- What do we learn?



- What can we learn from?
 - Parallel Corpora (human translations)
 - Monolingual corpora (books, wikipedia, the web...)
- What do we learn?

2020

Translation dictionary:

anok - pippat erok - total ghirok - hilat hihok - arrat izok - vat ok-drubel - at-drubel

ok-yurp - at-yurp ok-voon - at-voon ororok - bichat plok - rrat sprok - dat zanzanok - zanzanat

Translation Model



- What can we learn from?
 - Parallel Corpora (human translations)
 - Monolingual corpora (books, wikipedia, the web...)
- What do we learn?

2020

Translation dictionary:

anok - pippat erok - total ghirok - hilat hihok - arrat izok - vat ok-drubel - at-drubel

Word pair counts:

1 . erok 7 . lalok 2 . ok-drubel 2 . ok-voon 3 . wiwok 1 anok drok 1 anok ghirok

- ok-yurp at-yurp ok-voon - at-voon ororok - bichat plok - rrat sprok - dat zanzanok - zanzanat
 - Translation Model

- 1 hihok yorok
- 1 izok enemok
- 2 izok hihok
- 1 izok jok
- 1 izok kantok
- 1 izok stok
- 1 izok vok

Language Model





• Learning/Training Phase

def learn(parallel_data): # do something return parameters

$$L: (\Sigma_f^* \times \Sigma_e^*)^* \to \Theta$$





• Learning/Training Phase

• Inference Phase

2020

def learn(parallel_data): # do something return parameters

$$L: (\Sigma_f^* \times \Sigma_e^*)^* \to \Theta$$

def translate(French, parameters): # do something return English

 $T: \Sigma_f^* \times \Theta \to \Sigma_e^*$







• Learning/Training Phase

• Inference Phase



def learn(parallel_data): # do something return parameters

$$L: (\Sigma_f^* \times \Sigma_e^*)^* \to \Theta$$

def translate(French, parameters): # do something return English

$$T: \Sigma_f^* \times \Theta \to \Sigma_e^*$$

Using probability: $T(f, \theta) = \arg \max_{e \in \Sigma_e^*} p_{\theta}(e|f)$









2020



2020







• Formalizes...

2020







- Formalizes...
 - The concept of **models**

2020







- Formalizes...
 - The concept of **models**
 - The concept of **data**

2020







- Formalizes...
 - The concept of **models**
 - The concept of data
 - The concept of **learning**

2020







- Formalizes...
 - The concept of **models**
 - The concept of **data**
 - The concept of **learning**
 - The concept of inference (prediction)







- Formalizes...
 - The concept of **models**
 - The concept of **data**
 - The concept of **learning**
 - The concept of inference (prediction)
- Enables to model **ambiguity**









• We would like to **model** the **probability** of a **translation** given a **source** sentence





 We would like to model the probability of a translation given a source sentence

2020

p(English|Chinese) =



- We would like to **model** the **probability** of a translation given a source sentence
- We can use **Bayes Rule**

2020

p(English|Chinese) =



- We would like to **model** the **probability** of a **translation** given a **source** sentence
- We can use **Bayes Rule**

2020

p(English|Chinese) =

 $p(English) \times p(Chinese|English)$ p(Chinese)





- We would like to **model** the **probability** of a **translation** given a **source** sentence
- We can use **Bayes Rule**
 - Why would we want that?

2020

p(English|Chinese) =

 $p(English) \times p(Chinese|English)$ p(Chinese)





- We would like to **model** the **probability** of a **translation** given a **source** sentence
- We can use **Bayes Rule**
 - Why would we want that?

 $p(English) \times p(Chinese|English)$

p(Chinese)

language model

translation model

normalization (ensures we're working with valid probabilities).





How do we define *p*(*Chinese* | *English*)?

2020



How do we define *p*(*Chinese* | *English*)?

• IBM Models (Brown, Dellapietra, Dellapietra, and Mercer, 93')



How do we define *p(Chinese|English)*?

- IBM Models (Brown, Dellapietra, Dellapietra, and Mercer, 93')
- "We define a concept of word-by-word **alignment**"



How do we define *p*(*Chinese* | *English*)?

- IBM Models (Brown, Dellapietra, Dellapietra, and Mercer, 93')
- "We define a concept of word-by-word alignment"





How do we define *p(Chinese|English)*?

- IBM Models (Brown, Dellapietra, Dellapietra, and Mercer, 93')
- "We define a concept of word-by-word **alignment**"







Understanding Alignments

ROEE AHARONI



Understanding Alignments

Alignment function



 $a: \{1 \rightarrow 1, 2 \rightarrow 2, 3 \rightarrow 3, 4 \rightarrow 4\}$



Understanding Alignments

- Alignment function
- Reordering



 $a: \{1 \rightarrow 3, 2 \rightarrow 4, 3 \rightarrow 2, 4 \rightarrow 1\}$


Understanding Alignments

- Alignment function
- Reordering
- One-to-Many





Understanding Alignments

- Alignment function
- Reordering
- One-to-Many
- Dropping words



 $a: \{1 \rightarrow 2, 2 \rightarrow 3, 3 \rightarrow 4\}$



Understanding Alignments

0

- Alignment function
- Reordering
- One-to-Many
- Dropping words
- Inserting words



 $a: \{1 \rightarrow 1, 2 \rightarrow 2, 3 \rightarrow 3, 4 \rightarrow 0, 5 \rightarrow 4\}$





• Given a source sentence, how was the target sentence generated?



• Given a source sentence, how was the target sentence generated?





• Given a source sentence, how was the target sentence generated?

Although north wind howls , but sky still very clear . 虽然北风呼啸,但天空依然十分清澈。 ε





- Given a source sentence, how was the target sentence generated?
- Sample a length for the target sentence

Although north wind howls , but sky still clear . very 虽然北风呼啸,但天空依然十分清澈。 ε

 $p(English \ length|Chinese \ length)$





- Given a source sentence, how was the target sentence generated?
- Sample a length for the target sentence
- For each target position:

Although north wind howls , but sky still clear very 虽然北风呼啸,但天空依然十分清澈。 ε

 $p(English \ length|Chinese \ length)$





- Given a source sentence, how was the target sentence generated?
- Sample a length for the target sentence
- For each target position:
 - Sample an alignment (to a position in the source)

Although north wind howls , but sky still very clear 虽然北风呼啸,但天空依然十分清澈

p(Chinese word position)





- Given a source sentence, how was the target sentence generated?
- Sample a length for the target sentence
- For each target position:
 - Sample an alignment (to a position in the source)
 - Sample a word translation given this alignment

Although north wind howls , but sky still very clear . 虽然北风呼啸,但天空依然十分清澈。 ε

However

p(English word|Chinese word)



- Given a source sentence, how was the target sentence generated?
- Sample a length for the target sentence
- For each target position:
 - Sample an alignment (to a position in the source)
 - Sample a word translation given this alignment
 - Repeat until done



However,

p(English word|Chinese word)



- Given a source sentence, how was the target sentence generated?
- Sample a length for the target sentence
- For each target position:
 - Sample an alignment (to a position in the source)
 - Sample a word translation given this alignment
 - Repeat until done



However, the sky remained clear under the strong north wind.





$p(\mathbf{f}, \mathbf{a} | \mathbf{e}) = p(I|J) \prod_{i=1}^{I} p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1





$p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod_{i=1}^{I} p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1





ROEE AHARONI

sample alignment

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod_{i=1}^{I} p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1





ROEE AHARONI

sample alignment

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod_{i=1}^{n} p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1

sample word translation







• f - foreign sentence (Chinese)

ROEE AHARONI

sample alignment

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod_{i=1}^{n} p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1

sample word translation







- f foreign sentence (Chinese)
- a alignment

ROEE AHARONI

sample alignment

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1

sample word translation







- f foreign sentence (Chinese)
- a alignment
- e English sentence

ROEE AHARONI

sample alignment

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1

sample word translation







- f foreign sentence (Chinese)
- a alignment
- e English sentence
- I foreign sentence length

ROEE AHARONI

sample alignment

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1

sample word translation







- f foreign sentence (Chinese)
- a alignment
- e English sentence
- I foreign sentence length
- J English sentence length

ROEE AHARONI

sample alignment

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1

sample word translation







- f foreign sentence (Chinese)
- a alignment
- e English sentence
- I foreign sentence length
- J English sentence length
- a_i alignment of i-th foreign word

ROEE AHARONI

sample alignment

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1

sample word translation







- f foreign sentence (Chinese)
- a alignment
- e English sentence
- I foreign sentence length
- J English sentence length
- a_i alignment of i-th foreign word
- f_i foreign word in position i

ROEE AHARONI

sample alignment

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1

sample word translation







- f foreign sentence (Chinese)
- a alignment
- e English sentence
- I foreign sentence length
- J English sentence length
- a_i alignment of i-th foreign word
- f_i foreign word in position i
- e_{a_i} English word in position a_i

ROEE AHARONI

sample alignment

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod p(a_i|J) \cdot p(f_i|e_{a_i})$ i=1

sample word translation







2020

$p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod_{i=1}^{I} p(a_i|J) \cdot p(f_i|e_{a_i})$





• What are the parameters we should learn?

2020

$p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod_{i=1}^{r} p(a_i|J) \cdot p(f_i|e_{a_i})$





- What are the parameters we should learn?
 - Sentence length distributions

2020

$p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod p(a_i|J) \cdot p(f_i|e_{a_i})$

Sentence length distributions





- What are the parameters we should learn?
 - Sentence length distributions
 - Alignment distributions

Alignment distributions

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod p(a_i|J) \cdot p(f_i|e_{a_i})$

Sentence length distributions





- What are the parameters we should learn?
 - Sentence length distributions
 - Alignment distributions
 - Word translation distributions

Alignment distributions

 $p(\mathbf{f}, \mathbf{a}|\mathbf{e}) = p(I|J) \prod p(a_i|J) \cdot p(f_i|e_{a_i})$

Sentence length distributions

Word translation distributions







If we know the alignments, easy:



- If we **know** the **alignments**, easy:
 - *p*(*I*|*J*) learn by counting

Aligned Chinese sentences of length I p(|J) =

English sentences of length J



- If we **know** the **alignments**, easy:
 - p(I|J) learn by counting
 - Alignment distributions use uniform distribution





- If we **know** the **alignments**, easy:
 - p(I|J) learn by counting
 - Alignment distributions use uniform distribution
 - Word translation distributions again by counting





- If we **know** the **alignments**, easy:
 - p(I|J) learn by counting
 - Alignment distributions use uniform distribution
 - Word translation distributions again by counting
- But do we know the alignments?






 Parameters (translation) probabilities) and alignments are both unknown!



- Parameters (translation) probabilities) and alignments are both unknown!
- If we knew the alignments in the training data, we could calculate the parameters by counting (prev. slide)



- Parameters (translation) probabilities) and alignments are both unknown!
- If we knew the alignments in the training data, we could calculate the parameters by counting (prev. slide)
- If we knew the parameters, we could calculate the *expected* alignment counts



- Parameters (translation) probabilities) and alignments are both unknown!
- If we knew the alignments in the training data, we could calculate the parameters by counting (prev. slide)
- If we knew the parameters, we could calculate the *expected* alignment counts







• Dempster, Laird and Rubin (1977)



- Dempster, Laird and Rubin (1977)
- One of the most widely used algorithms in machine learning



- Dempster, Laird and Rubin (1977)
- One of the most widely used algorithms in machine learning
- Many applications, e.g. clustering



- Dempster, Laird and Rubin (1977)
- One of the most widely used algorithms in machine learning
- Many applications, e.g. clustering
- Main idea start randomly, and iterate until convergence:





- Dempster, Laird and Rubin (1977)
- One of the most widely used algorithms in machine learning
- Many applications, e.g. clustering
- Main idea start randomly, and iterate until convergence:
 - Calculate expected counts for missing data (expectation, or E-step) using current model





- Dempster, Laird and Rubin (1977)
- One of the most widely used algorithms in machine learning
- Many applications, e.g. clustering
- Main idea start randomly, and iterate until convergence:
 - Calculate expected counts for missing data (expectation, or E-step) using current model
 - Find new model parameters that maximize the data likelihood (maximization, or M-step)







• Start with all alignments equally likely





- Start with all alignments equally likely
- In each iteration:





- Start with all alignments equally likely
- In each iteration:
 - look at the entire dataset and sum the (expected) alignments we saw for each word and it's possible translations (E-step)





- Start with all alignments equally likely
- In each iteration:
 - look at the entire dataset and sum the (expected) alignments we saw for each word and it's possible translations (E-step)
 - Update the translation probabilities according to those global counts (M-step)...





- Start with all alignments equally likely
- In each iteration:
 - look at the entire dataset and sum the (expected) alignments we saw for each word and it's possible translations (E-step)
 - Update the translation probabilities according to those global counts (M-step)...
 - ...which will update the alignment counts





- Start with all alignments equally likely
- In each iteration:
 - look at the entire dataset and sum the (expected) alignments we saw for each word and it's possible translations (E-step)
 - Update the translation probabilities according to those global counts (M-step)...
 - ...which will update the alignment counts
- Repeat steps above





- Start with all alignments equally likely
- In each iteration:
 - look at the entire dataset and sum the (expected) alignments we saw for each word and it's possible translations (E-step)
 - Update the translation probabilities according to those global counts (M-step)...
 - ...which will update the alignment counts
- Repeat steps above
- Until convergence







• We need to compute:



- We need to compute:
 - Expectation step: probability of alignments

2020

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



- We need to compute:
 - Expectation step: probability of alignments
 - Maximization step: probability of word translations

2020

$p(a|\mathbf{e}, \mathbf{f})$ $t(e|f; \mathbf{e}, \mathbf{f})$





Apply the chain rule

2020

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



- Apply the chain rule
- Numerator IBM Model 1 definition

2020

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



- Apply the chain rule
- Numerator IBM Model 1 definition
- Denominator:

2020

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

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

= $\sum_{a(1)=0}^{l_f} \dots \sum_{a(l_e)=0}^{l_f} p(\mathbf{e}, a|\mathbf{f})$
= $\sum_{a(1)=0}^{l_f} \dots \sum_{a(l_e)=0}^{l_f} \frac{\epsilon}{(l_f+1)^{l_e}} \prod_{j=1}^{l_e} t(e_j)$





- Apply the chain rule
- Numerator IBM Model 1 definition
- Denominator:
 - Marginalize over all possible alignments

2020

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

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

= $\sum_{a(1)=0}^{l_f} \dots \sum_{a(l_e)=0}^{l_f} p(\mathbf{e}, a|\mathbf{f})$
= $\sum_{a(1)=0}^{l_f} \dots \sum_{a(l_e)=0}^{l_f} \frac{\epsilon}{(l_f+1)^{l_e}} \prod_{j=1}^{l_e} t(e_j)$





- Apply the chain rule
- Numerator IBM Model 1 definition
- Denominator:
 - Marginalize over all possible alignments
 - IBM model 1 definition

2020

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

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

= $\sum_{a(1)=0}^{l_f} \dots \sum_{a(l_e)=0}^{l_f} p(\mathbf{e}, a|\mathbf{f})$
= $\sum_{a(1)=0}^{l_f} \dots \sum_{a(l_e)=0}^{l_f} \frac{\epsilon}{(l_f+1)^{l_e}} \prod_{j=1}^{l_e} t(e_j)$











• We want to get rid of the exponential number of multiplications







- We want to get rid of the exponential number of multiplications
- Move the constants out







- We want to get rid of the exponential number of multiplications
- Move the constants out
- Last trick change sum of products to product of sums







• So finally we got:


So finally we got:

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

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



 $f_i)$

- So finally we got:
 - Probability of alignment given a sentence pair

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

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



 f_i

- So finally we got:
 - Probability of alignment given a sentence pair
 - Based on the translation probabilities alone

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

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



 f_i

- So finally we got:
 - Probability of alignment given a sentence pair
 - Based on the translation probabilities alone
 - We will use this to get the expected alignment counts

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

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



 f_i



• We want to collect alignment counts to compute the new translation parameters, using the **current** translation parameters



- We want to collect alignment counts to compute the new translation parameters, using the **current** translation parameters
- For each possible pair (e,f) in each example (**e**,**f**) sum the expected counts of this translation



- We want to collect alignment counts to compute the new translation parameters, using the **current** translation parameters
- For each possible pair (e,f) in each example (**e**,**f**) sum the expected counts of this translation

2020

$c(e|f; \mathbf{e}, \mathbf{f}) = \sum p(a|\mathbf{e}, \mathbf{f}) \sum^{e} \delta(e, e_j) \delta(f, f_{a(j)})$



- We want to collect alignment counts to compute the new translation parameters, using the **current** translation parameters
- For each possible pair (e,f) in each example (**e**,**f**) sum the expected counts of this translation

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



- We want to collect alignment counts to compute the new translation parameters, using the **current** translation parameters
- For each possible pair (e,f) in each example (**e**,**f**) sum the expected counts of this translation

2020

p(the|la) = 0.7 p(house|la) = 0.05p(the|maison) = 0.1 p(house|maison) = 0.8

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



- We want to collect alignment counts to compute the new translation parameters, using the **current** translation parameters
- For each possible pair (e,f) in each example (**e**,**f**) sum the expected counts of this translation

2020

p(





- We want to collect alignment counts to compute the new translation parameters, using the **current** translation parameters
- For each possible pair (e,f) in each example (**e**,**f**) sum the expected counts of this translation
- Maximization: after we do this over the entire corpus, sum and normalize to get the **new** parameters

2020

$p(\text{the}|\text{la}) = 0.7 \qquad p(\text{house}|\text{la}) = 0.05$ $p(\text{the}|\text{maison}) = 0.1 \qquad p(\text{house}|\text{maison}) = 0.8 \qquad 0.875 \qquad 0.875 \qquad 0.941$ $c(e|f;\mathbf{e},\mathbf{f}) = \sum_{a} p(a|\mathbf{e},\mathbf{f}) \sum_{i=1}^{l_e} \delta(e,e_j) \delta(f,f_{a(j)})$ the house t(o|f) l_e

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

 $t(e|f; \mathbf{e}, \mathbf{f}) = \frac{\sum_{(\mathbf{e}, \mathbf{f})} c(e|f; \mathbf{e}, \mathbf{f}))}{\sum_{f} \sum_{(\mathbf{e}, \mathbf{f})} c(e|f; \mathbf{e}, \mathbf{f}))}$







• Finally:



• Finally:

Input: set of sentence pairs (**e**, **f**) **Output:** translation

- 1: initialize t(e|f) u
- 2: while not conver
- // initialize 3:
- $\operatorname{count}(e|f) = 0$ 4:
- total(f) = 0 for 5:
- for all sentence 6:
- // compute 7:

8:

- for all words
- s-total(e)9:
- for all wo 10:
- s-total(11:
- end for 12:
- end for 13:

ce pairs (e , f)	14:
prob. $t(e f)$	15
iniformly	16
rged do	17:
0 for all a f	18:
o for all e, f	19:
or all f	20:
e pairs (e,t) do	21:
normalization	22
s e in e do	22
= 0	23
ords f in f do	24:
(a) = t(a f)	25
$e_{J} + = \iota(e_{J})$	26:

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











- Finally:
- You will implement this in Exercise 1

Input: set of sentence pairs (**e**, **f**) **Output:** translation

- 1: initialize t(e|f) u
- 2: while not conver
- // initialize 3:
- count(e|f) = 04:
- total(f) = 0 for 5:
- for all sentence 6:
- // compute 7:

8:

9:

- for all words
- s-total(e)
- for all wo 10:
- s-total(11:
- end for 12:
- end for 13:

ce pairs (e , f)	14:
prob. $t(e f)$	15
iniformly	16
rged do	17:
0 for all a f	18:
o for all e, f	19:
or all f	20:
e pairs (e,t) do	21:
normalization	22
s e in e do	22
= 0	23
ords f in f do	24:
(a) = t(a f)	25
$e_{J} + = \iota(e_{J})$	26:

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













• How can we measure the alignment quality?



- How can we measure the alignment quality?
- AER Och and Ney, 2000



AER(A, S, P) = (1

= (

					-			en
								1978
e								,
-	-	-			-		-	on
ciblo	-	-			-			a
SIDIE			,					enregistré
						×		1,122,000
dicted				\bigcirc				divorces
	-	-		-	-		-	sur
			,					le
$ A \cap S + A \cap P $		·						continent
$1 - \frac{ A + \beta + A + F }{ A + G }$								
$ A + S \qquad f$	-	-		-	-		-	
	in	978	ans	ced	000	nes	•	
3+3 1		ï	ic	OL(2,(tir		
$\left(1 - \frac{1}{3 + 4}\right) = \frac{1}{7}$			ner	div	,12			
5 1 4/ 1			Aı	Ŭ	Ŧ			



- How can we measure the alignment quality?
- AER Och and Ney, 2000
- Possible contains Sure

= Sure = Pos= Prec

$$AER(A, S, P) = \left($$

= (

					-			en
								1978
re			,					,
	-	-			-			on
sciblo	-	-		-	-			a
SIDIC								enregistré
								1,122,000
dicted		-		\bigcirc	-			divorces
	-	-		-	-			sur
			,	,				le
$(A \cap S + A \cap P)$					-			continent
$1 - \frac{ A + \beta + A + F }{ A + A }$	-	-						
$ A + S \qquad f$	-	-		-	-			
	i'n	978	ans	ced	000	nes	•	
(3+3) 1		ï	ic	OL	2,(tir		
$\begin{pmatrix} 1 - \frac{1}{3+4} \end{pmatrix} = \frac{1}{7}$			mer	div	,12			
			A		Ţ			



- How can we measure the alignment quality?
- AER Och and Ney, 2000
- Possible contains Sure
- Must hit all Sure to be perfect, ok to not cover all probable

= Sure = Pos= Prec

 $AER(A, S, P) = \Big($

=

					-			en
								1978
e								,
-	-	-			-		-	on
ciblo	-	-			-			a
SIDIE			,					enregistré
						×		1,122,000
dicted				\bigcirc				divorces
	-	-		-	-		-	sur
			,					le
$ A \cap S + A \cap P $		·						continent
$1 - \frac{ A + \beta + A + F }{ A + G }$								
$ A + S \qquad f$	-	-		-	-		-	
	in	978	ans	ced	000	nes	•	
3+3 1		ï	ic	OL(2,(tir		
$\left(1 - \frac{1}{3 + 4}\right) = \frac{1}{7}$			ner	div	,12			
5 1 4/ 1			Aı	Ŭ	Ŧ			



• IBM model 1

- IBM model 1
 - Generative model, based on:

- IBM model 1
 - Generative model, based on:
 - Alignments (latent variables)



- IBM model 1
 - Generative model, based on:
 - Alignments (latent variables)
 - Which are used to calculate translation parameters



- IBM model 1
 - Generative model, based on:
 - Alignments (latent variables)
 - Which are used to calculate translation parameters
 - Using Expectation Maximization



- IBM model 1
 - Generative model, based on:
 - Alignments (latent variables)
 - Which are used to calculate translation parameters
 - Using Expectation Maximization
- <u>Exercise 1</u> May 18th



- IBM model 1
 - Generative model, based on:
 - Alignments (latent variables)
 - Which are used to calculate translation parameters
 - Using Expectation Maximization
- <u>Exercise 1</u> May 18th



Turkish: Otel-imiz-in karşı-sın-da-ki dükkân-da gör-düğ-üm bir elbise-yi dene-mek iste-r-im. inverse order of morphemes and concepts I'd like to try on a suit I've see-n in a shop across the street from our hotel. English:



