

Natural Language Processing II

Khalil Sima'an
Universiteit van Amsterdam

Natural Language Processing II

Why Machine Translation?

Main Questions and General Approach?

Vanilla Treasures of Machine Translation (MT)

Major interest by industry!



Vanilla Treasures of Machine Translation (MT)

Major interest by industry!



amazon.com

YAHOO!



GET RICH TODAY: *Translate better than Google!*

Technological motivation

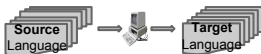
- Cultural, economic and societal impact
- Huge volume that never gets translated
- MT is enabling: Speed + Low cost



BUT why conduct research on MT (beside technology)?

Why Machine Translation (MT)?

Technological challenge...

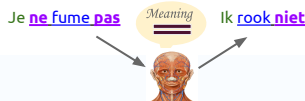


Scientific challenge Human Language Understanding

We never observe “meaning” in the wild.

But translation Data has two crucial properties

- Human *meaning preserving behavior*: $\text{Meaning}(I) == \text{Meaning}(O)$
- Both Input and Output observable.



Translation Data == Translation Equations == Meaning Equations

Motivation 1: Find the Latent Structure of Translation Equations

Motivation 2: How to Translate Correctly, i.e., Build new equations

The Structure of Equivalence?

Sentence-level translation equations

(De zonnestrallen die door het raam binnenkomen)
==
(The sun rays that infiltrate through the window)

But how are translation equations built-up? Important for generalization.

The Structure of Equivalence? Analogy

Sentence-level translation equations

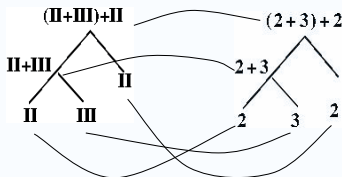
(De zonnestralen die door het raam binnenkomen) == (The sun rays that infiltrate through the window)

But how are translation equations built-up? Important for generalization.

Analogy: Decomposition of equations

- Two decimal alphabets;
- We know the “atomic units”, e.g., II=2, III=3
- Non-ambiguous translation
- One composition operator (+)
- No idioms, just composition

✦ Easy to decompose recursively



Recursive Structure of Translation Equivalence, How?

Translation Equivalence: Challenges

Induce mapping

Parallel Corpus: A large sample of source-target pairs of human translations.

| | | |
|---|---|---|
| I ran up a big bill. | Ik heb een grote rekening opgelopen . | Ambiguity |
| I ran up a big hill. | Ik rende een grote heuvel op . | Stochastic decisions |
| He destroyed them. | Hij richtte hen ten gronde . | Idioms: how to identify? |
| Je ne fume pas Ik rook niet | Je ne VP-F pas Ik VP-N niet . | Non-Contiguous mapping. |
| The president meets Saudi economic officials يسقبل الرئيس مسؤولين اقتصاديين سعوديين | The president meets a Saudi economic official يسقبل سعودي اقتصادي مسؤول الرئيس يستقبل | Morph. Variations Canonical forms? |
| 澳洲 是 与 北韩 有 邦交 的 少数 国家 之一。 Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi . Australia is with North Korea have dipl. rels. that few countries one of . Australia is one of the few countries that have diplomatic relations with North Korea. | Word Order Differences Mappings with permuted word order: huge space (n!). Example from (Chiang 2007) | |

Let us concentrate on Word Order for now

Structure of Translation Equivalence: Word Order

Parallel Corpus: A large sample of source-target pairs of human translations.

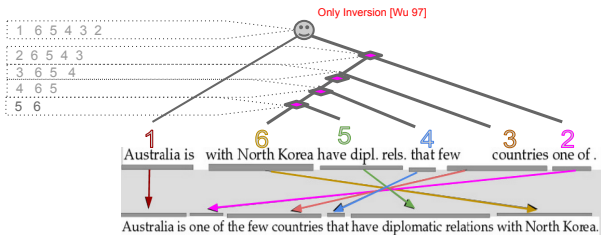
Induce word-level mapping: Many-to-Many Word Alignment induced as latent structure [existing].*

| | | | | | | | | | | | |
|------------------|--|----------------|-----------------|-----------------------------|------------------|-------------------------------|-----------------|----------------------|---------------------------|-----------------------|---|
| Source positions | 1 | 2 | 3 | 4 | 5 | 6 | | | | | |
| Source words | 澳洲 Aozhou Australia | 是 shi is | 与 yu with | 北韩 Beihan North Korea | 有 you have | 邦交 bangjiao dipl. rels. | 的 de that | 少数 shaoshu few | 国家 guojia countries | 之一 zhiyi one of | 。 |
| Alignments | | | | | | | | | | | |
| Target words | Australia is one of the few countries that have diplomatic relations with North Korea. | | | | | | | | | | |
| Target positions | 1 | 6 | 5 | 4 | 3 | | | | | 2 | |

* Up-to some encapsulation of idioms, morphology, unaligned words: First approximation available [Brown et al 1992; Och & Ney 2003].

Alignments \neq Sequence of individual alignments

Hierarchical Word Order (Surface Composition)

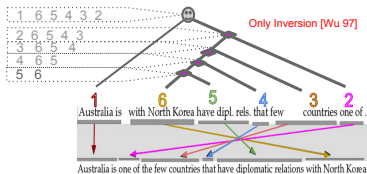


Hierarchical composition could have benefits: Long range reordering

The Questions of Translation Equivalence

Q1. How to learn translation equivalence over “words”?

- What are the units of equivalence?
- Which units map to which?
- What composition is needed to learn this?



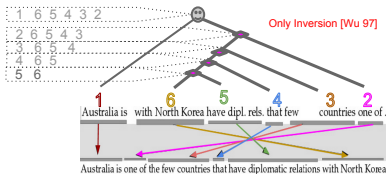
Q3. How to compose new translations from old ones to preserve meaning?

- Which representations?
- Which compositions preserve meaning?

Hierarchical vs. Sequential View: Applications

Q1. How to learn translation equivalence over “words”?

- What are the units of equivalence?
- Which units map to which?
- What composition is needed to learn this?



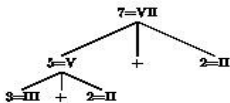
Q.2 How to learn word order from examples?

- Word order as first big challenge!
- Structure of equivalence?
- NLP II mostly about this!

Q3. How to compose new translations from old ones to preserve meaning?

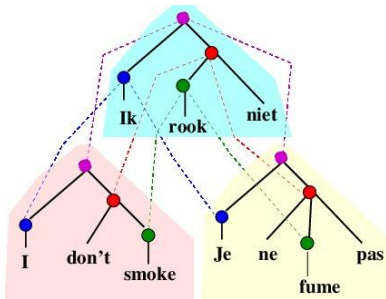
- Which representations?
- Which compositions preserve meaning?

Hierarchical Equivalence: Questions



Questions:

- **Lexicon:** Which word-translation pairs?
- **Structure:** Which composition structure?
- **Composition:** Which operators?



- . How to learn all this from parallel data?
- . How to deal with ambiguity?

The Structure of NLP II

1. How to learn a lexicon and mapping between words? Sequential view.
 - a. Word-based models and word alignments (IBM Models)
 - b. Inducing alignments and using them for extracting phrases, i.e., translation equations at any level, not only sentence level
2. How to evaluate Machine Translation system output?
3. How to learn hierarchical models based on Synchronous Grammars?
 - a. Synchronous grammars
 - b. Hierarchical phrased-based model
4. How obtain semantic representations from multilingual data?
5. How to learn models of word-order differences (reordering) between languages?
 - a. Permutations and their decomposition/factorization
 - b. Synchronous grammars and permutations
 - c. Learning from data

Statistical Machine Translation: First Steps

Translation between Languages

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

Human translators translate texts from a source language to a target language.

I don't smoke.

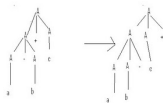
Je ne fume pas



Can we build a computer program that translates texts from one language to another?



What challenges will we face and how do we tackle them?



So many languages, so little time

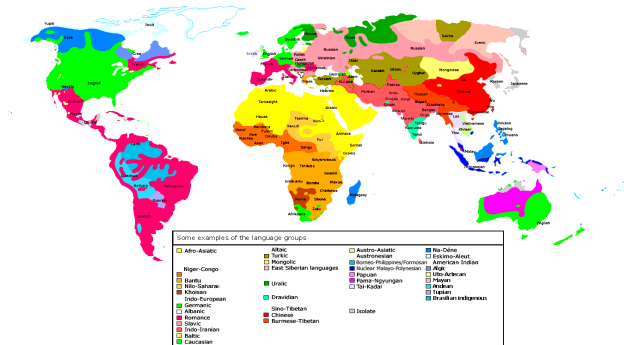
Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models



- Are the differences between languages arbitrary?
- Are there shared regularities between different languages?

How should we automatically translate?

History I: Premature Optimism and Failure

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

Some history on translation and speech recognition:

- During 50's and 60's
First computers; Chomsky's grammars; programming languages; big optimism and huge funding

`Translation is Easy: we can program this!!`

ALPAC (Automatic Language Processing Advisory Committee) Report 1966 (U.S. Government).

Failure: AI abandons NLP, NLP abandons Translation

- During 70's and 80's:
AI: "You need world-knowledge: build an ontology"
CS: Concentrate on Information Retrieval
Linguistics: We need better theory

History II: Renewed Optimism

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

- During 70's: A group of statisticians at IBM TJ Watson "digs up" an old idea ([Weaver 1948, 1949]):

When I look at an article in Russian, I say: 'This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.

Communication and Information Theory (Shannon, Weaver); Code breaking (Turing).

- During 80's: Success in ASR; Look at Translation
- During 90's: Success in parsing and Translation
- By 2006: Google introduces "Google Translate"!

Next: How good is statistical MT these days?

Modeling Human Translation Expertise

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

- Translators are Experts in Translation
- Humans: Study, work, acquire by experience . . .

Can we model “expertise acquisition” from experience?

- Observe and learn how humans translate?
- Use input-output translation examples: Parallel corpora
No access to what happens in between
- How do we build and select the correct translation?
Ambiguity is stalking us all the way.

How can we learn translation regularities from data?

Data and Statistical Models

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

Parallel corpus = a collection of text-chunks and their translations.

Parallel corpora are the by-product of *human translation*.
Every source chunk is paired with a target chunk.

| Dutch | English |
|---|---|
| De prijs van het huis is gestegen. | The price of the house has risen. |
| Het huis kan worden verkocht. | The house can be sold. |
| Als het de marktprijs daalt zullen sommige gezinnen een zware tijd doormaken. | If the market price goes down, some families will go through difficult times. |
| : | : |
| : | : |
| : | : |
| : | : |
| : | : |
| : | : |

- Hansards Canadian Parliament Proc. (English-French).
- European Parliament Proc. (23 languages).
- United Nations documents.
- Newspapers: Chinese-English; Arabic-English; Urdu-English.

The hidden structure of translation

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

How to model the translation mapping in the data?

The big cities will join forces if the prime minister maintains his present policy in the long run.

????

De grote steden zullen samen optrekken als de premier zijn huidige beleid op lange termijn blijft handhaven.

What is the nature of the mapping?

$$\text{"Translate(sentence) = } \sum_i \text{Translate(part}_i\text{)"} ??$$

- What are $part_i$ and $\text{Translate}(part_i)$ in the data?
- What is \sum_i ?
- How to model differences in word-order, morphology etc?
- What about ambiguity, idioms etc?

Probabilistic Modeling: Simple Noisy Channel

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

Source sentence $\mathbf{s} = s_1, \dots, s_n$

Target sentence $\mathbf{t} = t_1, \dots, t_n$

$$\arg \max_{\mathbf{t}} P(\mathbf{t} \mid \mathbf{s}) = \arg \max_{\mathbf{t}} P(\mathbf{t}) \times P(\mathbf{s} \mid \mathbf{t})$$

- **Target Language Model $P(\mathbf{t})=?$**

How regular is a given string \mathbf{t} in the target language?

$$P(\mathbf{t}) = \sum_{\mathbf{d}} P(\mathbf{t}, \mathbf{d})$$

Derivations \mathbf{d} : Finite-State / Context-Free Grammar

- **Translation Model $P(\mathbf{s} \mid \mathbf{t})=?$**

How to model the mapping $\mathbf{t} \rightarrow \mathbf{s}$?

This course: Learning translation models from data

Modeling Parallel Corpus Data

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

- How to represent the source sentence?
- How to represent the target sentence?
- How to model the mapping between these representations?
We need to model sentence pairs!!
- Is translation compositional?
- Some options: Probabilistic Synchronous Grammars, Probabilistic Tree Transducers, etc.
- What learning algorithms?
- How to automatically evaluate translation output?

Data and Models: Structure of lecture

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

- General statistical framework
- Word-based models: word alignments
- Phrase-based models: phrase-alignments
- Tree-based models: tree-alignments

Introduction to Statistical Machine Translation

Statistical Approach: Parallel Corpora

Natural Language Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based Models

Task: Translate a source sentence \mathbf{f} to a target sentence \mathbf{e} .
Data: Parallel corpus (source-target sentence pairs).





[Daily Briefing](#) / [Press Releases](#) / [Radio, TV, Photo, Documents, Maps](#) / [Publications](#) / [Stamp](#) / [Database](#) / [UN Work](#)
[Peace & Security](#) / [Economic & Social Development](#) / [Human Rights](#) / [Sustainable Africa](#) / [International Law](#)

Welcome to the United Nations



Secretary-General
Situation in the Middle East
UN Action in Iraq
Renewing the United Nations
UN Action against Terrorism
Issues on the UN Agenda
Civil Society & Business
UN Webcast
CyberSchoolBus

International Conference on Financing for Development, 2008
29 November - 2 December 2008

[Home](#) / [Recent Additions](#) / [Employment](#) / [Procurement](#) / [Comments](#) / [Q&A](#) / [UN System Sites](#) / [Index](#) / [Search](#)
 © 1996 - 2008 [UN.org](#) / [English](#) / [Français](#) / [Русский](#) / [Español](#)

[illegible]

Source-Channel Approach: IBM Models (1990's)

Parallel Corpus Example

Parallel corpus **C** = a collection of text-chunks and their translations.

Parallel corpora are the by-product of *human translation*.
Every source chunk is paired with a target chunk.

| Dutch | English |
|---|---|
| De prijs van het huis is gestegen. | The price of the house has risen. |
| Het huis kan worden verkocht. | The house can be sold. |
| Als het de marktprijs daalt zullen sommige gezinnen een zware tijd doormaken. | If the market price goes down, some families will go through difficult times. |
| . | . |
| . | . |
| . | . |
| . | . |
| . | . |
| . | . |

- Hansards Canadian Parliament Proc. (English-French).
- European Parliament Proc. (23 languages).
- United Nations documents.
- Newspapers: Chinese-English; Arabic-English; Urdu-English.
- TAUS corpora.

Generative Source-Channel Framework

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

Given source sentence \mathbf{f} , select target sentence \mathbf{e}

$$\arg \max_{\mathbf{e} \in E(\mathbf{f})} \{ P(\mathbf{e} | \mathbf{f}) \} = \arg \max_{\mathbf{e} \in E(\mathbf{f})} \{ \overbrace{P(\mathbf{e})}^{L.M.} \times \overbrace{P(\mathbf{f} | \mathbf{e})}^{T.M.} \}$$

Set $E(\mathbf{f})$ is the set of hypothesized translations of \mathbf{f} .

$P(\mathbf{f} | \mathbf{e})$: accounts for divergence in ...

- word order
- morphology
- syntactic relations
- idiomatic ways of expression
- :

How to estimate $P(\mathbf{e} | \mathbf{f})$? **Sparse-data problem!**

Inducing The Structure of Translation Data

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

e = Mary did not slap the green witch .

? ? ? ?

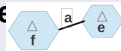
f = Maria no dio una bofetada a la bruja verde .



The latent structure of translation equivalence

Graphical representations Δ_f and Δ_e for **f** and **e**

Relation **a** between Δ_f and Δ_e



$$\arg \max_{\mathbf{e} \in E(\mathbf{f})} \{ P(\mathbf{e} \mid \mathbf{f}) \} =$$

$$\arg \max_{\mathbf{e} \in E(\mathbf{f})} \{ \sum_{\langle \Delta_f, \mathbf{a}, \Delta_e \rangle} P(\mathbf{e}, \Delta_f, \Delta_e, \mathbf{a} \mid \mathbf{f}) \}$$

The difficult question: Which $\Delta_{f/e}$ and **a** fit data best?

Structure in current models

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

$$\Delta_f \xrightarrow{\mathbf{a}} \Delta_e$$

In most current models structure of **reordering**:

- $\Delta_{f/e}$ are structures over word positions.
- \mathbf{a} is an **alignment** between groups of word positions in Δ_f and Δ_e .
- **Challenge:** Number of permutations of n words is $n!$

Structure shows translation units **composing** together

- What are the atomic translation units?
- How these compose together **efficiently**?
- How to put probs. on these structures?

Structure helps combat sparsity and complexity

Structure in Existing Models: Sketch

Natural
Language
Processing II

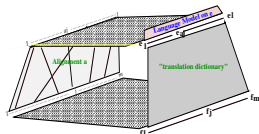
Dr Khalil
Sima'an

Challenges

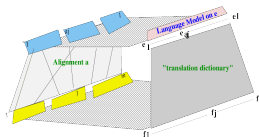
This course

Word-Based
Models

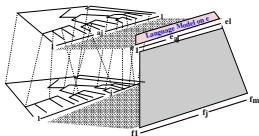
Word-based



Phrase-based



Tree-based



Problem: No sufficient stats to estimate $P(\mathbf{e} \mid \mathbf{f})$ from data

Word-Based Models: Word Alignments

Some History and References

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

Statistical models with word-alignments:

- Brown, Cocke, Della Pietra, Della Pietra, Jelinek, Lafferty, Mercer and Roossin. A statistical approach to machine translation. Computational Linguistics, 1990.
- Brown, Della Pietra, Della Pietra and Mercer. The mathematics of statistical machine translation: parameter estimation., Computational Linguistics, 1993.
- Och and Ney: A Systematic Comparison of Various Statistical Alignment Models. Computational Linguistics, 2003.

Word-Based Models and Word-Alignment

Natural
Language
Processing II

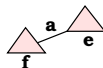
Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

a is a mapping between word positions.



- Δ_f and Δ_e are sequences of word positions.

$$\mathbf{e} = e_1^l = e_1 \dots e_l \text{ and } \mathbf{f} = f_1^m = f_1 \dots f_m$$

- A hidden word-alignment **a**:

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

- Each source position has a single link to a target position or to position zero

$$\mathbf{a} : \{pos_f\} \rightarrow (\{pos_e\} \cup \{0\})$$

- \mathbf{a}_i or $\mathbf{a}(i)$, i.e., word position in **e** with which **f_i** is aligned.

Word Alignment Example

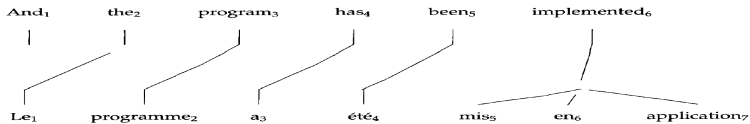
Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models



Word Alignment Example

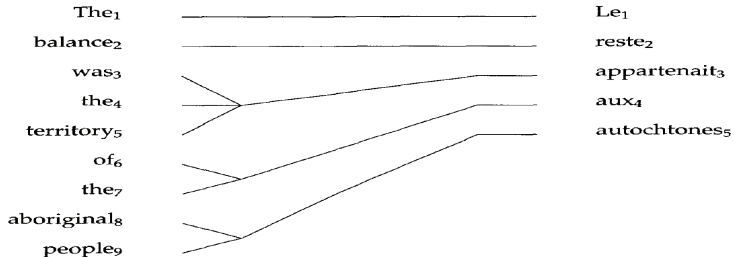
Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models



Word Alignment Example: Not covered in this setting

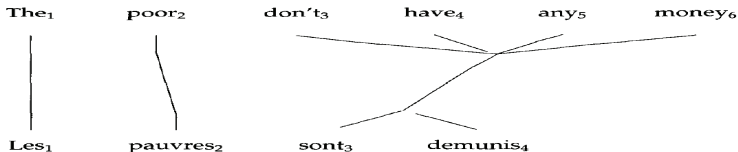
Natural
Language
Processing II

Dr Khalil
Sima'an

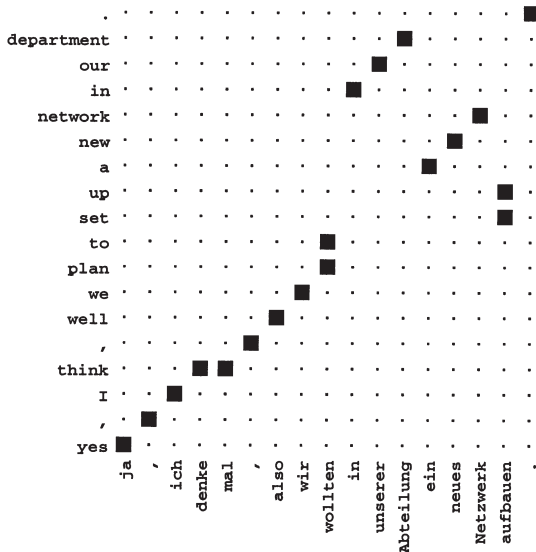
Challenges

This course

Word-Based
Models



Word Alignment Matrix Example



Translation model with word alignment

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

$$\arg \max_{\mathbf{e}} P(\mathbf{e} | \mathbf{f}) = \arg \max_{\mathbf{e}} P(\mathbf{e}) \times P(\mathbf{f} | \mathbf{e})$$

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

Questions

- How to parametrize the model?
How are \mathbf{e} , \mathbf{f} and \mathbf{a} composed from basic units?
- How to train the model?
How to acquire word alignment?
- How to translate with this model?
Decoding and computational issues (for second part)

Word-Alignment As Hidden Structure

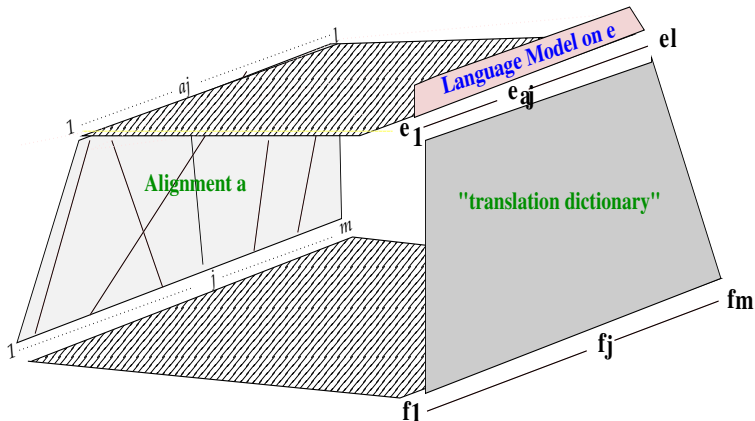
Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models



We need to decompose

- The alignment \mathbf{a} and the length m : $P(\mathbf{a} \mid \mathbf{e})$
- "Translation dictionary" $P(\mathbf{f} \mid \mathbf{e}, \mathbf{a})$

Word Alignment Models: General Scheme

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

Alignment of positions in **f** with positions in **e**:

$$\mathbf{a} = a_1^m = a_1 \dots a_m$$

Markov process over **a**

$$P(a_1^m, f_1^m \mid e_1^l) = P(m \mid \mathbf{e}) \times \prod_{j=1}^m P(a_j \mid a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) \times P(f_j \mid a_1^j, f_1^{j-1}, m, \mathbf{e})$$

In words: to generate alignment **a** and foreign sentence **f**

- 1 Choose a length m for **f**
- 2 Generate alignment a_j given the preceding alignments, words in **f**, m , and **e**
- 3 Generate word f_j conditioned on structure so far and **e**.

IBM models are obtained by simplifications of this formula.

IBM Model I

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

$$P(a_1^m, f_1^m \mid e_1 \dots e_l) = P(m \mid \mathbf{e}) \times \prod_{j=1}^m P(a_j \mid a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) \times P(f_j \mid a_1^j, f_1^{j-1}, m, \mathbf{e})$$

IBM Model I:

Length: $P(m \mid \mathbf{e}) \approx P(m \mid l) \approx \epsilon$ A **fixed** probability ϵ .

Align with **uniform** probability j with any a_j in \mathbf{e}_1^l or

NULL: $P(a_j \mid a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) \approx (l+1)^{-1}$

Note that a_j can be linked with l positions in \mathbf{e} or with NULL.

Lexicon: lexicon parameters $\pi_t(f \mid e)$

$$P(f_j \mid a_1^j, f_1^{j-1}, m, \mathbf{e}) \approx P(f_j \mid e_{a_j}) = \pi_t(f_j \mid e_{a_j})$$

Parameters: ϵ and $\{\pi_t(f \mid e) \mid \langle f, e \rangle \in \mathbf{C}\}$.

Sketch IBM Model I

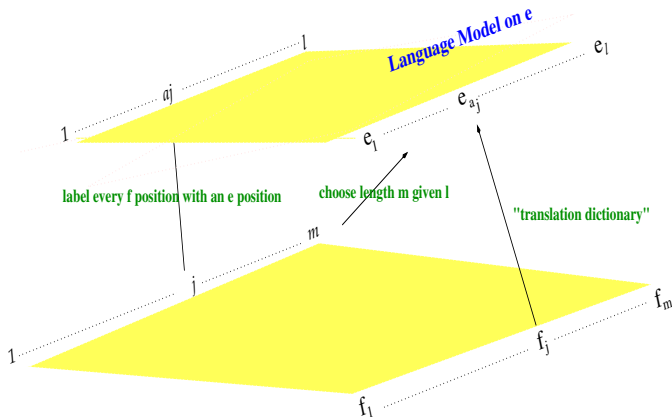
Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models



IBM Model I Explicit

IBM Model I altogether

$$\begin{aligned} P(\mathbf{f} | \mathbf{e}) &= \sum_{a_1^m} P(a_1^m, f_1^m | \mathbf{e}_1 \dots \mathbf{e}_l) \\ &= \frac{\epsilon}{(l+1)^m} \times \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m \pi_t(f_j | \mathbf{e}_{a_j}) \end{aligned}$$

Parameters: ϵ and $\{\pi_t(f | \mathbf{e}) \mid \langle f, \mathbf{e} \rangle \in \mathbf{C}\}$.

Fix ϵ , i.e., in practice put a uniform probability over a range $[1..m]$, for some natural number m .

Crucial step: Efficiency (trick A)

$$= \frac{\epsilon}{(l+1)^m} \times \prod_{j=1}^m \sum_{i=0}^l \pi_t(f_j | \mathbf{e}_i)$$

Questions regarding IBM Model I

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

- How to parametrize the model?
- How to train the model?
How to acquire word alignment?
- How to translate with this model?
Decoding and computational issues (for second part)

EM for Lexicon and Word Alignment Probs

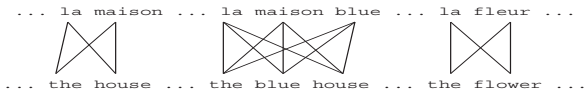
Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models



EM for Lexicon and Word Alignment Probs

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models



EM for Lexicon and Word Alignment Probs

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models



EM for Lexicon and Word Alignment Probs

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

... la maison ... la maison bleu ... la fleur ...
/ | | X | |
... the house ... the blue house ... the flower ...

IBM Model II

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

Extends IBM Model I at alignment probs:

$$P(a_1^m, f_1^m \mid e_1 \dots e_l) \approx \epsilon \times \prod_{j=1}^m \frac{P(a_j \mid a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) \times \pi_t(f_j \mid e_{a_j})}{P(a_j \mid a_1^{j-1}, f_1^{j-1}, m, \mathbf{e})}$$

IBM Model II: changes only one element in IBM Model I:

- IBM Model I does not take into account the position of words in both strings

$$P(a_j \mid a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) = P(a_j \mid j, l, m) := \pi_A(a_j \mid j, l, m)$$

Where $\pi_A(.|..)$ are parameters to be learned from data.

IBM Models III, IV and V concentrate on more complex alignments allowing, e.g., 1 – to – n (fertility)

IBM Model II Parameters

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

$$P(a_1^m, f_1^m \mid e_1 \dots e_l) \approx \epsilon \times \prod_{j=1}^m \pi_A(a_j \mid j, l, m) \times \pi_t(f_j \mid e_{a_j})$$

Parameters: $\{\pi_A(a_j \mid j, l, m)\}$ and $\{\pi_t(f_j \mid e_{a_j})\}$

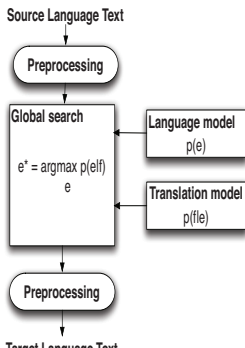
Estimation

Very similar to IBM Model I: EM estimation with the same complexity.

Translation Using EM Estimates

- Lexicon probability estimates: $\{\hat{\pi}_t(f_j | e_{a_j})\}$
- Alignment probabilities: $\{\hat{\pi}_A(a_j | j, m, l)\}$
- Translation Model + Language Model + Decoder

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



Viterbi Word-Alignment using EM estimates

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

After EM has stabilized on estimates

$$\{\hat{\pi}_t(f_j | e_{a_j})\} \quad \text{and} \quad \{\hat{\pi}_A(a_j | j, m, l)\}$$

For every $\langle \mathbf{f}, \mathbf{e} \rangle$ in \mathbf{C} apply the following

$$\arg \max_{a_1^m} P(a_1^m | f_1^m, e_1^l) =$$

$$\arg \max_{a_1^m} \prod_{j=1}^m \{\hat{\pi}_A(a_j | j, m, l)\} \pi_t(f_j | e_{a_j})$$

This can be done efficiently:

$$\text{solution } a_i = \arg \max_{j \in [0..l]} \hat{\pi}_A(a_j | j, m, l) \hat{\pi}_t(f_i | e_j)$$

HMM Alignment Model: General Form

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

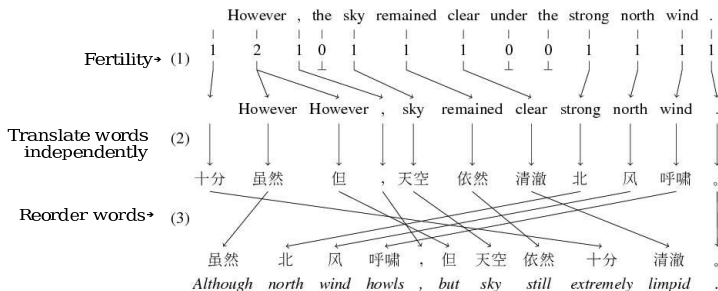
$$P(a_1^m, f_1^m \mid e_1 \dots e_l) \approx \epsilon \times \prod_{j=1}^m \frac{P(a_j \mid a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) \times \pi_t(f_j \mid e_{a_j})}{1}$$

- Words do not move independently of each other:
condition word movement on previous word movement

$$P(a_j \mid a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) \approx P(a_j \mid a_{j-1}, m)$$

IBM Model III (and IV): Example

- A hidden word-alignment \mathbf{a} : $P(\mathbf{f} | \mathbf{e}) = \sum_{\mathbf{a}} P(\mathbf{a}, \mathbf{f} | \mathbf{e})$



Estimate alignment + lexicon + reordering + fertility parameters.

Word-based Models (Och & Ney 2003)

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

Table 1

Overview of the alignment models.

| Model | Alignment model | Fertility model | E-step | Deficient |
|---------|-----------------|-----------------|---------------|-----------|
| Model 1 | uniform | no | exact | no |
| Model 2 | zero-order | no | exact | no |
| HMM | first-order | no | exact | no |
| Model 3 | zero-order | yes | approximative | yes |
| Model 4 | first-order | yes | approximative | yes |
| Model 5 | first-order | yes | approximative | no |
| Model 6 | first-order | yes | approximative | yes |

Word-Alignment As Hidden Structure: Sufficient?

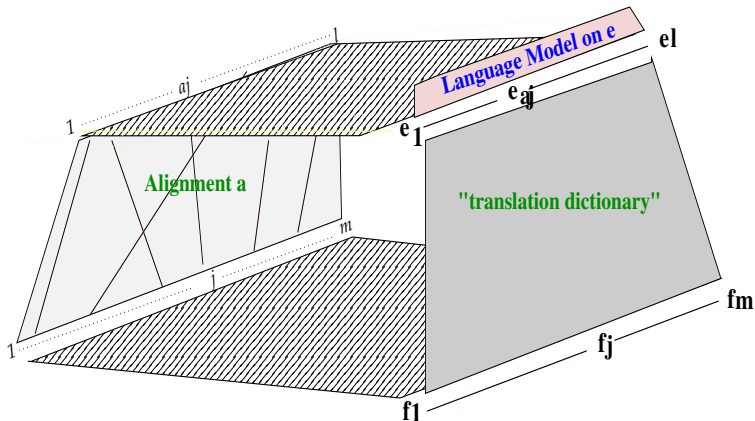
Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models



We assumed alignment between words and dictionary:

- Alignment \mathbf{a} and the length m : $P(\mathbf{a} \mid \mathbf{e})$
- Dictionary $P(\mathbf{f} \mid \mathbf{e}, \mathbf{a})$

Limitations of Word-based Models

Limitations of word-based translation:

- Many-to-one and many-to-many is common:
“Makes more difficult”/bemoeilijkt “Dat richtte (hen)
ten gronde”/”That destroyed (them)”
- Reordering takes place (often) by whole blocks.
Reordering individual words increases *ambiguity*.
“The (big heavy) cow/la vaca (pesada grande)”
- Translation works by “fixed expressions” (idiomatic).
Concatenating word-translations increases *ambiguity*.

Estimates of $P(\mathbf{f} \mid \mathbf{e})$ by word-based models are inaccurate.

Instead of words as basic events: multi-word events in
corpus.

NLP II topics

Natural
Language
Processing II

Dr Khalil
Sima'an

Challenges

This course

Word-Based
Models

We will cover literature (mostly articles) about

- Translation models: word-, phrase-, syntax-based
- Reordering models and synchronous grammars
- MT evaluation
- Paraphrasing and semantic models from parallel data
- Decoding algorithms