Machine Translation 1: Word alignments

Computational Linguistics

Alexander Koller

02 January 2018

slides contain material from mt-class.org

Google Translate



Google Translate

EL PAÍS	5	PORTADA INTERNACIO	NAL POI							
Google										Anmelden
Übersetze	r									×
Spanisch Deut	ch Englisch	Sprache erkennen	Ŧ	+	Deutsch	Englisch	Französisch	- Üb	ersetzen	

Joachim Löw, seleccionador de Alemania, ha anunciado este jueves la lista de los 30 jugadores preseleccionados para acudir al Mundial de Brasil, en la que destacan la ausencia del futuro portero del Barcelona Ter Stegen, y la incorporación de Sami Khedira, del Real Madrid. El medio, que siempre ha contado con la confianza del seleccionador, ya se ha recuperado de la rotura del ligamento cruzado y el interior de la rodilla derecha que se produjo durante un amistoso ante Italia en el mes de noviembre y que le ha mantenido apartado del terreno de juego durante siete meses.

de Sami Khedira, del Real Madrid. El medio, que siempre ha contado
 con la confianza del seleccionador, ya se ha recuperado de la rotura del ligamento cruzado y el interior de la rodilla derecha que se produjo durante un amistoso ante Italia en el mes de noviembre y que le ha mantenido apartado del terreno de juego durante siete meses.

Joachim Löw , Deutschland, am Donnerstag angekündigt, die Liste der 30 Spieler in die engere Wahl , die Weltmeisterschaft in Brasilien, die die Abwesenheit von zukünftigen Barcelona -Torhüter Ter Stegen, und der Einbau von Sami Khedira von Real Madrid gehören zu besuchen. Das Medium , das immer genossen hat, das Vertrauen des Trainers, und hat sich von der Kreuzbandriss und der Innenseite des rechten Knies , die bei einem Freundschaftsspiel gegen Italien im November aufgetreten erholt und er hat sich von der gehalten Feld für sieben Monate.

Automatic Translation



Google Pixel Buds, 2017





<u>Google's Pixel 2 event</u> in San Francisco on Wednesday had a lot of stuff to show off and most of it was more of the same: the next iteration of the flagship smartphone, new Home speakers and various ways of entwining them more deeply into your smart home, a new laptop that's basically a

Early History



One naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography.

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."

Warren Weaver to Norbert Wiener (1947)

Really Early History



Rosetta Stone, discovered 1799; deciphered 1822 (Champollion)

Really Early History

PIKAETEKAILITIKLEN

HNKATEIN

 $\frac{1}{2} = \frac{1}{2} + \frac{1$

Rosetta Stone, discovered 1799; deciphered 1822 (Champollion)

Classical view on translation



Types of MT systems

- What's it for?
 - fully automatic translation
 - support for human translators
- How does it work?
 - rule-based
 - statistical
 - neural
- Neural methods: see Language Technology II Here: elementary statistical methods

Corpora

- Learning translation models requires *parallel corpora*: text in one language with its translation in another.
- Popular parallel corpora:
 - Hansards (Canadian parliament): English/French
 - Europarl (European parliament): EU member languages
 - Literary texts with their translations (e.g. bible)

Step 1: Lexical Alignment

CHIER CALLER WALCAWER

習言が完全を行うに置い

EYTER

小VTI:S系上心生管型TEDT 1465-2 PPPP 学びは、1915年18月11日にあった。 1915年19月11日にあった。 1915年19月11日には、1915年には、1915年19月1日 1915年19月11日には、1915年19月1日 1915年19月11日

JLKANA F AN ANALESAN WIELEA HODALMH EMUAEIANAIAAU ENNALESAN WATARASY NYOLETALENAENAI ERAELENMMENLINA TAEN TOLETAEMAKIN NALEN ARENE I MENKAGA DEPERMITO ANTIN NEALTHNER TENENENENEN TOMENALA TAIL OF THIS HANK A TEINING MINING THIS INTO THIS IN THIS HANK A TEINING MINING THIS INTO TAILUTIES A TO A TAIL A

11311 0 Amiz 1 370 3400 الدين المراجع ا مراجع مراجع المراجع الم בווע איום וועמוב לי אוב 52-31422-42-212-212-212-2112 E.J. HILE.Y גיון ווזאיונייתן נציבוליות אינאו בין ואראבאראל אווואליאלווואליאליוויור באווויור באווואר Eranoal 102 mZ (1024-510)-3, m2 10. 23 1641 12 2 3 11/1. 15 1 1 / 1111 11C . 111 (Y - 1 / 2 - 1) 2 S 1 bile

Step 1: Lexical Alignment



EYTER

A NOXYPATAITAITON ATATAN NO AAMAH ANNKATAKIA TOIELAENILAITOTIE NAYTHIA ANNKATAKIA TOIELAENILAITOTIE NAYTHIA ANNKATAKIA TOIELAENILAITOITANTANCIITA

Lexical Translation

- We want to learn a model P(e | f):
 - e = "English" word (target language)
 - f = "French" word (original, foreign language)
- Gives a naive translation model for P(e | f).
 (Boldface e, f are English, Foreign sentences.)
- Linked to idea of word alignments.
 - alignments often independently useful (e.g. parse tree projection)

Word alignments



Alignment

• Alignments can be visualized by drawing links between two sentences, and they are represented as vectors of positions:



$$\mathbf{a} = (1, 2, 3, 4)$$

Reordering

• Words may be reordered during translation.





Word Dropping

• A source word may not be translated at all ("1" does not occur as a_i for any English position i)



$$a = (2, 3, 4)$$

Word Insertion

- Words may be inserted during translation
 - English "just" does not have an equivalent
 - record this by aligning with special NULL token at "position 0"



 $\mathbf{a} = (1, 2, 3, 0, 4)$

One-to-many Translation

• A Foreign word may translate into *more than one* English word.



$$\mathbf{a} = (1, 2, 3, 4, 4)$$

Many-to-one Translation

• *More than one* Foreign word may *not* translate into a single English word (can't represent this).





- Model P($\mathbf{a}, \mathbf{e} \mid \mathbf{f}, \mathbf{m}$) = P($\mathbf{e} \mid \mathbf{a}, \mathbf{f}, \mathbf{m}$) * P($\mathbf{a} \mid \mathbf{f}, \mathbf{m}$).
 - obtain $P(\mathbf{e} \mid \mathbf{f}, \mathbf{m})$ by marginalizing \mathbf{a} out \rightarrow translation
 - obtain $P(\mathbf{a} \mid \mathbf{f}, m)$ by marginalizing \mathbf{e} out \rightarrow compute alignments



- Model P($\mathbf{a}, \mathbf{e} \mid \mathbf{f}, \mathbf{m}$) = P($\mathbf{e} \mid \mathbf{a}, \mathbf{f}, \mathbf{m}$) * P($\mathbf{a} \mid \mathbf{f}, \mathbf{m}$).
 - obtain $P(\mathbf{e} \mid \mathbf{f}, \mathbf{m})$ by marginalizing \mathbf{a} out \rightarrow translation
 - obtain $P(\mathbf{a} \mid \mathbf{f}, m)$ by marginalizing \mathbf{e} out \rightarrow compute alignments



- Model P($\mathbf{a}, \mathbf{e} \mid \mathbf{f}, \mathbf{m}$) = P($\mathbf{e} \mid \mathbf{a}, \mathbf{f}, \mathbf{m}$) * P($\mathbf{a} \mid \mathbf{f}, \mathbf{m}$).
 - obtain $P(\mathbf{e} \mid \mathbf{f}, \mathbf{m})$ by marginalizing \mathbf{a} out \rightarrow translation
 - obtain $P(\mathbf{a} \mid \mathbf{f}, m)$ by marginalizing \mathbf{e} out \rightarrow compute alignments



- Model P($\mathbf{a}, \mathbf{e} \mid \mathbf{f}, \mathbf{m}$) = P($\mathbf{e} \mid \mathbf{a}, \mathbf{f}, \mathbf{m}$) * P($\mathbf{a} \mid \mathbf{f}, \mathbf{m}$).
 - obtain $P(\mathbf{e} \mid \mathbf{f}, \mathbf{m})$ by marginalizing \mathbf{a} out \rightarrow translation
 - obtain $P(\mathbf{a} \mid \mathbf{f}, m)$ by marginalizing \mathbf{e} out \rightarrow compute alignments



- Model P($\mathbf{a}, \mathbf{e} \mid \mathbf{f}, \mathbf{m}$) = P($\mathbf{e} \mid \mathbf{a}, \mathbf{f}, \mathbf{m}$) * P($\mathbf{a} \mid \mathbf{f}, \mathbf{m}$).
 - obtain $P(\mathbf{e} \mid \mathbf{f}, m)$ by marginalizing \mathbf{a} out \rightarrow translation
 - obtain $P(\mathbf{a} \mid \mathbf{f}, m)$ by marginalizing \mathbf{e} out \rightarrow compute alignments



- Model P($\mathbf{a}, \mathbf{e} \mid \mathbf{f}, \mathbf{m}$) = P($\mathbf{e} \mid \mathbf{a}, \mathbf{f}, \mathbf{m}$) * P($\mathbf{a} \mid \mathbf{f}, \mathbf{m}$).
 - obtain $P(\mathbf{e} \mid \mathbf{f}, \mathbf{m})$ by marginalizing \mathbf{a} out \rightarrow translation
 - obtain $P(\mathbf{a} \mid \mathbf{f}, m)$ by marginalizing \mathbf{e} out \rightarrow compute alignments



- Model P($\mathbf{a}, \mathbf{e} \mid \mathbf{f}, \mathbf{m}$) = P($\mathbf{e} \mid \mathbf{a}, \mathbf{f}, \mathbf{m}$) * P($\mathbf{a} \mid \mathbf{f}, \mathbf{m}$).
 - obtain $P(\mathbf{e} \mid \mathbf{f}, \mathbf{m})$ by marginalizing \mathbf{a} out \rightarrow translation
 - obtain $P(\mathbf{a} \mid \mathbf{f}, m)$ by marginalizing \mathbf{e} out \rightarrow compute alignments

IBM Model 1



- Simplifying assumptions:
 - The alignment decisions for the *m* English words are independent.
 - The alignment distribution for each a_i is uniform over all source words and NULL.
 - The English words are generated independently, conditioned only on their aligned Foreign words.

for each
$$i \in [1, 2, ..., m]$$

 $a_i \sim \text{Uniform}(0, 1, 2, ..., n)$
 $e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$

IBM Model 1

for each
$$i \in [1, 2, ..., m]$$

 $a_i \sim \text{Uniform}(0, 1, 2, ..., n)$
 $e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$

$$P(e_i, a_i \mid \mathbf{f}, m) = P(a_i \mid \mathbf{f}, m) \cdot P(e_i \mid a_i, \mathbf{f}, m) = \frac{1}{n+1} \cdot P(e_i \mid f_{a_i})$$

$$P(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^{m} P(e_i, a_i \mid \mathbf{f}, m) = \prod_{i=1}^{m} \frac{1}{n+1} \cdot P(e_i \mid f_{a_i})$$

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



Example

das				
e	t(e f)			
the	0.7			
that	0.15			
which	0.075			
who	0.05			
this	0.025			

Haus				
e	t(e f)			
house	0.8			
building	0.16			
home	0.02			
household	0.015			
shell	0.005			

ıst				
e	t(e f)			
is	0.8			
's	0.16			
exists	0.02			
has	0.015			
are	0.005			

kle	_	
e	t(e f)	[f)
small	0.4	(e
little	0.4	P
short	0.1	
minor	0.06	<u> </u>
petty	0.04	t(e



P(e, a | f, m) = 1/5 * P(Haus|house) * 1/5 * P(ist|is) * 1/5 * P(small|klein) = 1/125 * 0.8 * 0.8 * 0.4 = 0.002

a = (2, 3, 4)

Computing best alignments

- Assume that we know parameters P(e | f) and we are given e and f. What is alignment a that maximizes P(a | e, f)?
- Because of independence of $a_1, ..., a_m$, can choose best aligned word in **f** for each word in **e** separately.

$$a_i^* = \arg \max_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$
$$= \arg \max_{a_i=0}^n p(e_i \mid f_{a_i})$$

Training $p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^{m} \frac{1}{1+n} p(e_i \mid f_{a_i})$

- Parameters of our model: translation probs P(e | f) for any two words e and f.
- If we could observe alignments, then we could just do MLE: C(e aligned with f) / C(f)
- Because we usually only have raw parallel text, we need to use EM.
 - estimate counts from estimate of P
 - re-estimate P from estimated counts

EM: An Example

P(e f)	house	blue
maison	0.5	0.5
bleue	0.5	0.5



 $p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^{m} \frac{1}{1+n} p(e_i \mid f_{a_i})$

1. Compute P(e, a | f) for each alignment of each sentence pair.

 $P(\mathbf{e}_1, \mathbf{a}_{11} | \mathbf{f}_1) = 1/9 * 1/2 * 1/2 = 1/36$

$$P(\mathbf{e}_1, \mathbf{a}_{12} | \mathbf{f}_1) = 1/9 * 1/2 * 1/2 = 1/36$$

 $P(\mathbf{e}_2, \mathbf{a}_2 \mid \mathbf{f}_2) = 1/2 * 1/2 = 1/4$

(note: these are not really all alignments)

EM: An Example

P(e f)	house	blue
maison	0.5	0.5
bleue	0.5	0.5



2. Normalize $P(\mathbf{e}, \mathbf{a} \mid \mathbf{f})$ to yield $P(\mathbf{a} \mid \mathbf{e}, \mathbf{f})$. $P(\mathbf{a} \mid \mathbf{e}, \mathbf{f}) = \frac{P(\mathbf{a}, \mathbf{e} \mid \mathbf{f})}{P(\mathbf{e} \mid \mathbf{f})} = \frac{P(\mathbf{a}, \mathbf{e} \mid \mathbf{f})}{\sum_{\mathbf{a}'} P(\mathbf{a}', \mathbf{e} \mid \mathbf{f})}$

 $P(\mathbf{a}_{11} \mid \mathbf{e}_1, \mathbf{f}_1) = 1/2$

3. collect expected counts

	tc	house	blue
$P(\mathbf{a}_{12} \mid \mathbf{e}_1, \mathbf{f}_1) = 1/2 \longrightarrow$	maison	3/2	1/2
	bleue	1/2	1/2

 $P(\mathbf{a}_2 \mid \mathbf{e}_2, \mathbf{f}_2) = 1$

EM: An Example

4. Normalize expected counts C(e, f) by total expected counts C(f) to obtain revised translation probs P(e | f).

expected counts

revised translation probs

tc	house	blue	P(e f)	house	blue
maison	3/2	1/2	 maison	3/4	1/4
bleue	1/2	1/2	bleue	1/2	1/2

EM: Round Two

P(e f)	house	blue
maison	3/4	1/4
bleue	1/2	1/2



$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^{m} \frac{1}{1+n} p(e_i \mid f_{a_i})$$

 $P(e_1, a_{11} | f_1) = 1/9 * 1/4 * 1/2 = 1/72$

$$P(e_1, a_{12} | f_1) = 1/9 * 3/4 * 1/2 = 3/72$$

$$P(e_2, a_2 | f_2) = 1/2 * 3/4 = 3/8$$

EM: Round Two

P(e f)	house	blue
maison	3/4	1/4
bleue	1/2	1/2





EM: Round Two

exp	ected cou	ints	revised	translatio	on probs
tc	house	blue	P(e f)	house	blue
maison	7/4	1/4	 maison	7/8	1/8
bleue	1/4	3/4	bleue	1/4	3/4

After many iterations:

P(e f)	house	blue
maison	≈ 1	≈ 0
bleue	≈ 0	≈ 1

Efficient computation

• Computation of P(**a** | **e**, **f**) in E-step is tricky:

$$P(a_i = j \mid \mathbf{e}, \mathbf{f}) = \frac{P(a_i = j, \mathbf{e} \mid \mathbf{f})}{P(\mathbf{e} \mid \mathbf{f})} = \frac{\sum_{\mathbf{a}:a_i = j} \prod_{i'=1}^m P(e_{i'} \mid f_{a_{i'}})}{\sum_{\mathbf{a}} \prod_{i'=1}^m P(e_{i'} \mid f_{a_{i'}})}$$

- Summation over **a** is exponential in sentence length.
- By clever use of law of distributivity, can rewrite this term so it can be computed in quadratic time.
 See Lopez tutorial on website. (Note flipped e and f.)

Extensions

- IBM Model 2: P(**a**) not uniform, but implements *reordering model* that prefers alignments in which words stay close to their original position.
- Model 3: adds *fertility model* that predicts the number of English words to which a given f will be aligned. Can't do EM, approximate with sampling.
- Models 4-5: more complicated reordering models.
- Implemented in GIZA++ and successor tools.

Conclusion

- Machine translation: one of the most useful and most challenging disciplines of NLP.
- Today: word alignments.
 - IBM Model 1
 - computing best alignments
 - EM training
 - advanced models
- Next time: let's actually translate something.