

Statistical Machine Translation

presentation: Adam Lopez
slides: Chris Callison-Burch

Various approaches

- Word-for-word translation
- Syntactic transfer
- Interlingual approaches
- Controlled language
- Example-based translation
- Statistical translation

Advantages of SMT

- Data driven
- Language independent
- No need for staff of linguists of language experts
- Can prototype a new system quickly and at a very low cost

Statistical machine translation

- Find most probable English sentence given a foreign language sentence
- Automatically align words and phrases within sentence pairs in a parallel corpus
- Probabilities are determined automatically by training a statistical model using the parallel corpus

Parallel corpus

what is more , the relevant cost dynamic is completely under control .	im übrigen ist die diesbezügliche kostenentwicklung völlig unter kontrolle .
sooner or later we will have to be sufficiently progressive in terms of own resources as a basis for this fair tax system .	früher oder später müssen wir die notwendige progressivität der eigenmittel als grundlage dieses gerechten steuersystems zur sprache bringen .
we plan to submit the first accession partnership in the autumn of this year .	wir planen , die erste beitrittspartnerschaft im herbst dieses jahres vorzulegen .
it is a question of equality and solidarity .	hier geht es um gleichberechtigung und solidarität .
the recommendation for the year 1999 has been formulated at a time of favourable developments and optimistic prospects for the european economy .	die empfehlung für das jahr 1999 wurde vor dem hintergrund günstiger entwicklungen und einer für den kurs der europäischen wirtschaft positiven perspektive abgegeben .
that does not , however , detract from the deep appreciation which we have for this report .	im übrigen tut das unserer hohen wertschätzung für den vorliegenden bericht keinen abbruch .

Probabilities

- Find most probable English sentence given a foreign language sentence

$$p(e|f)$$

Probabilities

- Find most probable English sentence given a foreign language sentence

$$p(e|f)$$
$$\hat{e} = \arg \max_e p(e|f)$$

Probabilities

- Find most probable English sentence given a foreign language sentence

$$p(e|f)$$
$$\hat{e} = \arg \max_e p(e|f)$$
$$p(e|f) = \frac{p(e)p(f|e)}{p(f)}$$

Probabilities

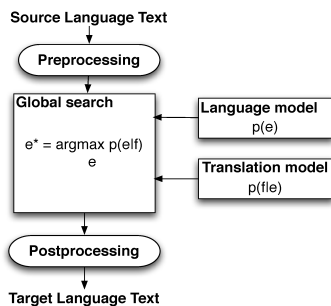
- Find most probable English sentence given a foreign language sentence

$$p(e|f)$$
$$\hat{e} = \arg \max_e p(e|f)$$
$$p(e|f) = \frac{p(e)p(f|e)}{p(f)}$$
$$\hat{e} = \arg \max_e p(e)p(f|e)$$

What the probabilities represent

- $p(e)$ is the "Language model"
 - Assigns a higher probability to fluent / grammatical sentences
 - Estimated using monolingual corpora
- $p(f|e)$ is the "Translation model"
 - Assigns higher probability to sentences that have corresponding meaning
 - Estimated using bilingual corpora

For people who don't like equations



Language Model

- Component that tries to ensure that words come in the right order
- Some notion of grammaticality
- Standardly calculated with a trigram language model, as in speech recognition
- Could be calculated with a statistical grammar such as a PCFG

Trigram language model

- $p(\text{l like bungee jumping off high bridges}) =$

Trigram language model

- $p(\text{l like bungee jumping off high bridges}) =$
 $p(\text{l} \mid \langle s \rangle \langle s \rangle)^*$

Trigram language model

- $p(\text{l like bungee jumping off high bridges}) =$
 $p(\text{l} \mid \langle s \rangle \langle s \rangle)^*$
 $p(\text{like} \mid \langle s \rangle \text{l})^*$

Trigram language model

- $p(\text{l like bungee jumping off high bridges}) =$
 $p(\text{l} \mid \langle s \rangle \langle s \rangle)^*$
 $p(\text{like} \mid \langle s \rangle \text{l})^*$
 $p(\text{bungee} \mid \text{l like})^*$

Trigram language model

- $p(\text{l like bungee jumping off high bridges}) =$
 $p(\text{l} \mid \langle s \rangle \langle s \rangle)^*$
 $p(\text{like} \mid \langle s \rangle \text{l})^*$
 $p(\text{bungee} \mid \text{l like})^*$
 $p(\text{jumping} \mid \text{like bungee})^*$

Trigram language model

- $p(\text{l like bungee jumping off high bridges}) =$
 $p(\text{l} \mid \langle s \rangle \langle s \rangle)^*$
 $p(\text{like} \mid \langle s \rangle \text{l})^*$
 $p(\text{bungee} \mid \text{l like})^*$
 $p(\text{jumping} \mid \text{like bungee})^*$
 $p(\text{off} \mid \text{bungee jumping})^*$

Trigram language model

- $p(\text{l like bungee jumping off high bridges}) =$
 $p(\text{l} \mid \langle s \rangle \langle s \rangle)^*$
 $p(\text{like} \mid \langle s \rangle \text{l})^*$
 $p(\text{bungee} \mid \text{l like})^*$
 $p(\text{jumping} \mid \text{like bungee})^*$
 $p(\text{off} \mid \text{bungee jumping})^*$
 $p(\text{high} \mid \text{jumping off})^*$

Trigram language model

- $p(\text{l like bungee jumping off high bridges}) =$
 $p(\text{l} \mid \langle s \rangle \langle s \rangle)^*$
 $p(\text{like} \mid \langle s \rangle \text{l})^*$
 $p(\text{bungee} \mid \text{l like})^*$
 $p(\text{jumping} \mid \text{like bungee})^*$
 $p(\text{off} \mid \text{bungee jumping})^*$
 $p(\text{high} \mid \text{jumping off})^*$
 $p(\text{bridges} \mid \text{off high})^*$

Trigram language model

- $p(\text{l like bungee jumping off high bridges}) =$
 $p(\text{l} \mid \langle s \rangle \langle s \rangle)^*$
 $p(\text{like} \mid \langle s \rangle \text{l})^*$
 $p(\text{bungee} \mid \text{l like})^*$
 $p(\text{jumping} \mid \text{like bungee})^*$
 $p(\text{off} \mid \text{bungee jumping})^*$
 $p(\text{high} \mid \text{jumping off})^*$
 $p(\text{bridges} \mid \text{off high})^*$
 $p(\langle /s \rangle \mid \text{high bridges})^*$

Trigram language model

- $p(\text{l like bungee jumping off high bridges}) =$
 $p(\text{l} \mid \langle s \rangle \langle s \rangle)^*$
 $p(\text{like} \mid \langle s \rangle \text{l})^*$
 $p(\text{bungee} \mid \text{l like})^*$
 $p(\text{jumping} \mid \text{like bungee})^*$
 $p(\text{off} \mid \text{bungee jumping})^*$
 $p(\text{high} \mid \text{jumping off})^*$
 $p(\text{bridges} \mid \text{off high})^*$
 $p(\langle /s \rangle \mid \text{high bridges})^*$
 $p(\langle /s \rangle \mid \text{bridges} \langle /s \rangle)^*$

Calculating Language Model Probabilities

- Unigram probabilities

$$p(w_1) = \frac{\text{count}(w_1)}{\text{total words observed}}$$

Calculating Language Model Probabilities

- Bigram probabilities

$$p(w_2 \mid w_1) = \frac{\text{count}(w_1 w_2)}{\text{count}(w_1)}$$

Calculating Language Model Probabilities

- Trigram probabilities

$$p(w_3|w_1w_2) = \frac{\text{count}(w_1w_2w_3)}{\text{count}(w_1w_2)}$$

Calculating Language Model Probabilities

- Can take this to increasingly long sequences of n-grams
- As we get longer sequences it's less likely that we'll have ever observed them

Backing off

- Sparse counts are a big problem
- If we haven't observed a sequence of words then the count = 0
- Because we're multiplying the n-gram probabilities to get the probability of a sentence the whole probability = 0

Backing off

$$.8 * p(w_3|w_1w_2) + .15 * p(w_3|w_2) + .049 * p(w_3) + .001$$

- Avoids zero probs

Translation model

- $p(f|e)$... the probability of some foreign language string given a hypothesis English translation
- f = Ces gens ont grandi, vécu et oeuvré des dizaines d'années dans le domaine agricole.
- e = *Those people have grown up, lived and worked many years in a farming district.*
- e = *I like bungee jumping off high bridges.*

Translation model

- How do we assign values to $p(f|e)$?

$$p(f|e) = \frac{\text{count}(f, e)}{\text{count}(e)}$$

- Impossible because sentences are novel, so we'd never have enough data to find values for all sentences.

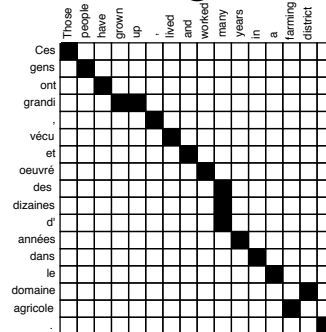
Translation model

- Decompose the sentences into smaller chunks, like in language modeling

$$p(f|e) = \sum_a p(a, f|e)$$

- Introduce another variable a that represents alignments between the individual words in the sentence pair

Word alignment



Alignment probabilities

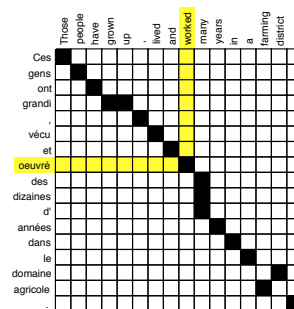
- So we can calculate translation probabilities by way of these alignment probabilities

$$p(f|e) = \sum_a p(a, f|e)$$

- Now we need to define $p(a, f|e)$

$$p(a, f|e) = \prod_{j=1}^m t(f_j|e_i)$$

Calculating $t(f_j|e_i)$



- Counting! I told you probabilities were easy!

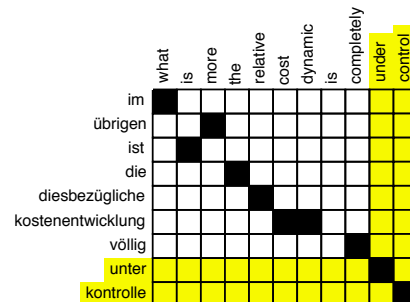
$$= \frac{\text{count}(f_j, e_i)}{\text{count}(e_i)}$$

- worked... fonctionné, travaillé, marché, oeuvré
- 100 times total 13 with this f. 13%

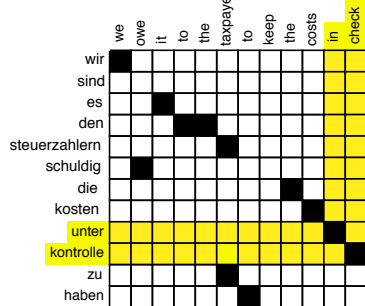
Calculating $t(f_j|e_i)$

- Unfortunately we don't have word aligned data, so we can't do this directly.
- OK, so it's not quite as easy as I said.
- Tomorrow's lecture will describe how word alignments are obtained using Expectation Maximization.

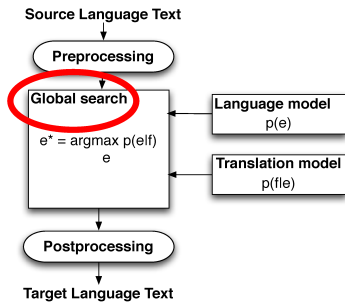
Phrase Translation Probabilities



Phrase Translation Probabilities



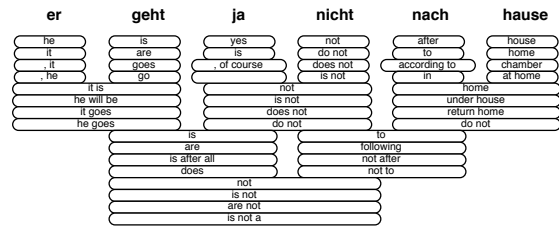
"Diagram Number 1"



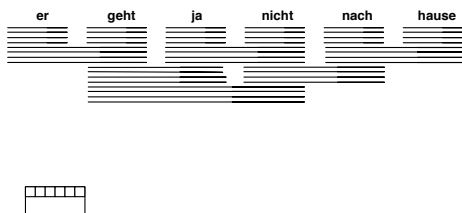
The Search Process AKA "Decoding"

- Look up all translations of every source phrase
- Recombine the target language phrases that maximizes the translation model probability * the language model probability
- This search over all possible combinations can get very large so we need to find ways of limiting the search space

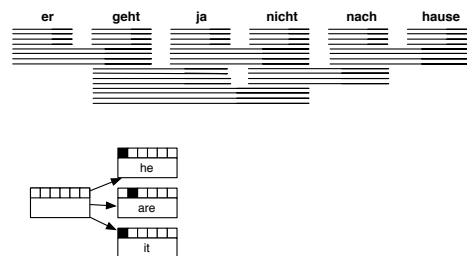
Translation Options



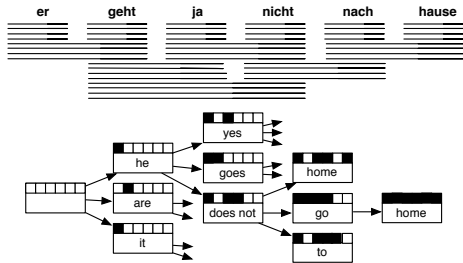
Search



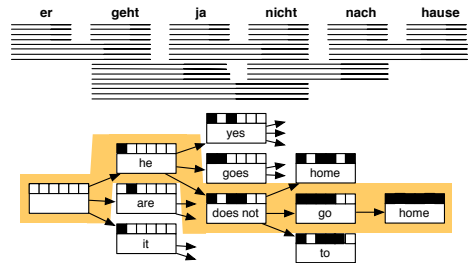
Search



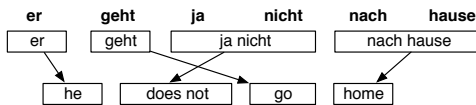
Search



Search



Best Translation



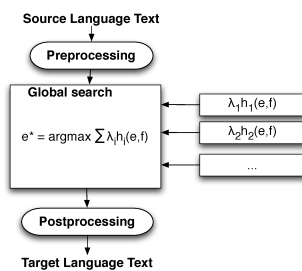
The Search Space

- In the end the item which covers all of the source words and which has the highest probability wins!
- That's our best translation

$$\hat{e} = \arg \max_e p(e)p(f|e)$$
- And there was much rejoicing!

Alternative models

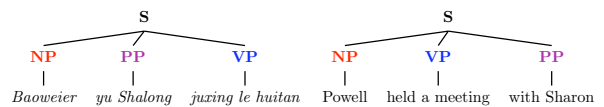
Linear models



Alternative models

Tree-based models

S → **NP**⁽¹⁾ **PP**⁽²⁾ **VP**⁽³⁾, **NP**⁽¹⁾ **VP**⁽³⁾ **PP**⁽²⁾
NP → Baoweier, Powell
PP → yu Shalong, with Sharon
VP → juzing le huitan, held a meeting



Wrap-up: SMT is data driven

- Learns translations of words and phrases from parallel corpora
- Associate probabilities with translations empirically by counting co-occurrences in the data
- Estimates of probabilities get more accurate as size of the data increases

Wrap-up: SMT is language independent

- Can be applied to any language pairs that we have a parallel corpus for
- The only linguistic thing that we need to know is how to split into sentences, words
- Don't need linguists and language experts to hand craft rules because it's all derived from the data

Wrap-up: SMT is cheap and quick to produce

- Low overhead since we aren't employing anyone
- Computers do all the heavy lifting / statistical analysis of the data for us
- Can build a system in hours or days rather than months or years

More Information

- <http://www.statmt.org> - papers, tutorials, etc.
- Statistical Machine Translation. In *ACM Computing Surveys* 40(3), Aug 2008.
At <http://homepages.inf.ed.ac.uk/alopez>
BibTeX at <http://github.com/alopez/smtbib>