



Human Language Technology: Applications to Information Access

Lesson 6: Translation Models

November 3, 2016

EPFL Doctoral Course EE-724 Andrei Popescu-Belis Idiap Research Institute, Martigny

Generative modeling for statistical MT *(reminder)*

- Noisy channel model
 - generation of French sentence *f* considered as a transmission of English sentence *e* into French
 - goal: given f, what is the most likely e?
 - this will actually produce a translation *e* of *f*
- Principle (using Bayes' theorem)
 - learn English language model: P(e)
 - learn (reverse) translation model: P(f|e)
 - decode source sentence: find* most likely e given f

 $\underset{e \in English}{\operatorname{argmax}} P(e \mid f) = \underset{e \in English}{\operatorname{argmax}} \left(P(f \mid e) \cdot P(e) \right)$

*find: easier said than done!

Translation modeling for MT

- Main question: how do we model *P*(*e*|*f*)?
 - given a foreign (French) sentence f (which we don't understand) and an English translation candidate e,
 - how do we compute (estimate) the probability that e is really a translation of f?
- Plan of the lesson
 - IBM Model 1
 - expectation-maximization (EM) algorithm
 - quick overview of IBM Models 2-5
 - phrase-based translation models
 - sentence and word alignment algorithms

Learning a translation model

- Estimating parameters of generative model from data
 - be able to estimate P(e|f) for any e, f
 - *P*(*e*|*f*) or *P*(*f*|*e*) are equivalent because the method is the same
- Data = parallel corpus: pairs of source and target (human-translated) sentences
- IBM Models 1 and 2
 - IBM Model 1: alignments, EM algorithm, implementation
 - IBM Models 2-5: incremental complexification
- Evaluation of TMs: perplexity

Alignment function: a

- Word mapping from a French sentence to an English one
 - DEFINITION: *a*(*i*) is the position of the French word that is translated into the English word which is in position *i*
- Example of a correct alignment
 - <u>French</u>: il₁ a_2 gravi₃ a_4 montagne₅ enneigée₆ \rightarrow <u>English</u>: he₁ climbed₂ the₃ snowy₄ mountain₅ *Alignment*: a(1)=1, a(2)=3, a(3)=4, a(4)=6, a(5)=5
- Remarks
 - not injective; not surjective; needs <u>NULL</u> token (0) in French for English words that are not translations of a French word
- Another example
 - <u>French</u>: NULL₀ je₁ ne₂ veux₃ pas₄ me₅ taire₆ \rightarrow <u>English</u>: I₁ do₂ not₃ want₄ to₅ shut₆ up₇ *Alignment*: a(1)=1, a(2)=0, a(3)=2, a(4)=3, a(5)=0, a(6)=6, a(7)=6

Example of difficult word (and even sentence) alignments (Jurafsky & Martin 1999, p. 473)



10

Learning to estimate P(e|f)

- Intuitive idea
 - estimate P(e|f) using the probabilities that
 a word in e gets translated into a word in f
 - these probabilities cannot be computed directly
 - unless we have word-aligned training data (but we don't, for now)
- Better idea
 - introduce the alignment functions as latent variables to be estimated as well
 - word translations and alignments are related:

 $P(e|f) = \sum_{a} P(e, a|f) \text{ but also } P(a|e, f) = P(e, a|f) / P(e|f)$ (applying the chain rule)

Use of the EM algorithm

(Dempster, Laird, Rubin 1977)

- Expectation maximization
 - iterative method to estimate model parameters together with latent variables which have reciprocal dependencies
 - applicable when
 - equations for parameters and variables cannot be solved directly
 - derivative of likelihood function (of parameters) cannot be obtained
- Iterate the E and M steps
 - Expectation: using current estimate of parameters, compute a likelihood function over latent variables (and parameters)
 - Maximization: re-compute that parameters so that they maximize the likelihood that was found
- Initialization: e.g. with uniform probabilities for parameters

EM applied to IBM Model 1

• Reminder:

 $P(e|f) = \sum_{a} P(e, a|f)$ and P(a|e, f) = P(e, a|f) / P(e|f)

- We can make *P(e, a | f)* more explicit
 - assume it depends only on word translation probabilities $P(e_i|f_i)$, noted $t(.|.) \rightarrow \text{``Model 1''}$
 - assume these are independent
 - note E and F lengths in words of sentences e and f

- then: $P(e, a | f) = (\varepsilon \prod_{j=1..E} t(e_j | f_{a(j)})) / (F + 1)^E$

• ε is a normalization factor

Making P(e|f) more explicit (E step)

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

$$\begin{split} &= \sum_{a(1)=0..F} \dots \sum_{a(E)=0..F} P(e, a \mid f) \\ &= \sum_{a(1)=0..F} \dots \sum_{a(E)=0..F} (\varepsilon/(F+1)^E) \prod_{j=1..E} t(e_j \mid f_{a(j)}) \\ &= (\varepsilon/(F+1)^E) \sum_{a(1)=0..F} \dots \sum_{a(E)=0..F} \prod_{j=1..E} t(e_j \mid f_{a(j)}) \\ &= (\varepsilon/(F+1)^E) \prod_{j=1..E} \sum_{i=0..F} t(e_j \mid f_i) \end{split}$$

Also: $P(a | e, f) = \prod_{j=1..E} (t(e_j | f_{a(j)}) / \sum_{i=0..F} t(e_j | f_i))$

EM applied to IBM Model 1 (again)

- Better use the word-specific t(.|.) values for EM rather than the full-sentence probability P(e|f)
 - they can be estimated directly using counts
 - they suffice to compute P(e|f) given the above formulae
 - much fewer than P(e|f) values when training on large corpora
- Making the EM steps more explicit
- 1. Start with uniform $t(e_i | f_i)$ values
 - a. <u>Expectation step</u>: compute all *P*(*a*|*e*, *f*) using all current *t*(.|.) values
 - b. <u>Maximization step</u> (of likelihood): update the t(.|.) values (improve estimates) using counts and P(a|e, f) values
- 2. Iterate (1) and (2) with a guarantee of convergence!

Estimating *t*(.|.) values using *counts*

- New notation: say *we* and *wf* are English and French *words*
- Define *count(we, wf, e, f)* as the number of times *wf* was translated into *we* according to all alignments of *e* and *f* weighed by their probability → M step:

 $t(we | wf) = (\sum_{(e,f)} count(we, wf, e, f)) / (\sum_{we} \sum_{(e,f)} count(we, wf, e, f))$ (the denominator is the normalization factor)

• Collecting counts ($\delta(..., ...$) is Kronecker's function) $\rightarrow \mathbf{E}$ step: $count(we, wf, e, f) = \sum_{a} P(a \mid e, f) \sum_{j=1..E} \delta(we, e_j) \delta(wf, f_{a(j)})$ which are expressed in terms of t(.|.) using last line of slide 14: $= (t(we \mid wf) \sum_{j=1..E} \delta(we, e_j) \sum_{i=1..F} \delta(wf, f_i)) / (\sum_{i=0..F} t(we \mid f_i))$

A simple implementation [from Koehn 2010, page 91]

```
initialize t(we|wf) uniformly
while <not converged>
    for all we, wf do
         count(we, wf)=0
         total(wf)=0
    end for
    for₁ all e, f do
         for, all we \in e do
              subtotal(we) = 0
              for<sub>3</sub> all wf \in f do
                   subtotal(we) += t(we | wf)
              end for<sub>3</sub>
         end for<sub>2</sub>
```

```
for<sub>4</sub> all we \in e do
                   for<sub>5</sub> all wf \in f do
                          count(we, wf) += t(we|wf)/
                                       subtotal(we)
                          total(wf) += t(we|wf)/
                                       subtotal(we)
                   end for<sub>5</sub>
             end for<sub>4</sub>
      end for<sub>1</sub>
      for_6 all wf do
            for<sub>7</sub> all we do
                   t(we|wf) = count(we|wf)/
                          total(wf)
             end for<sub>7</sub>
      end for<sub>6</sub>
end while
                                                      17
```

When do we stop?

• Quality of a translation model: perplexity

- if s are sentences from a given corpus $\log_2 PP = -\sum_s \log_2 P(e_s | f_s)$

- Iterate until PP stops decreasing (<converge>)
- IBM Model 1
 - EM guarantees that Model 1 will converge towards a global minimum of perplexity

IBM Model 2

- In Model 1, alignment probabilities were factored out from the formula for P(e, a | f) and from EM
 We had: P(e, a | f) = (ε / (F + 1)^E) Π_{j=1..E} t(e_j | f_{a(j)}))
 (the big cat → le gros chat) and
 (the big cat → chat gros le) have the same probability!
- Alignment probability distribution $P_a(i|j, E, F)$
 - probability that it is the French word in position *i* that corresponds to the English word in position *j*

Therefore, $P(e, a | f) = \varepsilon' \prod_{j=1..E} t(e_j | f_{a(j)}) P_a(a(j) | j, E, F)$

IBM Model 2 (continued)

• Equations are transformed similarly to Model 1

 $P(e|f) = \varepsilon \prod_{j=1..E} \sum_{i=0..F} t(e_j|f_i) P_a(i|j, E, F)$

- Computation of (fractional) counts for word translations $count(we, wf, e, f) = \sum_{j=1..E} \sum_{i=1..F} (t(we | wf) P_a(i | j, E, F) \delta(we, e_j) \delta(wf, f_i) / (\sum_{k=0} E t(we | f_k) P_a(k | j, E, F))$
- Computation of counts for alignments

 $count_{a}(i, j, e, f) = t(e_{j}|f_{j}) P_{a}(i|j, E, F) / (\sum_{k=0..F} t(e_{j}|f_{k}) P_{a}(k|j, E, F))$

• Training algorithm similar to Model 1 + start with initial values of $t(e_j|f_j)$ obtained by some iterations of Model 1

IBM Models 1-3

- IBM Model 1: lexical translation
- IBM Model 2: added absolute alignment model
- IBM Model 3: added fertility and absolute distortion models
 - probability that one English word generates several French ones
 - insertion of NULL token with a fixed probability after each word
 - probability of distortion instead of alignment $P_d(j|i, E, F)$
 - exhaustive count collection no longer possible → sample the alignment space to find optimum by hill climbing, from Model 2

IBM Models 4 and 5

- IBM Model 4: adds relative distortion model
 - introduces notion of "cept" (sort of phrase)
 - models distortion based on the position of a word in a cept (e.g. initial or not), possibly also on word class (POS-based or empirically)
- IBM Model 5: deal with a deficiency
 - deficiency = alignments can have "superposed" words
 - avoid spreading probability over such alignments
 - less commonly used than Models 1-4

Using IBM Models for word alignment: for instance in the GIZA++ tool

- Results of IBM Models after EM algorithm = probabilities for lexical translation and alignment
- Can be used to determine the most probable word alignment for each sentence pair (= the "Viterbi alignment")
 - Model 1, for each word e_i select the most likely word f_i (using $t(e_i|f_i)$)
 - Model 2, same but maximize $t(e_i|f_j) P_a(j|i, E, F)$
 - Models 3-5, no closed form expression: start with Model 2, then use heuristics
- Many other methods exist for *word alignment*
 - generative: train HMMs on linking probabilities, then use Viterbi decoding or another dynamic programming method
 - discriminative: structured prediction, feature functions, etc.
 - still, for phrase-based translation, IBM Models 1-4 perform well

Improving word alignments: towards phrase-based translation models

- For a given translation direction, the IBM models can find one-to-one alignments, multiple-to-one, one-to-zero, but *never one-to-multiple*
 - still, for a correct alignment, we might need both
 - {Paul} {was waiting} {inside} \leftrightarrow {Paul} {attendait} {à l'intérieur}
- Solution: symmetrization, by running the algorithm in both directions
 - consider the intersection of the two sets of alignment points, or their union, or enrich intersection with some points of the union, etc.
- IBM Models are no longer used as translation models (to estimate P(e|f)) but only to produce (1) probabilities of word translation, and (2) word alignments that will help learning phrase-based TMs
 - lower IBM models are used as steps to learn higher ones $(1 \rightarrow 4)$

Phrase-based translation models

- The goal remains the same
 - given a foreign (French) sentence f, look for the English sentence e which maximizes P(e|f), with $\operatorname{argmax}_{e} P(e|f) = \operatorname{argmax}_{e} P_{\mathsf{TM}}(f|e) P_{\mathsf{LM}}(e)$
- But take a different approach to compute P_{TM}: consider each sentence *e* and *f* as made of <u>phrases</u>

 $e = \{\underline{e}_1, ..., \underline{e}_i, ..., \underline{e}_M\}$ and $f = \{\underline{f}_1, ..., \underline{f}_j, ..., \underline{f}_N\}$

- phrases are non-empty ordered sets of contiguous words
- phrases cover entirely the sentence

Note on the word `phrase'

- originally, a linguistic notion (noun phrase, verb phrase)
- here, just a set of words, no linguistic motivation
- in French: `groupe de mots', not `phrase'

Phrase-based translation probability

$$P_{\mathsf{TM}}(f|e) = \prod_{i=1..M} P(f_i|\underline{e}_i) \cdot d(\mathsf{START}(f_i) - \mathsf{END}(f_{i-1}) - 1)$$

 $P(f_i | \underline{e}_i)$ is the probability that \underline{e}_i is translated into f_i

d is a "distance-based reordering model"

START(f_i) is the position of the first word of phrase f_i END(f_i) is the position of the last word of phrase f_i

- e.g., if phrases \underline{e}_{i-1} and \underline{e}_i are translated in sequence by phrases \underline{f}_{i-1} and \underline{f}_i , then START(\underline{f}_i) = END(\underline{f}_{i-1})+1, and we have d(0)
- a simple and efficient function: $d(x) = \alpha^{|x|}$ (adjust α if needed)

How do we estimate $P(\underline{f}|\underline{e})$ i.e. the probabilities of phrase translations?

- Consider aligned sentence pairs

 (can be computed, e.g. with the Gale and Church algorithm)
- 2. Perform word alignment for each pair
 - e.g. with the IBM Models
- 3. For each sentence pair, extract all phrase pairs that are "consistent" with the word alignment
 - possibly with some filtering, e.g. by length
- 4. Estimate *P* from counts of phrase pairs

Phrases consistent with alignment points



Formal definitions

- Definition of "consistent with an alignment"
 - a phrase pair (f, e) is consistent with an alignment A iff
 - all words f_i from f that have alignment points have them with words in \underline{e} (not outside \underline{e}): $\forall e_i \in \underline{e}$, $(e_i, f_i) \in A \implies f_i \in \underline{f}$
 - and vice-versa
 - and the pair includes at least one alignment point
- Definition of counts
 - extract all phrase pairs from the corpus
 - how many times each pair was extracted? count(<u>e</u>, <u>f</u>)
 - then estimate: $P(\underline{f}|\underline{e}) = \text{count}(\underline{e}, \underline{f}) / \sum_{\underline{f}_i} \text{count}(\underline{e}, \underline{f}_i)$

Remarks

- Number of extracted phrases: quadratic in the number of words, for each sentence pair
 - approximation: $|e| \approx |f|$ and "correct" alignments
 - note: unaligned words generate many more phrases
 - advantage: a lot of phrases to choose from
 - disadvantage: large memory/disk requirements
 - solution: remove long phrases and/or those seen only once
- Comparison of phrase-based and IBM models
 - phrase-based have simpler formulation
 - but they require word alignments
 - so IBM models (or others) are still needed to find the Viterbi alignment for each sentence pair

Towards log-linear "translation models"

- Translate *f* = find *e* which maximizes the product of three terms
 - probabilities of inverse phrase translations $P_{tm}(\underline{f}_i | \underline{e}_i)$
 - reordering model for each phrase $d(\text{START}(\underline{f}_i) \text{END}(\underline{f}_{i-1}) 1)$
 - language model for each word $P_{\text{Im}}(e_k | e_1, ..., e_{k-1})$
- Terms can be weighted: no longer Bayesian, but more efficient
 more components can be added + weights can be tuned
- So, now we maximize

 $\Pi_{i=1..M} \left(P_{\text{tm}}(\underline{f}_i | \underline{e}_i)^{\lambda_{\text{tm}}} \cdot d(\text{START}(\underline{f}_i) - \text{END}(\underline{f}_{i-1}) - 1)^{\lambda_{\text{rm}}} \right) \cdot \prod_{k=1..|e|} P_{\text{lm}} \left(e_k | e_1, ..., e_{k-1} \right)^{\lambda_{\text{lm}}}$ which can be expressed as $\exp(\sum \lambda_i h_i(e))$ using $h(...) = \log P(...)$

Translate f = find e (vector of features) that maximizes
 a weighted sum of feature functions (trained separately)

Extensions to the model

- Additional useful factors determined empirically
 - 1. Direct translation probabilities (use both e|f and f|e)
 - 2. Lexical weighting of phrase pairs
 - re-estimate likelihood of a pair based on the translation probabilities of the words that compose it
 - because rare phrases might get high $P_{tm}(f_i | \underline{e}_i)$ scores
 - 3. Word penalty: multiply by ω for each word
 - $\omega < 1$ favors shorter translations, $\omega > 1$ favors longer ones
 - tune the value of $\boldsymbol{\omega}$
 - 4. Phrase penalty: ρ factor, similar to ω
- Reordering model can be improved
- Phrase-based models can be trained directly using EM
 - but results are not better than the word-based approach

Conclusions

- Many methods for translation modeling
- We presented two principal approaches
 IBM Models and Phrase-based
- We showed the importance of word alignment
- Many variants and extensions exist
- More complex models: syntax-based = hierarchical
 must learn tree-based translation models
- Missing elements:
 - LANGUAGE MODELING, or how to estimate P(e)
 - DECODING, or how to search for e

References

- Philipp Koehn, Statistical Machine Translation, Cambridge University Press, 2010

 chapters 4 and 5
- Jörg Tiedemann, *Bitext Alignment*, Morgan & Claypool, 2011
- P. Brown, S. Della Pietra, V. Della Pietra, and R. Mercer, "The mathematics of statistical machine translation: parameter estimation", *Computational Linguistics*, 19(2), p. 263-311, 1993
 - introduced IBM Models 1-5
- Franz Josef Och and Hermann Ney, "A Systematic Comparison of Various Statistical Alignment Models", *Computational Linguistics*, 29(1), p. 19-51, 2003
 - review of past state-of-the-art algorithms for word alignment
- Dempster, A.P. and Laird, N.M. and Rubin, D.B., "Maximum likelihood from incomplete data via the EM algorithm", *Journal of the Royal Statistical Society Series B*, p. 1-38, 1977.
 - introduced the EM algorithm

Practical work

- Install the Moses MT system, build a phrasebased translation model
 - Sections 1, 2, 4 (up to 4.2 included) of 'TP-MTinstructions'
 - optionally: Section 3 to verify that Moses works
- Goal: train Moses on a domain or language pair of your own, examine the translation models (size and "perceived quality"): 4.2