

School of Computing and Information Systems
The University of Melbourne
COMP90042 WEB SEARCH AND TEXT ANALYSIS (Semester 1, 2019)

Workshop exercises: Week 12

Discussion

1. What aspects of human language make automatic translation difficult?

The whole gamut of linguistics, from lexical complexity, morphology, syntax, semantics etc. In particular if the two languages have very different word forms (e.g., consider translating from an morphologically light language like English into Turkish, which has very complex morphology), or very different syntax, leading to different word order. These raise difficult learning problems for a translation system, which needs to capture these differences in order to learn translations from bi-texts, and produce these for test examples.

2. For the following “bi-text”:

Language A	Language B
green house	casa verde
the house	la casa

- (a) What is the logic behind **IBM Model 1** for deriving word alignments?
 - The core idea is that we are going to have a translation table which stores the probability of translating a word from the target language into every possible word in the source language (again, this is the wrong direction due to the noisy channel model).
 - The probability of a sentence can then be treated as a **uni-gram** probability, conditioned on how the tokens in the two sentences are aligned. Or, essentially, the product of all of the corresponding probabilities from the translation table.
- (b) Work through the first few iterations of using the **Expectation Maximisation** algorithm to build a translation table for this collection. Check your work by comparing to the `WSTA_N21_machine_translation.ipynb` output.
 - We need to establish the direction of translation before we begin (although it isn't important in this particular example): let's say, we're trying to translate language B into language A. Consequently, we want the alignments where we're translating A into B. (Again, the opposite direction, due to the “noisy channel” model.)
 - We're going to initialise our translation table T with **uniform** values: every word from A is equally likely to be translated as every word from B:
 - The other thing we'll want to establish is the set of possible alignments. This can be done exhaustively by hand, because the “sentences” under consideration are so short; this is not practical for longer sentences, however.

T	casa	la	verde	Total
green	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	1
house	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	1
the	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	1

- I'm going to follow the notebook in ignoring the possibility of alignment words from B with the null element of A. We can deal with the converse — where words in A align with the null element of B by simply not aligning them to anything. (Proper models also deal with the former.)
- Consequently, each of the two sentences (called I and II below) has four (2^2) possible alignments (where every token of B is accounted for), namely:
 - Ia: green aligns with casa and house aligns with verde
 - Ib: green aligns with verde and house aligns with casa
 - Ic: green aligns with casa and verde (house implicitly aligns to null)
 - Id: house aligns with casa and verde
 - IIa: the aligns with la and house aligns with casa
 - IIb: the aligns with casa and house aligns with la
 - IIc: the aligns with la and casa
 - IId: house aligns with la and casa
- Now, we're going to calculate the **expected** likelihood of each of these possible alignments, according to the following formula:

$$\hat{P}(F, A|E) = \frac{\epsilon}{(I+1)^J} \prod_{j=1}^J t(f_j|e_{a_j})$$

- A close inspection might lead us to say that the (+1) should be excluded, because we're neglecting the null term from the A tokens, but it doesn't actually matter, as we'll see in a moment.
- For Ia, we observe the following:

$$\begin{aligned} \hat{P}(F, A|E) &= \frac{\epsilon}{(I+1)^J} t(\text{casa}|\text{green})t(\text{verde}|\text{house}) \\ &= \frac{\epsilon}{(2+1)^2} \left(\frac{1}{3}\right)\left(\frac{1}{3}\right) = \frac{\epsilon}{9 \cdot 9} \end{aligned}$$

- Because our translation table is uniform, every calculation will look the same.
- Now, we're going to make a **maximum** likelihood estimate of each entry in our translation table. We do this by summing the expected probability of the alignment for each possible translation.
- For green:
 - It aligns with casa in Ia ($\frac{\epsilon}{9 \cdot 9}$) and Ic (same), to give a total of $\frac{\epsilon}{9 \cdot 9}$.
 - It aligns with verde in Ib ($\frac{\epsilon}{9 \cdot 9}$) and Ic (same), to give a total of $\frac{\epsilon}{9 \cdot 9}$.
 - It never aligns with la, because they don't appear in a sentence together.

T	casa	la	verde	Total
green	$\frac{\epsilon}{99}$	0	$\frac{\epsilon}{99}$	$\frac{\epsilon}{99}$
house	$\frac{\epsilon}{99}$	$\frac{\epsilon}{99}$	$\frac{\epsilon}{99}$	$\frac{\epsilon}{99}$
the	$\frac{\epsilon}{99}$	$\frac{\epsilon}{99}$	0	$\frac{\epsilon}{99}$

- Let's summarise our likelihoods in the (un-normalised) translation table.
- We will now normalise the rows so that they look like probabilities. Doing this causes all of the $\frac{\epsilon}{9}$ terms to vanish; consequently, we will just ignore them for the rest of the steps below.
- After simplifying, here is the new translation table:

T	casa	la	verde	Total
green	$\frac{1}{2}$	0	$\frac{1}{2}$	1
house	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$	1
the	$\frac{1}{2}$	$\frac{1}{2}$	0	1

- At this point, it perhaps isn't obvious that this table will give us better alignment estimates, but it does:
- For Ia, we observe the following (ignoring the ϵ term):

$$\begin{aligned}\hat{P}(F, A|E) &= t(\text{casa}|\text{green})t(\text{verde}|\text{house}) \\ &= \left(\frac{1}{2}\right)\left(\frac{1}{4}\right) = \frac{1}{8}\end{aligned}$$

- For Ib:

$$\begin{aligned}\hat{P}(F, A|E) &= t(\text{verde}|\text{green})t(\text{casa}|\text{house}) \\ &= \left(\frac{1}{2}\right)\left(\frac{1}{2}\right) = \frac{1}{4}\end{aligned}$$

- For Ic:

$$\begin{aligned}\hat{P}(F, A|E) &= t(\text{casa}|\text{green})t(\text{verde}|\text{green}) \\ &= \left(\frac{1}{2}\right)\left(\frac{1}{2}\right) = \frac{1}{4}\end{aligned}$$

- For Id:

$$\begin{aligned}\hat{P}(F, A|E) &= t(\text{casa}|\text{house})t(\text{verde}|\text{house}) \\ &= \left(\frac{1}{2}\right)\left(\frac{1}{4}\right) = \frac{1}{8}\end{aligned}$$

- The calculations for II are similar.
- Updating the alignment counts for green gives us:
 - It aligns with casa in Ia ($\frac{1}{8}$) and Ic ($\frac{1}{4}$), to give a total of $\frac{3}{8}$.
 - It aligns with verde in Ib ($\frac{1}{4}$) and Ic (same), to give a total of $\frac{1}{2}$.
 - It never aligns with la.
- Our un-normalised counts (neglecting the ϵ terms) are now:

T	casa	la	verde	Total
green	$\frac{3}{8}$	0	$\frac{1}{2}$	$\frac{7}{8}$
house	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{3}{4}$
the	$\frac{1}{8}$	$\frac{1}{2}$	0	$\frac{5}{8}$

T	casa	la	verde	Total
green	$\frac{3}{7}$	0	$\frac{4}{7}$	1
house	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	1
the	$\frac{1}{7}$	$\frac{4}{7}$	0	1

- We can see that we have correctly observed that green is most likely to be verde, house to be casa, and the to be la. Summarising the normalised probabilities:
- Further iterations will continue to improve these counts, and to observe that Ib and IIa are the most likely alignments.