

Machine Translation

Assignment 1

Handed out: February 1st, due: February 15th

IBM Model 1

Implement the EM algorithm for IBM Model 1 (see class six) in your favorite programming language. Test it on the corpora available on the Wiki <http://www.statmt.org/mtm2/?n=Main.SummerSchool> for a summer school we organised one year.

- a toy corpus
- a small segment of the Europarl corpus <http://www.statmt.org/europarl.tgz> (French-English or German-English)

Your program should output two different things:

- A table containing the word translation probabilities that were learned (note: think of an efficient data structure for such a sparse matrix)
- The most likely alignment for each sentence pair in the training data

Please return a report containing:

- 1 page description of your source code
- 1 page excerpt of the word translation table
- 1 page excerpt of Viterbi alignments
- your source code

Pseudo-code of IBM Model 1 as presented in the lecture:

Require: set of sentence pairs (e , f)	14: {collect counts}
Ensure: translation prob. $t(e f)$	15: for all words e in e do
1: initialize $t(e f)$ uniformly	16: for all words f in f do
2: while not converged do	17: $\text{count}(e f) \ += \frac{t(e f)}{\text{s-total}(e)}$
3: {initialize}	18: $\text{total}(f) \ += \frac{t(e f)}{\text{s-total}(e)}$
4: $\text{count}(e f) = 0$ for all e, f	19: end for
5: $\text{total}(f) = 0$ for all f	20: end for
6: for all sentence pairs (e , f) do	21: end for
7: {compute normalization}	22: {estimate probabilities}
8: for all words e in e do	23: for all foreign words f do
9: $\text{s-total}(e) = 0$	24: for all English words e do
10: for all words f in f do	25: $t(e f) = \frac{\text{count}(e f)}{\text{total}(f)}$
11: $\text{s-total}(e) \ += \ t(e f)$	26: end for
12: end for	27: end for
13: end for	28: end while