

CMPE 409 - Machine Translation

Middterm Project

Deadline: 23:00, May 17

In your assignment, explain your codes with *comments*. Without *comments*, your assignment will not be marked.

Problem

In this assignment you are asked to implement the IBM Module I (Figure I) with python programming language.

Input: set of sentence pairs (\mathbf{e}, \mathbf{f})	14: // collect counts
Output: translation prob. $t(e f)$	15: for all words e in \mathbf{e} do
1: initialize $t(e f)$ uniformly	16: for all words f in \mathbf{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 (\mathbf{e}, \mathbf{f}) do	21: end for
7: // compute normalization	22: // estimate probabilities
8: for all words e in \mathbf{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 \mathbf{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

Figure 1: EM training algorithm for IBM Model 1

In your code test the following tasks:

- In each iterations print out:

- Iterate total five times over the two pairs sentence that we have explained in the lecture
 - $s\text{-total}(e)$ values for each pairs
 - expected counts: $\text{count}(e|f)$
 - $\text{total}(f)$
 - estimate probabilities: $t(e|f)$
- Test your code with five pairs Turkish-to- English sentences (5 parallel sentence, you can write these sentence yourself). Just record result of the last iteration
- Compare your results with the Python IBM modules of the NLTK library
- Write a short report that contains the result that you have gotten from previous tasks.

Submission

- Submit your source code with a **readme.txt** file. Your code should include comments. Be sure you have understand it.
- Submit your report: It contains the data you have gotten from in previous section.

Hint: Look at the lecture notes and examples as references