



## Exercise Sheet

### Phrase Extraction, Language Models, Noisy Channel

(Thanks Costanza Conforti)

Please turn in your writeup as a PDF named as usual, like: `fraser_alexander_ex3.pdf`

#### Exercise 1. Phrase Extraction

Given the English-German sentence pair below, with the relative word alignment, extract all phrases consistent with the word alignment<sup>1</sup>.

- How many distinct contiguous English and German phrases does this example contain? How many phrases are extracted?
- In some cases, given an English phrase, it is not possible to extract matching phrases in German. Find at least one example in this sentence pairs.
- Which is the effect of unaligned words in the number of extracted phrases?
- How many phrases can be extracted from a sentence of length  $n$ ?

	michael	geht	davon	aus	,	dass	er	im	haus	bleibt
michael	■									
assumes		■	■	■						
that						■				
he							■			
will										■
stay										■
in								■		
the									■	
house									■	

```

Input: word alignment A for sentence pair (e, f)
Output: set of phrase pairs BP
1: for e_start = 1 ... length(e) do
2:   for e_end = e_start ... length(e) do
3:     // find the minimally matching foreign phrase
4:     (f_start, f_end) = (length(f), 0)
5:     for all (e, f) in A do
6:       if e_start ≤ e ≤ e_end then
7:         f_start = min(f, f_start)
8:         f_end = max(f, f_end)
9:       end if
10:    end for
11:    add extract(f_start, f_end, e_start, e_end) to set BP
12:  end for
13: end for
function extract(f_start, f_end, e_start, e_end)
1: return [] if f_end == 0 // check if at least one alignment point
2: // check if alignment points violate consistency
3: for all (e, f) in A do
4:   return [] if e < e_start or e > e_end
5: end for
6: // add phrase pairs (incl. additional unaligned f)
7: E = {}
8: f_e = f_start
9: repeat
10:  f_e = f_end
11:  repeat
12:    add phrase pair (e_start - e_end, f_e - f_e) to set E
13:    f_e++
14:  until f_e aligned
15:  f_e--
16: until f_e aligned
17: return E
    
```

<sup>1</sup>Both the pseudo-code and the example in this exercise are taken from Koehn, Philipp. *Statistical machine translation*. Cambridge University Press, 2009

**Exercise 2.** Language Model, Noisy Channel<sup>2</sup>

- (a) Calculate the parameters  $p(e|e-1)$  of a Bigram Language Model from a corpus containing the following four sentences:

START the house is small  
 START the house is little  
 START the building is small  
 START the home building is small

- (b) Given the sentence  $\mathbf{f}$  = “das Haus ist klein” and the following parameters  $p(f|e)$ :

the		house		building		is		exists		little		small	
der	0.3	Haus	0.5	Gebäude	0.7	ist	0.7	ist	0.2	klein	0.7	klein	0.8
die	0.3	Heim	0.1	Haus	0.3	existiert	0.2	existiert	0.7	kurz	0.3	kurz	0.2
das	0.4	Gebäude	0.4			hat	0.1	hat	0.1				

calculate  $p(e|f)$  of the following translations:

$\mathbf{e}$  = “the building is little”

$\mathbf{e}$  = “the house exists small”

using the Language Model of point (a) and IBM Model 1 as translation model.

Recover that

$$\operatorname{argmax}_e p(e|f) = \operatorname{argmax}_e \frac{p(f|e)p(e)}{p(f)} = \operatorname{argmax}_e p(f|e)p(e) \quad (1)$$

---

<sup>2</sup>Originally conceived by Laura Jehl, PhD candidate at the University of Heidelberg, modified by Costanza Conforti