



Exercise Sheet

Evaluating MT output

(Thanks to Costanza Conforti, Laura Jehl)

Please submit your writeup in a single PDF file called `lastname_firstname_ex4.pdf` (e.g., `fraser_alexander_ex4.pdf`).

Exercise 1. Given the following output translations of 2 MT Systems, with the corresponding reference¹:

- SYSTEM A - Israeli officials responsibility of airport safety
- SYSTEM B - airport security Israeli officials are responsible
- REFERENCE - Israeli officials are responsible for airport security

(a) Calculate the **Position-independent error rate** (PER) of both system.

Recall that

$$PER = 1 - \frac{\text{correct} - \max(0, \text{output_length} - \text{reference_length})}{\text{reference_length}} \quad (1)$$

What do you observe?

(b) Then, calculate the **Word Error Rate** (WER) of both systems.

WER employs the *Levenshtein* distance, which is defined as the minimum number of editing steps needed to match two sequences. Considered editing steps are:

- *substitution* = replace one word with another
- *insertion* = add word
- *deletion* = drop word

To calculate the WER, proceed as follows: compute the Levenshtein distance by filling the matrix below, containing words from SYSTEM B output on one axis and reference words on the other, as in has been done in the example below for SYSTEM A output. Remember that the Levenshtein distance between two strings of length $|a|$ and $|b|$ is given by:

Given the Levenshtein distance, calculate the WER of both systems using the formula:

$$WER = \frac{\text{substitutions} + \text{insertions} + \text{deletions}}{\text{reference_length}} \quad (2)$$

What do you observe?

¹The examples of this exercise, as well as the Figure, are taken from Koehn, Philipp. *Statistical machine translation*. Cambridge University Press, 2009

$$\text{lev}_{a,b}(i, j) = \begin{cases} \max(i, j) & \text{if } \min(i, j) = 0, \\ \min \begin{cases} \text{lev}_{a,b}(i-1, j) + 1 \\ \text{lev}_{a,b}(i, j-1) + 1 \\ \text{lev}_{a,b}(i-1, j-1) + 1_{(a_i \neq b_j)} \end{cases} & \text{otherwise.} \end{cases}$$

		Israeli	officials	responsibility	of	airport	safety
	0	1	2	3	4	5	6
Israeli	1	0	1	2	3	4	5
officials	2	1	0	1	2	3	4
are	3	2	1	1	2	3	4
responsible	4	3	2	2	2	3	4
for	5	4	3	3	3	3	4
airport	6	5	4	4	4	3	4
security	7	6	5	5	5	4	4

		airport	security	Israeli	officials	are	responsible
	0	1	2				
Israeli	1	1	2				
officials	2						
are	3						
responsible	4						
for	5						
airport	6						
security	7						

(c) Finally, calculate the **BLEU** score of both sentences.

BLEU score combines n -gram precision with a *brevity penalty*, defined as $\min(1, \frac{\text{output_length}}{\text{reference_length}})$.

Using the formula:

$$\text{BLEU-N} = \min(1, \frac{\text{output_length}}{\text{reference_length}}) \prod_{i=1}^n \text{precision}_i \quad (3)$$

calculate the BLEU score with maximum order n for n -grams to be matched equal to 1, 2, 3 and 4. (Note that BLEU-4 is standardly used in the literature)

Metric	SYSTEM A	SYSTEM B
precision (1-gram)	/6	/6
precision (2-gram)	/5	/5
precision (3-gram)	/4	/4
precision (4-gram)	/3	/3
brevity penalty (4-gram)	/7	/7
BLEU-1		
BLEU-2		
BLEU-3		
BLEU-4		