Improving Coverage of Translation Memories with Language Modelling

Vít Baisa, Josef Bušta, Aleš Horák

Natural Language Processing Centre Faculty of Informatics Masaryk University

RASLAN 2014

Introduction

- Translation memories:
 - used in computer-aided translation systems,
 - manually built,
 - relatively small and focused,
 - usually in-house and not for (even academical) use.
- Our goal is to expand a TM to increase its coverage.
- We work with $En \leftrightarrow Cs$ language pair.

Related work

- TM are understudied resources,
- related topic: machine translation, example-based machine translation (EBMT),
- papers focused on searching and matching algorithms for CAT systems,
- ► *WeBiText*: TM from bilingual Canadian websites.
- TransSearch: EBMT system, Hansard corpus; linguistically motivated segments.

Methods for expanding TM

- Subsegment generation,
- subsegment combination,
- subsegment lexicalization,
- machine translation of subsegments.

Subsegment generating

- Use TM to train translational model (MGIZA),
- build word alignment for pairs from TM, it can be displayed in a word matrix,
- generate new subsegments and score them,
- resulting pairs can be added directly to TM or
- can be combined together.



Join & substitute

- we can make new segments by two methods
- join we can concatenate several subsegments and check the result against a language model trained on large monolingual corpus
- substitute sometimes a segment to be translated and a segment in a TM differ only in one or two words in the middle, in that case we can translate it using another subsegment but at the same time locate position of this in-the-middle words using the word matrix; these words may be translated automatically (using a dictionary, another subsegment from a TM) or simply left out for manual translation
- overlapping and non-overlapping join and substitute can be distinguished

Overlapping substitute example

new subsegment	Provozovatelé musí dodržovat zvláštní pravidla pro výzkumné
its translation	Operators shall comply with the special rules on research
from subsegments	Provozovatelé musí vytvářet zvláštní pravidla pro výzkumné musí dodržovat zvláštní
their translations	Operators shall create the special rules on research shall comply with the special

Subsegment lexicalizatoin

- generalization of the previous method
- all segments are tokenized and lemmatized
- searching and matching operations work on lemmata
- Ijoin concatenation of two different segments from TM and TM^{sub} on lemmata; when concatenating into new resulting segments, appropriate word form (case, gender and number) is generated
- Isubstitute substitution of a part of target segment with another segment using lemmata

With this method we expect increasing the recall (coverage) but at the same time not decreasing the translation accuracy of original segments from *TM*. So it is partially rule-based method.

Machine translation of subsegments, example

A sentence from TM: Návod na použití desinfekčního přípravku najdete na konci této brožury

A manual translation:

You can find instructions for use of disinfectant at the end of this brochure

A sentence for translation:

Návod na použití kartáče na vlasy najdete na konci této brožury

Not in TM: kartáče na vlasy

Google Translate returns: *hairbrush* (after lemmatization).

 \rightarrow Substitute the translation in the existing segment from TM.

Evaluation: subsegments generation & combination

We used a sample of TM and a testing document provided by a Czech translation services provider; as evaluation metrics we used the one used by MemoQ (CAT system).

	TM		TM ^{sub}		TM ^{NS}	
	Seg	%	Seg	%	Seg	%
matches	23	0.4	165	0.51	0	0
	TM ^{OJ}		TM ^{<i>NJ</i>}		TM ^{all}	
	ΤM	1 <i>01</i>	TN	1 ^{NJ}	TN	∧ ^{a∥}
	TM Seg	1 ^{0J} %	TN Seg	1 ^{NJ} %	TN Seg	1 ^{all} %

 $\begin{array}{ll} \mathsf{TM}^{OJ}-\mathsf{overlapping\ join} & \mathsf{TM}^{NJ}-\mathsf{non-overlapping\ join} \\ \mathsf{TM}^{sub}-\mathsf{generated\ subsegments} & \mathsf{TM}^{NS}-\mathsf{substitute} \\ & \mathsf{TM}-\mathsf{translation\ memory} \\ \mathsf{TM}^{all}=\mathsf{TM}+\mathsf{TM}^{sub}+\mathsf{TM}^{NS}+\mathsf{TM}^{OJ}+\mathsf{TM}^{NJ} \end{array}$

Evaluation: subsegments generation & combination

For an independent comparison, we also present our results for DGT translation memory; as evaluation metrics we used the one used by MemoQ (CAT system).

	ТМ		TM ^{sub}		TM ^{NS}	
	Seg	%	Seg	%	Seg	%
matches	31	0.03	276	0.25	58	0.45
	TM ^{OJ}					
	ΤM	101	ΤN	1 ^{NJ}		1 ^{all}
	T№ Seg	1 ^{0J} %	TN Seg	1 ^{NJ} %	TN Seg	Л ^{аШ} %

 $\begin{array}{ll} \mathsf{TM}^{OJ}-\mathsf{overlapping\ join} & \mathsf{TM}^{NJ}-\mathsf{non-overlapping\ join} \\ \mathsf{TM}^{sub}-\mathsf{generated\ subsegments} & \mathsf{TM}^{NS}-\mathsf{substitute} \\ & \mathsf{TM}-\mathsf{translation\ memory} \\ \mathsf{TM}^{all}=\mathsf{TM}+\mathsf{TM}^{sub}+\mathsf{TM}^{NS}+\mathsf{TM}^{OJ}+\mathsf{TM}^{NJ} \end{array}$

Evaluation: METEOR score

			TM ^s		
feature	TM ^{sub}	ТМ ^{ОЈ}	ΤΜ ^{ŊJ}	TM ^{NS}	TM ^{all}
precision	0.60	0.63	0.70	0.66	0.61
recall	0.67	0.74	0.74	0.71	0.68
f1	0.64	0.68	0.72	0.68	0.64
METEOR score	0.31	0.37	0.38	0.38	0.31
		L	JG I-IVI I		
feature	TM ^{sub}	TM ^{OJ}	TM ^{NJ}	TM ^{NS}	TM ^{all}
feature	TM ^{sub} 0.76	۲M ^{OJ} 0.93	TM ^{NJ} 0.91	TM ^{NS} 0.81	TM ^{all} 0.80
feature precision recall	TM ^{sub} 0.76 0.78	TM ^{OJ} 0.93 0.86	TM ^{NJ} 0.91 0.88	TM ^{NS} 0.81 0.85	TM ^{all} 0.80 0.81
feature precision recall f1	TM ^{sub} 0.76 0.78 0.77	TM ^{OJ} 0.93 0.86 0.89	TM ^{NJ} 0.91 0.88 0.89	TM ^{NS} 0.81 0.85 0.83	TM ^{all} 0.80 0.81 0.81

Conclusion and the future

- The preliminary results are promising,
- we are working on improvement of the first two methods,
- the rest of methods will be implemented,
- we expect a higher coverage.
- More detailed evaluation using biggger data, comparison with other techniques.