A Study of Style in Machine Translation: Controlling the Formality of Machine Translation Output

Xing Niu, Marianna Martindale, and Marine Carpuat CLIP Lab, University of Maryland



2017.09.11 @ EMNLP

- "Anybody hurt?"
- "Is someone wounded?"

A Human Translation Service*

it more "Hey Dude"	or "Dear Sir"?	one of the content	Order total	\$520.8
Informal	•	one of the content.	Estimated delivery 15 hou	irs. ()
Informal			Quality Policy	
Friendly Business	nsiator		Updated on 03/16/2017	
Formal Other			Payment method: 🔹 🖲	Credit card 🛛 🔘 PayP
ossible instructions	Voice Links	Casual, romantic, funny, serious etc. To your website, screen shots or other docs.	Pay & Confi	rm Order
	Purpose & Audience	This is going to my most important client etc.	View Full	Ouote

*gengo.com

Motivation

- "Anybody hurt?"
- "Is someone wounded?"
 - same literal meaning
 - for different audience or environment
- Goal: controlling formality in machine translation
 - ... by asking what is the expected level of formality
 - Prior work looked at other aspects of style
 - politeness in German (Sennrich et al. 2016), gender traits (Rabinovich et al. 2017)

- How to control formality in machine translation?
 - Re-ranking-based Formality-Sensitive Machine Translation (FSMT)

Formality-Sensitive MT



Formality-Sensitive MT

• *n*-best list re-ranking with a new feature: $f(h; \ell) = |\text{Formality}(h) - \ell|$

- Formality(*h*): the sentence-level formality score [-1,1] of translation hypothesis *h*.
- ℓ : the desired formality level.

- How to control formality in machine translation?
 - Re-ranking-based Formality-Sensitive Machine Translation (FSMT)
- How to score sentential formality?
 - Evaluating existing formality modeling methods

Formality Modeling

- Lexical formality models based on vector space models and formal/informal seed words:
 - SimDiff (Brooke et al. 2010) compares words to formal vs. informal seeds.
 - Support Vector Machine (SVM) finds a hyperplane that separates seeds.
 - Projecting a word in the high-dimensional space to a one-dimensional score using Principal Component Analysis (PCA) or Densifier (Rothe et al. 2016).

Formality Modeling – Intrinsic Evaluation

- sentential formality = weighted average of lexical formality
 - Comparing sentential scores with human annotations (11,263 sentences)
 - Lahiri (2015); Pavlick and Tetreault (2016)
 - Metrics: Spearman correlation.

Formality Modeling – Intrinsic Evaluation

	Latent Semantic Analysis (LSA)	Word2vec	Training data for vector space models: ICWSM 2009 Spinn3r
SimDiff	0.660	0.654	1.0 billion words nom blogs
SVM	0.657	0.585	
PCA	0.656	0.663	
Densifier	0.664	0.644	

- Models are close in performance.
- Densifier-LSA is selected as a representative for our FSMT system.

- How to control formality in machine translation?
 - Re-ranking-based Formality-Sensitive Machine Translation (FSMT)
- How to score sentential formality?
 - Evaluating existing formality modeling methods
- How effective is the FSMT system?
 - Automatic + Human evaluations

Formality-Sensitive MT – Evaluation

Data: MultiUN + OpenSubtitles2016 (French->English)

- 3 FSMT systems, with different desired formality
 - Low (ℓ=-0.4) | Neutral (ℓ=0) | High (ℓ=+0.4)

FSMT – Automatic Evaluation (BLEU)

Desired formality	Informal test set	Neutral test set	Formal test set
None (baseline)	39.74 🦱	40.17	47.97
Low	40.27 🦯 🗸	39.65 +	0.3 47.76
Neutral	38.70	40.46 🦊	47.84
High	37.58	39.53	47.97 🔶
		Shorter ser	ntence \rightarrow larger impact.

Formal sentences (MultiUN) are sufficiently different.

• Best: when desired formality level matches reference.

FSMT – Automatic Evaluation (BLEU)

Desired formality	Informal test set	Neutral test set	Formal test set
None (baseline)	39.74	40.17	47.97
Low	40.27	39.65	47.76
Neutral	38.70	40.46	47.84
High	37.58	39.53	47.97

- Δ BLEU \approx 3 \rightarrow translation quality difference is large.
- BLEU scores conflate translation errors and stylistic mismatch.

FSMT – Human Assessment

- 42 random affected translation pairs for $\ell = \pm 0.4$
- 15 volunteers
- Changes in formality:
 - not impact on adequacy
 - small impact on fluency
- According to formality judgment:
 - FSMT impacts 22/42 examples.
 - FSMT correctly yields 21/22 examples w.r.t. formality.
 - more formal output for ℓ =+0.4 than ℓ =-0.4.

FSMT – Human Assessment (Examples)

е	Examples	Comments
-0.4	anybody hurt ?	
+0.4	is someone wounded ?	annotated as more formal
-0.4	and then he ran away .	
+0.4	and then he escaped .	annotated as more formal
-0.4	he shot himself in the middle of it .	
+0.4	he committed suicide in the middle of it .	annotated as more formal
-0.4	to move things forward .	
+0.4	in order to move the process forward.	annotated as more formal
-0.4	how do you do ?	annotated as more formal
+0.4	how are you?	

- How to control formality in machine translation?
 - Re-ranking-based Formality-Sensitive Machine Translation (FSMT)
- How to score sentential formality?
 - Evaluating existing formality modeling methods
 - Empirical comparison \rightarrow similar performance
- How effective is the FSMT system?
 - Effective in controlling language formality without loss in translation quality
 - Based on automatic evaluation and human assessment

Code: <u>https://github.com/xingniu/computational-stylistic-variations</u>

