

An Open Source Toolkit for Word-level Confidence Estimation in Machine Translation

Christophe Servan, Ngoc Tien Le, Ngoc Quang Luong,
Benjamin Lecouteux and Laurent Besacier



4 Dec. 2015

Outline

1 Introduction

- Confidence Estimation (CE)
- Framework

2 WCE Toolkit

- Overview
- Features

3 Experiments

- Tasks
- State-of-the-Art
- Decision threshold
- Feature Selection

4 Demo

5 Conclusion and perspectives

Confidence Estimation (CE)

How much I trust the MT system outputs

Confidence Measures for Machine Translation

- Several levels:
 - ▶ Word
 - ▶ Sentence
 - ▶ Document
- Classification task:
 - ▶ Binary label (e.g.: Good/Bad)
 - ▶ Multiple labels (e.g.: kind of error)

⇒ can produce a score associated with the label.

⇒ Word-level Confidence Estimation with binary labels (and scores)

Framework

Computer Aided Translation (CAT) Tool

- help professional translators
- highlight which part to focus on
- speed up the post-editing process

⇒ already part of several CAT tools CasmaCAT/MateCAT, libelleX (Lingua & Machina)...

Toolkit overview

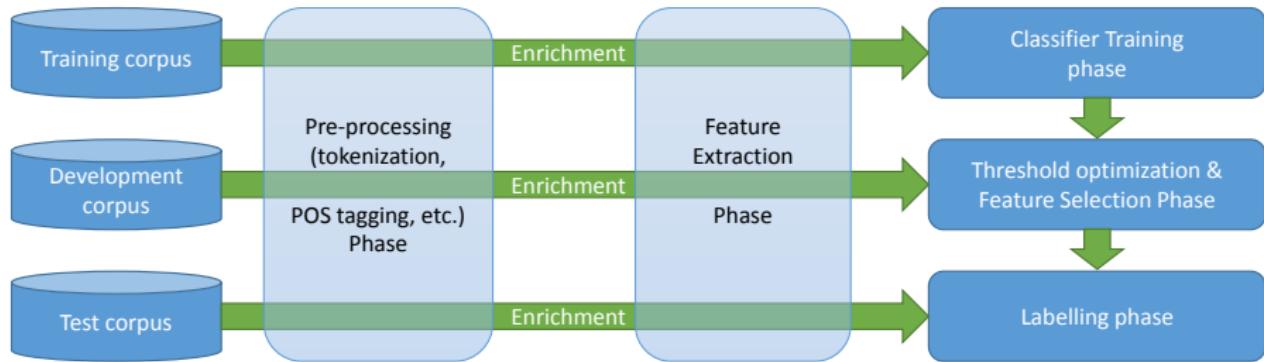


Figure: Pipeline of our Word-level Confidence Estimation tool

- Reproducible experiments (WMT14 QE shared task)
- Easy to add new features (written in python)
- External toolkits used: GIZA++, TreeTagger, SRILM, BabelNet, Berkeley parser, Wapiti...

Features Extracted

1 <i>Proper Name</i>	16 WPP Exact*
2 Unknown Stem	17 WPP Any*
3 # of Word Occurrences	18 WPP Min*
4 # of Stem Occurrences	19 WPP Max*
5 Source POS	20 Nodes*
6 <i>Source Word</i>	21 Numerical
7 Source Stem	22 Punctuation
8 Target POS	23 Stop Word
9 Target Word	24 Target Backoff Behaviour
10 Target Stem	25 <i>Constituent Label</i>
11 <i>Word context Alignments</i>	26 <i>Distance To Root</i>
12 <i>POS context Alignments</i>	27 <i>Polysemy Count – Target</i>
13 Stem context Alignments	28 Occur in Bing Translator
14 Longest Target N-gram Length	29 Occur in Google Translate
15 Longest Source N-gram Length	

*: internal features

Tasks for Word-level Confidence Estimation

two sets of experiments:

- en-es:

- ▶ WMT14 QE shared task
- ▶ State-of-the-Art evaluation task
- ▶ MT outputs only (*i.e.*: 1-best only)

- fr-en:

- ▶ in-house data
- ▶ SMT outputs with N -best list

Task protocol:

- Word-level Confidence Estimation
(one label per word)
- Binary labels (*OK* or *BAD*)
- Precision, Recall, F-Measure for each label and Mean F-Measure.

Amount of data (# lines):

Corpus	en-es	fr-en
Train	650	10K
Test	208	881
N -Best	1	1000

State-of-the-Art: WMT14 Quality Estimation task

The *en-es* Word-level Quality Estimation task¹

Systems	M-F	F(<i>bad</i>)
FBK-UPV-UEDIN-1	62.00	48.73
LIMSI	60.55	47.32
→ <i>Our toolkit</i>	60.76	47.17
LIG-1	63.55	44.47
LIG-2	63.77	44.11
FBK-UPV-UEDIN-2	62.17	42.63

¹Without internal features

Decision threshold optimization (oracle)

Task	Threshold	Label	P	R	F	M-F
en-es (WMT14)	Default	Good	71.24	77.73	74.35	60.76
		Bad	51.82	43.28	47.17	
	Optimized	Good	71.42	76.82	74.03	60.87
		Bad	51.49	44.45	47.71	
fr-en	Default	Good	84.45	90.22	87.24	64.96
		Bad	50.10	37.16	42.67	
	Optimized	Good	85.60	85.65	85.62	65.59
		Bad	45.61	45.50	45.56	

Set of features (no threshold optimization)

The *fr-en* task only:

Features	Labels	P	R	F	M-F
<i>Base.</i>	<i>Good</i>	81.97	92.22	86.80	58.64
	<i>Bad</i>	44.17	23.28	30.48	
<i>+ mod.</i>	<i>Good</i>	83.21	90.99	86.92	62.00
	<i>Bad</i>	47.24	30.53	37.09	
<i>+ new</i>	<i>Good</i>	83.55	90.11	86.70	62.65
	<i>Bad</i>	46.75	32.86	38.60	
<i>+ MT</i>	<i>Good</i>	84.45	90.22	87.24	64.96
	<i>Bad</i>	50.10	37.16	42.67	

Sequential Forward Selection algorithm (*fr-en* task only)

Rank	Feature	Rank	Feature
1	Stem context Alignments	16	Stop Words
2	WPP Exact	17	Nodes
3	<i>Word context Alignments</i>	18	# of Stem Occurrences
4	WPP Max	19	Numeric
5	WPP Any	20	Unknown Stem
6	WPP Min	21	Target Word
7	<i>POS context Alignments</i>	22	Source POS
8	Occur in Google Translate	23	<i>Polysemy Count – Target</i>
9	Longest Target <i>N</i> -gram Length	24	<i>Source Word</i>
10	Occur in Bing Translate	25	<i>Constituent Label</i>
11	Source Stem	26	Punctuation
12	Target Backoff Behaviour	27	Target Stem
13	Longest Source <i>N</i> -gram Length	28	<i>Proper Name</i>
14	# of Word Occurrences	29	Target POS
15	<i>Distance To Root</i>		

GitHub exploration:

<https://github.com/besacier/WCE-LIG>

Conclusion and perspectives

Contributions:

- Open-Source, Flexible and State-of-the-Art WCE Toolkit
- Decision Threshold and Feature Selection (SFS)
- Available on GitHub:
<https://github.com/besacier/WCE-LIG>

Perspectives:

- More flexible!
- Feature Selection with
Sequential Floating Forward Selection algorithm
- Add new features based on Word Embeddings
(monolingual and bilingual)
- Propose a version of the WCE toolkit for Speech Translation

Thanks for your attention!

Remarks & questions are welcome!

contact: christophe.servan@imag.fr & laurent.besacier@imag.fr