Automatic Normalisation of Modern French (i.e. from the 17th century)



Rachel Bawden Jonathan Poinhos Benoît Sagot Philippe Gambette







LREC - 2022



Eleni Kogkitsidou Simon Gabay











What is normalisation?

Application of a predefined convention to smooth out variation. We choose contemporary French as the norm.

Why normalise modern French texts?

- Spelling conventions were not yet fixed, so there is a high degree of variability between texts
- The texts look a lot like contemporary French, but the differences make it hard to apply standard French tools to them

au lieu que si j'eusse esté sage, je me fusse servy du changement qui paroissoit en elle, pour aider à me guerir; mon amour en redoubla, & je me conduisois si mal, que la

Accuracy (%)

Original	88.93
Normalised	91.10

PoS-tagging results when fine-tuning and testing French language model CamemBERT on 17th c. texts



Modern French vs. Contemporary French

- Some trivial changes (e.g. long s)
- spellings, other changes indicative of language change, etc.)



Many non-trivial changes (segmentation differences, introduction of "classical"

Previous work on normalisation

- Word lists, rules and edit-based approaches

 - E.g. Levenshtein distance is a strong baseline (Pettersson et al., 2013)
- MT approaches
 - and Erjavec, 2013; Petters- son et al., 2013b; Domingo and Casacuberta, 2021)
 - (Domingo and Casacuberta, 2018), but not always the case (Gabay and Barrault, 2020)

• Replacing words by others depending on predefined lists or predefined correspondences (manual or automatic) (Baron and Rayson, 2009; Bollmann et al., 2011, Porta et al., 2013).

 Most previous work has focused on character-based (learning correspondences of letters) in a word), largely evaluated on normalisation of individual words (Vilar et al., 2007; Scherrer

Statistical MT (SMT) vs. neural MT (NMT): SMT can be superior if little data is available

4

Contributions of this paper

1. New benchmark for the task: **parallel training data** for the normalisation of modern French into contemporary French

2. Development of **normalisation models** for Modern French: ruled-based, statistical, MT-inspired (statistical and neural)

3. **Evaluation metric** (symmetrised word accuracy) adapted to MT-inspired models and comparison of all models





- Sentence-aligned parallel dataset (modern-contemporary French)
- Genres: Caractères, comédie, tale, correspondence, law, fables, journalism, philosophy, poetry, novel, memoir novel, theology, tragedy, travel
 - Genres only available in test set: medicine, physics

		#unique tokens		#unique O	OV tokens
	#sents	ModFr	Fr	ModFr	Fr
Train	17.9k	264.3k	263.7k	_	_
Dev	2.4k	40.4k	40.3k	1.7k	1.3k
Test	5.7k	86.4k	86.2k	3.6k	2.5k

https://freem-corpora.github.io/



Normalisation methods

1. Simple rule-based method (regular-expression-based)

- Manually written based on simple corpus statistics (some purely typographic and others lexical)
- E.g. $\mathbf{i} \rightarrow \mathbf{s}, \tilde{\mathbf{o}} \rightarrow \mathbf{om}$ (if followed by m, b or p) and on (if not)

2. Statistical alignment-based model (ABA)

More details coming up!

3. MT-based approaches:

More details coming up!

* Optional lexicon-based post-processing step

- To be applied after the other 3 methods
- Replace words that match modulo certain regular changes (e.g. accents, long s)



Alignment-based approach (ABA)

- Word-level translation rules learned from an aligned training corpus
- Character-level transformation rules manually designed by observing frequent transformations
- For each word not recognised as being contemporary French:
 - replace by the word in the word-level transformation rule if it exists
 - apply all possible combinations of character-level transformation rules, keep the first word existing in contemporary French, keep the original word otherwise







Advantages:

- Flexible word segmentation: allows for word merging or splitting
- Words are normalised in context (helpful in some cases and even necessary in others):

	Normalisation example 1	Normalisation example 2
nostre 'our'	quel malheur est le nôtre	Les larmes sont trop peu pour pleurer notre mal
	'what woe is ours '	'The tears are too few to cry (for) our pain'
appellez 'call'	N'appelez point des yeux le Galant à votre aide	Royaumes, par nous vulgairement appelés Siam
	'Do not call the Galant for help with your eyes'	'kingdoms, known popularly by us as Siam'

MT-style approaches



- Statistical MT (SMT) (1) Moses
- Neural MT (NMT) Fairseq: (2) LSTM and (3) transformer

Extensive hyper-parameter searches:

- Subword segmentation (using sentencepiece and BPE):
 - Best subword segmentation with a vocabulary of 500 (SMT) and 1000 NMT
- Size of the networks (e.g. number of layers, embedding dimensions, etc.)
 - Best models were smaller than the base models used
- Learning rate and batch size

MT-style approaches



Evaluation

- Most commonly used metric = word/token-level accuracy
- - Not necessarily a one-to-one token-level alignment

Symmetrised accuracy:

Ref: Puisqu' Achille combat, nous allons triompher **Hyp**: Puisqu' Achile combat, et oui nous allons triompher

According to reference tokenisation

Puisqu'	Achille	combat	/	nous	allons	triompher
Puisqu'	Achile	combat	, <u> et oui</u>	nous	allons	triompher

According to hypothesis tokenisation

Puisqu'	Achile	combat	,	et	oui	nous
Puisqu'	Achille	combat				nous

• In need of a reproducible implementation and one that is adapted to sentence-level normalisation:

• Hallucinations need to be penalised (risk of them being all associated with a single reference token)

Accuracy = 6/7 = 0.86

Symmetrised acc = 0.82

11

allons	triompher
allons	triompher

Accuracy = 7/9 = 0.78

Results

Method	WordAcc (sym)	WordAcc (sy
Identity	72.73	43.00
+ Lefff	86.12	64.84
Rule-based	89.05	60.22
+ Lefff	90.85	66.51
ABA	95.14	69.50
+ Lefff	95.44	73.54
SMT	97.10±0.02	76.64±0.18
+ Lefff	97.24±0.02	78.37±0.20
LSTM	96.14±0.08	76.69±0.70
+ Lefff	96.25±0.10	78.35±0.79
Transformer	95.89±0.08	75.73±0.38
+ Lefff	96.01±0.09	77.51±1.00

ym) OOV

- Baselines already strong
- Best model = SMT
- Neural models do better on OOV words
- Postprocessing (+ Lefff):
 - Helps all methods
 - SMT+Lefff leads to best OOV scores





Similarity of the approaches

<u>ر</u>	Rule-based	ABA	SMT	LSTM
Fransforme	89.42	95.02	96.90	97.47
LSTM -	90.00	95.54	97.20	100.00
SMT -	90.16	96.17	100.00	97.20
ABA -	91.70	100.00	96.17	95.54
Rule-based	100.00	91.70	90.16	90.00

Transformer

100.00	
97.47	
96.90	

95.02

89.42

- Neural methods most similar (LSTM, Transformer)
- SMT is most similar to Transformer
- ABA most similar to rulebased



How zealous/conservative are the models?

Over-modifications



What is better? This is actually task-dependent:

- As an aid for manual normalisation: conservative better
- For a down-stream annotation task (e.g. PoS tagging): zealous better

Under-modifications





What sort of differences are there?

Comparison of the best rule-based approach (ABA+Lefff) and best MT approach (SMT+Lefff)

- ABA is less robust to inadequacies in the training corpus
 - E.g. succeeds with auoient \rightarrow avaient, but not avoient \rightarrow avaient (whereas SMT succeeds)
 - Lacks some rules (e.g. dealing double consonants)
- SMT is in general more "creative":
 - Some more creative errors (quite say to spot): ma pelsée 'pensée' -> pmentsée
 - Language model effect can be too strong (removes some determiners)
 - But handles ambiguity better (because it is contextual)
 - ABA+Lefff: Car enfin n'attends pas que mes feux redoublez,
 - SMT+Lefff: Car enfin n'attends pas que mes feux redoublés,
- The approaches appear to be complementary potential for combining the two!



Conclusion and perspectives

- French (dataset, baselines and state-of-the-art models)
- different advantages
 - Potential for combining them
- Aim to facilitate and encourage research on Modern French

New benchmark for the normalisation of Modern French into contemporary

Comparison of different approaches (rule-based and MT-inspired) with

• Further experiments required to test which model is best in different scenarios (i.e. to aid manual normalisation or for downstream tasks).

Our code and models are freely available: https://github.com/rbawden/ModFr-Norm

Further information on the FreEm project page: https://freem-corpora.github.io

Thank you very much!