

# Unsupervised Lemmatization as Embeddings-Based Word Clustering

Rudolf Rosa      Zdeněk Žabokrtský

Charles University, Faculty of Mathematics and Physics,  
Institute of Formal and Applied Linguistics  
Malostranské náměstí 25, 118 00 Prague 1, Czech Republic  
{rosa, zabokrtsky}@ufal.mff.cuni.cz

## Abstract

We focus on the task of unsupervised lemmatization, i.e. grouping together inflected forms of one word under one label (a lemma) without the use of annotated training data. We propose to perform agglomerative clustering of word forms with a novel distance measure. Our distance measure is based on the observation that inflections of the same word tend to be similar both string-wise and in meaning. We therefore combine word embedding cosine similarity, serving as a proxy to the meaning similarity, with Jaro-Winkler edit distance. Our experiments on 23 languages show our approach to be promising, surpassing the baseline on 23 of the 28 evaluation datasets.

## 1 Introduction

The task of lemmatization is to assemble inflections of the same word into one group, represented by a designated form of the word called lemma. It is a classical NLP task, potentially useful e.g. for information retrieval or machine translation.

The standard approach is to use supervised machine learning, exploiting a dataset of word and lemma pairs to train a lemmatizer (Kondratyuk et al., 2018). However, such datasets are available roughly for only 1% of world’s languages. An alternative is to use stemming (Lovins, 1968; Porter, 2001), which is typically rule-based, i.e. does not need annotated training data, but on the other hand is usually language-specific, and also tends to cluster the word forms too coarsely (many different lemmas may share the same stem).

As we want to cover many languages while keeping inflections of different lemmas in separate groups, we instead propose to perform unsupervised lemmatization as word form clustering.

The first step is to employ a suitable word form distance measure. We propose a measure combining (a) string similarity, implemented using edit

distance, and (b) meaning similarity, for which we use word embedding similarity as a proxy (Sec. 2).

We then precompute distances of probable inflections of the same words, and cluster them with agglomerative clustering (Sec. 3). We leave the last step of selecting a representant from each of the clusters as its lemma for future work.

We evaluate our setup on 28 datasets for 23 languages, finding that it outperforms the baseline on 23 of the datasets, but also identifying many of its limitations that yet need to be addressed (Sec. 4).

We make our source codes available together with this paper.

## 2 Word form distance measure

We propose a word form distance measure which combines string similarity with word embedding similarity, designed to assess word forms belonging to the same lemma as more similar than word forms belonging to different lemmas.

### 2.1 String similarity

For string similarity, we use the Jaro-Winkler (JW) edit distance (Winkler, 1990) from the `pyjarowinkler` Python package (implemented as a similarity).<sup>1</sup> Unlike Levenshtein (1966) edit distance, JW gives more importance to the beginnings of the strings than to their ends. We find this to be advantageous, as most of the inflection usually happens at the end of the word.<sup>2</sup>

To compute the edit distance of a pair of strings, we average their JW with JW of their *simplified variants*; the simplification consists of lowercasing, transliteration to ASCII using the Unidecode

<sup>1</sup><https://pypi.org/project/pyjarowinkler/>

<sup>2</sup>Based on the feature 26A of the WALS database by Dryer and Haspelmath (2013), most world’s languages, and in particular practically all of the languages in our data set, have a strongly or weakly suffixing inflectional morphology, whereas prefixing morphology is rare.

library,<sup>3</sup> and deletion of non-initial vowels (a e i o u y). This makes the measure somewhat softer, paying less attention to differences that tend to have lower importance in our setting.

## 2.2 Word embedding similarity

As shown by Mikolov et al. (2013), cosine similarity of word embeddings tends to capture morphological, syntactic, and semantic similarities of words. This motivates our use of word embedding similarity as a proxy to word meaning similarity.

We use the FastText word embeddings (Grave et al., 2018), which have the additional benefit of employing subword embeddings, thus also implicitly capturing string similarity to some extent.<sup>4</sup>

## 2.3 Combined distance measure

Our distance measure is based on a multiplication of the two similarities shifted to the  $[0, 1]$  interval:

$$\text{dist}(a, b) = 1 - JW(a, b) \cdot \frac{\cos(a, b) + 1}{2} \quad (1)$$

## 3 Clustering

We apply agglomerative clustering from Scikit-learn (Pedregosa et al., 2011), with average linkage.<sup>5</sup> The algorithm starts by assigning each word form to a separate cluster. In each step, it then merges the pair of clusters with the lowest average distance of their elements. The standard stopping criterion is to preset the final number of clusters to form. As we have not thought of a way to estimate the number of lemmas, we instead stop the algorithm once the cluster distance raises above a threshold  $t$ ; we use  $t = 0.4$ .

The algorithm only assigns a cluster to word forms that are part of the training vocabulary. If we encounter an out-of-vocabulary word form at test time, we perform a single clustering step with it: we assign it to the closest cluster if it is closer than  $t$ , otherwise, we put it into a new cluster.

### 3.1 Stem-based hyperclusters

Theoretically, we would like to compute the distances of all pairs of word forms. In practice, with

<sup>3</sup><https://pypi.org/project/Unidecode/>

<sup>4</sup>In our exploratory experiments on a Czech language dataset, we observed the accuracies to rise by approximately 20 percentage points when we substituted word2vec embeddings with FastText embeddings.

<sup>5</sup>Average linkage is recommended by the manual; we also tried single and complete linkage, but the results were worse.

Vocabulary size	OOV rate
1,000	50.3%
10,000	27.6%
100,000	8.8%
1,000,000	1.5%

Table 1: Proportion of test-data word forms that are not part of the vocabulary, for cs\_pdt.

a vocabulary of  $10^5$  word forms, computing the  $10^{10}$  distances would use prohibitive amounts of time and memory. Therefore, we use a stemming approach to pre-partition the space into hyperclusters, and run the clustering algorithm on each such hypercluster separately; word forms with different stems thus cannot be clustered into the same cluster. In this work, we define the stem of a word form as the first  $K$  characters of its *simplified variant* (see sec. 2.1); we use  $K = 3$ .<sup>6</sup>

Such crude stemming obviously separates some forms of the same lemma into separate hyperclusters (short words, irregular inflections, suppletives...), making it impossible for our method to reach the correct clustering. We intend to address this more properly in future work, especially by utilizing unsupervised morphological splitting to get better stems. However, some of the weak points of our approach, such as suppletives, probably cannot be easily resolved.

## 4 Experiments

### 4.1 Data

For the experiments reported in this paper, we use the pretrained word embedding dictionaries available from the FastText website.<sup>7</sup> The word embeddings had been trained on Wikipedia<sup>9</sup> and Common Crawl<sup>10</sup> texts with the FastText tool, “using CBOW with position-weights, in dimension 300, with character n-grams of length 5, a window of size 5 and 10 negatives”. We limit our vocabulary to  $N$  most frequent words, i.e. the first  $N$  words stored in the embedding dictionary; we use  $N = 100,000$ . We found that with a smaller dictionary, the method is more efficient computa-

<sup>6</sup>This reduces the time and memory complexity roughly 1,000 times; each experiment then uses about 1 GB and 1 h.

<sup>7</sup><https://fasttext.cc/>

<sup>8</sup>A drawback of using pre-trained word embeddings is that they typically use different tokenization than the evaluation data, necessarily leading to occasional problems.

<sup>9</sup><https://www.wikipedia.org/>

<sup>10</sup><http://commoncrawl.org/>

tionally, but the results are worse due to very high rates of out-of-vocabulary items (see Table 1).

We evaluate on treebanks from the Universal Dependencies 2.3 (Nivre et al., 2018). We use the *dev* part for evaluation, assigning each of the words in this part of the treebank to a cluster and then evaluating the clusters against the gold-standard lemmas (repeated word forms are used repeatedly in the evaluation, i.e. the evaluation is token-based instead of type-based). We used only treebanks that satisfy the following criteria:

- Contain at least 100,000 tokens.
- Lemmas are annotated (semi-)manually.<sup>11</sup>
- There are pretrained FastText embeddings available for the language.

This results in a set of 28 treebanks for 23 languages (a subset of the total 129 treebanks for 76 languages), listed in Table 2. For more information on the treebanks, please consult the UD webpage.<sup>12</sup>

We used the *cs\_pdt* treebank to tune the method and set its hyperparameters.

## 4.2 Evaluation

We evaluate the clustering using the standard Scikit *v-measure*, which penalizes both clusters containing forms of multiple lemmas as well as forms of a single lemma scattered in multiple clusters. The results are listed in Table 2.

As a baseline, we choose the better-performing of these two approaches for each dataset: either taking the form as the lemma, or taking the first 5 characters of the form as the lemma.<sup>13</sup>

The upper bound is an oracle, always selecting the gold standard lemma if it is located in the same hypercluster. Due to the stem-based hyperclustering (sec. 3.1), the oracle does not reach 100%; this constitutes one of the strongest limitations of our approach.<sup>14</sup>

We also express the performance of our method as error reduction on the scale from baseline (0%) to upper bound (100%).

Treebank	Baseline	Our	Upp.	Err.red.
ar_padt	form 4.19	3.90	2.93	23.1
ca_ancora	form 4.65	4.35	3.32	22.3
cs_cac	form5 3.56	2.25	1.14	54.0
cs_fictree	form5 4.82	4.08	2.68	34.6
cs_pdt	form5 4.93	3.41	1.65	46.6
da_ddt	form 2.32	2.16	1.55	21.2
en_ewt	form 2.29	2.22	1.78	13.8
es_ancora	form 3.99	3.38	2.25	34.7
et_edt	form5 4.78	4.31	2.54	20.9
fa_seraji	form 8.99	8.76	7.44	14.8
fr_gsd	form 4.12	3.81	2.70	22.0
hi_hdtb	form 4.18	3.58	2.83	44.3
hr_set	form5 4.04	2.87	1.71	50.2
it_isdt	form 4.27	3.71	2.78	37.8
it_postwita	form 3.60	4.07	2.37	-38.0
ja_gsd	form 1.64	1.93	1.41	-123.1
ko_kaist	form 0.14	2.41	0.11	-6392.8
la_ittb	form5 6.53	6.97	3.85	-16.4
la_proiel	form5 6.92	7.42	4.20	-18.4
lv_lvttb	form5 3.90	3.39	2.10	28.0
no_bokmaal	form 2.79	2.22	1.48	43.6
no_nynorsk	form 2.73	2.52	1.48	16.7
pl_lfg	form5 3.68	3.06	1.84	33.6
pt_bosque	form 3.57	3.17	2.55	39.0
ro_nonstd	form5 8.13	7.95	5.64	7.2
sk_snk	form5 2.87	2.01	0.63	38.2
uk_iu	form 2.66	1.94	0.88	40.7
ur_udtb	form 3.95	3.79	2.65	12.3
Average	4.08	3.77	2.45	-210.3
Median	3.97	3.40	2.31	22.7

Table 2: Results of form clustering, measured in % of  $1 - vmeasure$  (expressing the error, i.e. lower is better). Baseline (either full form or prefix of form of length 5), our system, and oracle upper bound. Last column is error reduction on the scale from baseline to upper bound, in %.

<sup>11</sup>I.e. the “Lemmas” feature in the treebank Readme file is “manual native”, “converted from manual”, “converted with corrections”, or “automatic with corrections”.

<sup>12</sup><https://universaldependencies.org/>

<sup>13</sup>For languages with little inflection, using the form is usually better; for highly inflectional languages, the 5 character prefix usually performs better.

<sup>14</sup>We have tried to selectively remerge at least some of the hyperclusters by using a coarser but more efficient merging strategy; however, we have not been successful with this approach so far.

### 4.3 Discussion

For 23 of the 28 datasets, our method achieves a positive error reduction; the median error reduction is 23%. Because of the extreme result for Korean, the average does not make much sense here.

The results are worst for Korean and Japanese, which are analytical languages with practically no inflection, making the “form” baseline very close to the upper bound – the clusters should mostly have the size of 1. As our hyperparameters are not tuned for this, and the whole idea of lemmatizing these languages is questionable, our results are very low here. We also observe deteriorations for a treebank of Italian Tweets and for treebanks of historical Latin, which are all known to be very hard datasets.

On all other datasets, we observe an improvement over the baseline, with an error reduction typically between 10% and 35%. The performance is especially good for Slavic languages (cs, hr, pl, sk, uk), where the error reduction is often around 50%. This is most probably due to the hyperparameters being tuned on the cs\_pdt treebank.

The threshold  $t$  controls the balance between too conservative and too eager cluster merging. Being too conservative typically leaves some string-wise distant inflections, such as different verb tenses, in separate clusters. Being too eager tends to also merge word forms related by derivation rather than inflection, e.g. adjectives and their corresponding adverbs. We have tried to avoid such merges by using part-of-speech (POS) tags to further separate the word forms, which seemed promising with supervised POS. However, we do not want to rely on supervised POS in an unsupervised method, and we observed poor results once we moved to the (rather noisy) unsupervised POS of Mareček et al. (2016).

We are convinced that the proposed approach is still too weak, lacking the means to reliably separate true inflections from other similar word forms. For this, we believe, it will be necessary to actively look for regularities in the clusters to try to recognize inflectional paradigms in them, and then refine the clusters by encouraging them to match the paradigms.

As the assignment of word forms to clusters is context-independent, the approach also cannot deal with homonymy. We believe this could be solved by switching from context-independent word embeddings to contextual word embeddings,

Distance	Average	Median
$JW$	8.17	7.92
$cos$	4.39	3.87
$JW \cdot cos$	3.77	3.40

Table 3: Comparison of the word form similarities, in % of  $1 - vmeasure$  of the clustering. Average and median over the 28 datasets.

such as BERT (Devlin et al., 2018).

In Table 3, we compare the combined distance measure with each of the two components used alone. The results show that combining the edit distance with the embedding similarity is stronger than using any of the measures alone. In fact, only for two datasets,  $cos$  was slightly better than  $JW \cdot cos$ . The embedding similarity alone performs much better than the edit distance alone. We believe that this is at least partially due to Fast-Text embeddings’ incorporation of subword information, which means that their similarity also captures string similarity to some extent.

## 5 Conclusion

In this work, we have suggested an approach to unsupervised lemmatization based on agglomerative clustering, using a word form distance measure combining string similarity (edit distance) and meaning similarity (word embeddings). The evaluation showed the approach to be promising, surpassing the baseline on most of the evaluation datasets. At the same time, it has many obvious weak points that need to be addressed in future.

## References

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Matthew S. Dryer and Martin Haspelmath, editors. 2013. *WALS Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig.
- Edouard Grave, Piotr Bojanowski, Prakhara Gupta, Armand Joulin, and Tomas Mikolov. 2018. Learning word vectors for 157 languages. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)*.
- Daniel Kondratyuk, Tomáš Gavenčiak, Milan Straka, and Jan Hajič. 2018. Lemmatag: Jointly tagging and lemmatizing for morphologically rich languages with BRNNs. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing EMNLP 2018*, pages 4921–4928, Stroudsburg, PA,

USA. ACL's special interest group on linguistic data and corpus-based approaches to NLP, Association for Computational Linguistics.

Vladimir I Levenshtein. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady*, volume 10, pages 707–710.

Julie Beth Lovins. 1968. Development of a stemming algorithm. *Mech. Translat. & Comp. Linguistics*, 11(1-2):22–31.

David Mareček, Zhiwei Yu, Daniel Zeman, and Zdeněk Žabokrtský. 2016. [Deltacorporus 1.1](#). LINDAT/CLARIN digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.

Joakim Nivre et al. 2018. [Universal dependencies 2.3](#). LINDAT/CLARIN digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.

F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.

Martin F Porter. 2001. Snowball: A language for stemming algorithms.

William E. Winkler. 1990. [String comparator metrics and enhanced decision rules in the fellegi-sunter model of record linkage](#). In *Proceedings of the Section on Survey Research Methods (American Statistical Association)*, pages 354–359.