

Morfessor 2.0: Toolkit for Statistical Morphological Segmentation – Model

Probabilistic Model definition

- Full description in (Virpioja, 2012; Virpioja et al., 2013)

- Generative model

$$p(A, W | \theta)$$

analyses words parameters

The model generates pairs of words and analysis (the segmentation of a word into morphs)

- Tokenization function

$$a = \phi(w; \theta)$$

- Cost derivation

$$\theta_{\text{MAP}} = \arg \max_{\theta} p(\theta) p(D | \theta)$$

$$L(\theta, D) = -\log p(\theta) - \log p(D | \theta)$$

prior data likelihood

The data (D) is a list of (non-segmented) words to learn the model from in unsupervised manner.

Data Likelihood

$$\begin{aligned} \log p(D | \theta) &= \sum_{j=1}^N \log p(W = w_j | \theta) \\ &= \sum_{j=1}^N \log \sum_{a \in \Phi(w_j)} p(A = a | \theta), \end{aligned}$$

Morfessor Baseline assumes independence of words. Also, only valid tokenisations of need to be considered. Morfessor selects only one tokenisation (analysis) for each word at a time, by introducing a hidden variable Y .

$$\begin{aligned} \log p(D | \theta, Y) &= \sum_{j=1}^N \log p(y_j | \theta) \\ &= \sum_{j=1}^N \log p(m_{j1}, \dots, m_{j|y_j|}, \#_w | \theta) \end{aligned}$$

selected analysis

Prior

(Creutz and Lagus, 2007) The parameters of Morfessor Baseline encode the properties of the morph lexicon:

$$p(\theta) = p(\mu) \times \mu! \times p(\text{properties}(m_1), \dots, \text{properties}(m_{\mu})).$$

#morphs #morph permutations

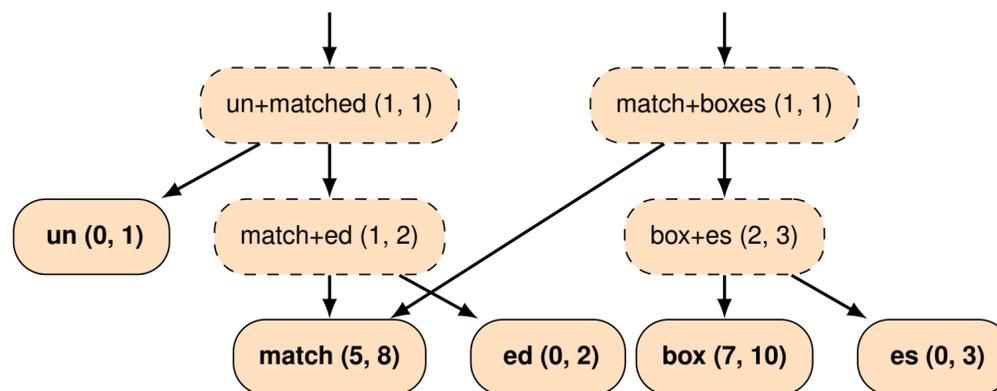
$$p(\sigma_i) = p(L = |\sigma_i|) \prod_{j=1}^{|\sigma_i|} p(C = \sigma_{ij})$$

morph length prior character distribution

Algorithm

(Creutz and Lagus, 2002)

```
function LOCALBATCHTRAIN(D, ε)
  θ, Y ← INITMODEL(DW)
  Lold ← ∞
  Lnew ← L(DW, θ, Y)
  while Lnew < Lold - ε do
    J ← RANDOMPERMUTATION(1, ..., N)
    for j ∈ J do
      θ, Y ← LOCALSEARCH(wj, D, θ, Y)
    Lold ← Lnew
    Lnew ← L(D, θ, Y)
  return θ, Y
```



unmatched, matchboxes, matched, boxes, match, and box

Likelihood weighting and Semi-supervised training

Likelihood weighting with α (Virpioja et al., 2011)

$$L(\theta, D) = -\log p(\theta) - \alpha \log p(D | \theta)$$

α can be determined in different ways, e.g using a development set, or some explicit knowledge like average morph length. Higher α reduces segmentation, lower α increases segmentation.

Semi-supervised (Kohonen et al., 2010)

$$L(\theta, D) = -\log p(\theta) - \alpha \log p(D | \theta) - \beta \log p(D_A | \theta)$$

annotated data

For semi-supervised learning another term is added to the cost, the likelihood of a set of annotations coming from the model. Also here a weight β is introduced to control the effect.

References

- Creutz, M. and Lagus, K. (2002). Unsupervised discovery of morphemes. In Maxwell, M., editor, *Proceedings of the ACL-02 Workshop on Morphological and Phonological Learning*, pages 21–30, Philadelphia, PA, USA. Association for Computational Linguistics.
- Creutz, M. and Lagus, K. (2007). Unsupervised models for morpheme segmentation and morphology learning. *ACM Transactions on Speech and Language Processing*, 4(1):3:1–3:34.
- Kohonen, O., Virpioja, S., and Lagus, K. (2010). Semi-supervised learning of concatenative morphology. In *Proceedings of the 11th Meeting of the ACL Special Interest Group on Computational Morphology and Phonology*, pages 78–86, Uppsala, Sweden. Association for Computational Linguistics.
- Virpioja, S. (2012). *Learning Constructions of Natural Language: Statistical Models and Evaluations*. PhD thesis, Aalto University.
- Virpioja, S., Kohonen, O., and Lagus, K. (2011). Evaluating the effect of word frequencies in a probabilistic generative model of morphology. In Pedersen, B. S., Nešpore, G., and Skadiņa, I., editors, *Proceedings of the 18th Nordic Conference of Computational Linguistics (NODALIDA 2011)*, volume 11 of *NEALT Proceedings Series*, pages 230–237. Northern European Association for Language Technology, Riga, Latvia.
- Virpioja, S., Smit, P., Grönroos, S., and Kurimo, M. (2013). Morfessor 2.0: Python implementation and extensions for Morfessor Baseline. Report 25/2013 in Aalto University publication series SCIENCE + TECHNOLOGY, Aalto University, Finland.

Morfessor 2.0: Toolkit for Statistical Morphological Segmentation – Codebase

Previous: Morfessor 1.0

- Written in perl
- Limited utf-8 support
- Older codebase, unsuitable for extensions

Morfessor 2.0

- Complete rewrite
- Inclusion of many previously published features
- Extensible for new features and algorithms

Usage

- Library interface
 - Directly use Morfessor from python scripts
- Command line interface
 - Run training evaluation and segmentation from the command line
 - Almost complete coverage of Morfessor functionality

Features

- On-line training
- Training speed-up with random skips
- Frequency threshold and dampening for words in training data
- Possibility to weight training data likelihoods
- Optimization of data likelihood weight based on development set
- Optimization of data likelihood weight based on average morph type length
- Optimization of data likelihood weight based on average morphs / word

Items in grey will be released in Morfessor 2.1

Implementation

- Python
 - Runs on Python2, Python3 and Pypy interpreters
 - Best performance: Pypy
- Unit-agnostic code
 - Split words into morphs, or sentences into phrases

Demo

- Build on top of library interface
- Will be available online in future

Distribution

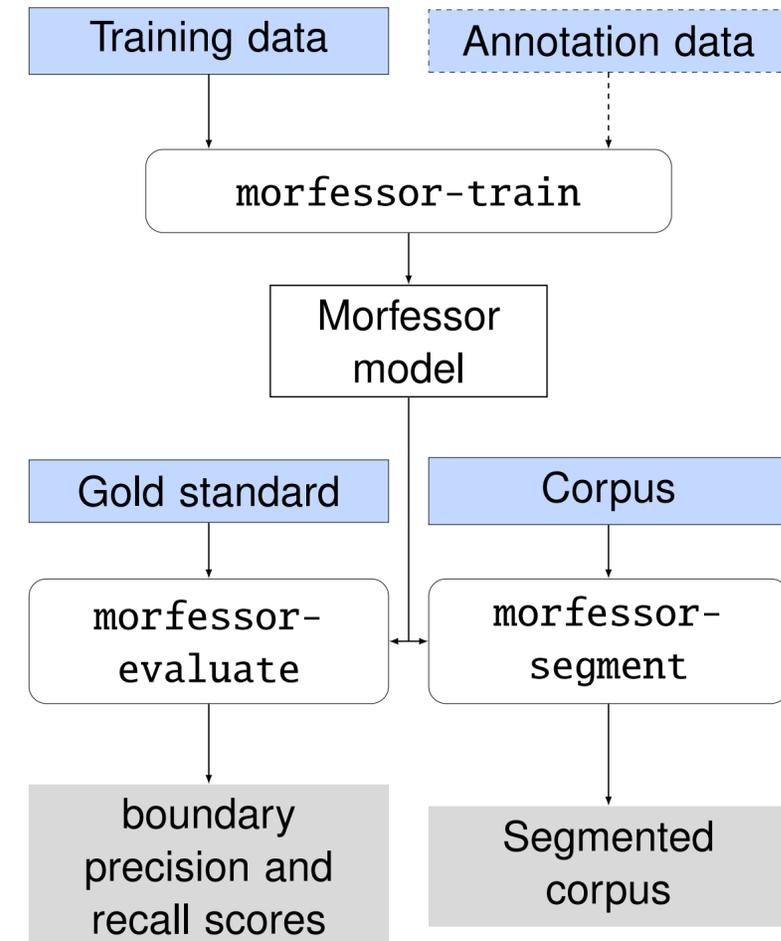
- Source code available at GitHub
- Packages available from the PYthon Package Index (Pypi)
- BSD Open Source License

Upcoming developments

- Morfessor 2.0 based model with morphotactic constraints (categories)
- Data selection techniques

Acknowledgements

The authors have received funding from the EC's 7th Framework Programme (FP7/2007–2013) under grant agreement n°287678 and the Academy of Finland under the Finnish Centre of Excellence Program 2012–2017 (grant n°251170) and the LASTU Programme (grants n°256887 and 259934). The experiments were performed using computer resources within the Aalto University School of Science "Science-IT" project.



Links

- Homepage
<http://www.cis.hut.fi/projects/morpho/>
- GitHub
<https://github.com/aalto-speech/morfessor>
- Pypi
`pip install morfessor`
<https://pypi.python.org/pypi/Morfessor>
- Documentation
<http://morfessor.readthedocs.org/en/latest/>

Morfessor 2.0: Toolkit for Statistical Morphological Segmentation – Related Projects

Aalto University - Morpho project

The work on Morfessor is funded through multiple projects. The work on Morfessor 2.0 were done in the scope projects mentioned below.

For updates on Morfessor, go to <http://www.cis.hut.fi/projects/morpho/> and subscribe to the mailing list!

COIN



CompuBrain: Computational modelling of brain's language

A project between O.V. Lounasmaa Laboratory and Departments of Information and Computer Science and Signal Processing and Acoustics at Aalto University, funded by Academy of Finland. Aims for brain-based models of language by bringing together expertise on neuroimaging and computational modelling of language.

Morpho challenge 2014

After the challenges in 2005 and 2007-2010 Aalto University plans to organize a new Morpho challenge.

Possible Tasks:

- Generation of new words for language modeling and spelling correction
- Bilingual morphological analysis for machine translation.

Starts in 2014, as soon as task and data are ready. More information <http://research.ics.aalto.fi/events/morphochallenge/>. Subscribe to the mailing list!

Simple4All - EU Framework 7 Programme

