

Filling the Gap: Decoding of Word Embeddings for Generation of Coherent New Words

Safa AlSaidi, Amandine Decker, Stephanie Monteiro

M2 — Software Project



1 Towards building a decoder

2 Our Encoder

3 Method

4 Our Approach

5 Future Leads

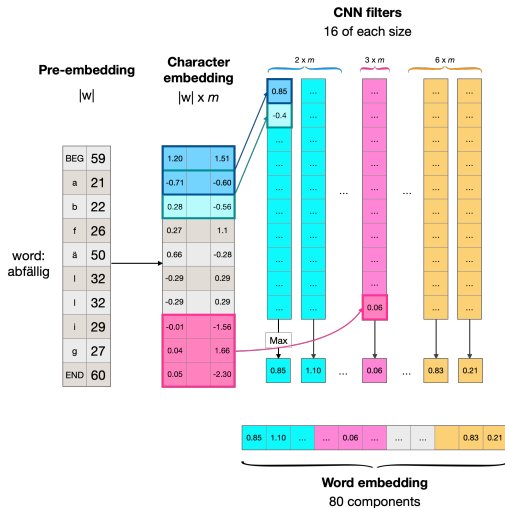
6 Applicability

What we want to work on

- ❶ Build a decoder that transforms *embedded vectors* to *words*
- ❷ Work on morphological inflection/derivation
 - ▶ root, affixes: prefixes, suffixes..
- ❸ Apply it to 11 different languages
 - ▶ based on two datasets [Cotterell et al., 2016, Karpinska et al., 2018]
 - ▶ improve regression task
 - $A : B :: C : X \xrightarrow{X=?} A : B :: C : D$

- 1 Towards building a decoder
- 2 Our Encoder**
- 3 Method
- 4 Our Approach
- 5 Future Leads
- 6 Applicability

Pre-existing embedding model



This model was inspired by [Kim et al., 2016]

- 1 Towards building a decoder
- 2 Our Encoder
- 3 Method**
- 4 Our Approach
- 5 Future Leads
- 6 Applicability

Method

Datasets:

- SIGMORPHON 2016 [Cotterell et al., 2016] and the Japanese Bigger Analogy Test Set [Karpinska et al., 2018].

Tools:

- Python libraries (mainly PyTorch)
- Grid'5000

Our repository on Github:

https://github.com/Safa-98/morphological_embeddings_decoder

- 1 Towards building a decoder
- 2 Our Encoder
- 3 Method
- 4 Our Approach**
- 5 Future Leads
- 6 Applicability

Our Approach

- ❶ Initial intuition: find a method to reverse the embedding model
- ❷ Usual decoders: RNN, GRU, LSTM
- ❸ One possible approach: *Teacher Forcing*
 - ▶ <https://rajatvd.github.io/Generating-Words-From-Embeddings/>

- 1 Towards building a decoder
- 2 Our Encoder
- 3 Method
- 4 Our Approach
- 5 Future Leads**
- 6 Applicability

Future Leads

- ❶ Auto-encoder
- ❷ Multilingual decoder/auto-encoder
- ❸ Decoder based on subwords/morphemes
 - ▶ <https://github.com/colingoldberg/morphemes>

- 1 Towards building a decoder
- 2 Our Encoder
- 3 Method
- 4 Our Approach
- 5 Future Leads
- 6 Applicability**

Where can it be applied?

- Analogy solving (educational application)
 - ▶ “play”:“build”::“replay”: $X \rightarrow X =$ “rebuild”
- General structure: apply to other types of data
- Multilingual model: underrepresented languages

شكرا جزيلا




Thank you

Merci

អរគុណ

Obrigado

References I

-  Cotterell, R., Kirov, C., Sylak-Glassman, J., Yarowsky, D., Eisner, J., and Hulden, M. (2016).
The sigmorphon 2016 shared task—morphological reinflection.
In *Proceedings of the 2016 Meeting of SIGMORPHON*, Berlin, Germany.
Association for Computational Linguistics.
-  Karpinska, M., Li, B., Rogers, A., and Drozd, A. (2018).
Subcharacter Information in Japanese Embeddings: When Is It Worth It?
In *Proceedings of the Workshop on the Relevance of Linguistic Structure in Neural Architectures for NLP*, pages 28–37, Melbourne, Australia.
Association for Computational Linguistics.
-  Kim, Y., Jernite, Y., Sontag, D., and Rush, A. (2016).
Character-aware neural language models.
In *Proceedings of the AAAI Conference on Artificial Intelligence*, 30(1).