

Partial Colexifications Improve Concept Embeddings

Arne Rubehn & Johann-Mattis List

Chair for Multilingual Computational Linguistics, University of Passau, Germany



Introduction

We learn meaningful, low-dimensional representations of concepts by applying graph embedding techniques on colexification networks.

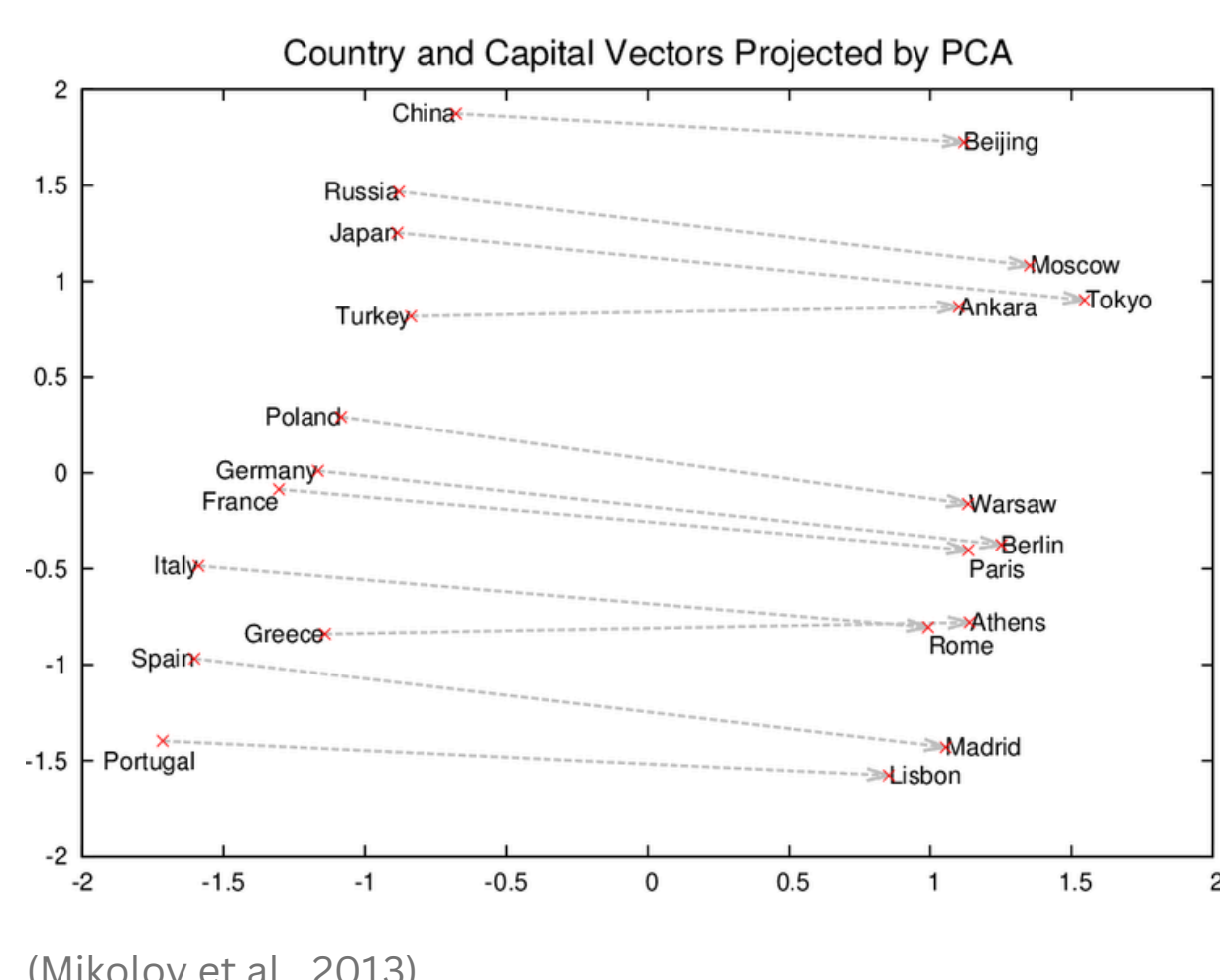
Motivation

- State-of-the-art methods in computational historical linguistics are generally unable to model or reconstruct semantic change
- It is notoriously hard to quantify meaning cross-linguistically
- Modeling semantic relations between *concepts* reliably in a “computer-friendly” way could help bridging the gap

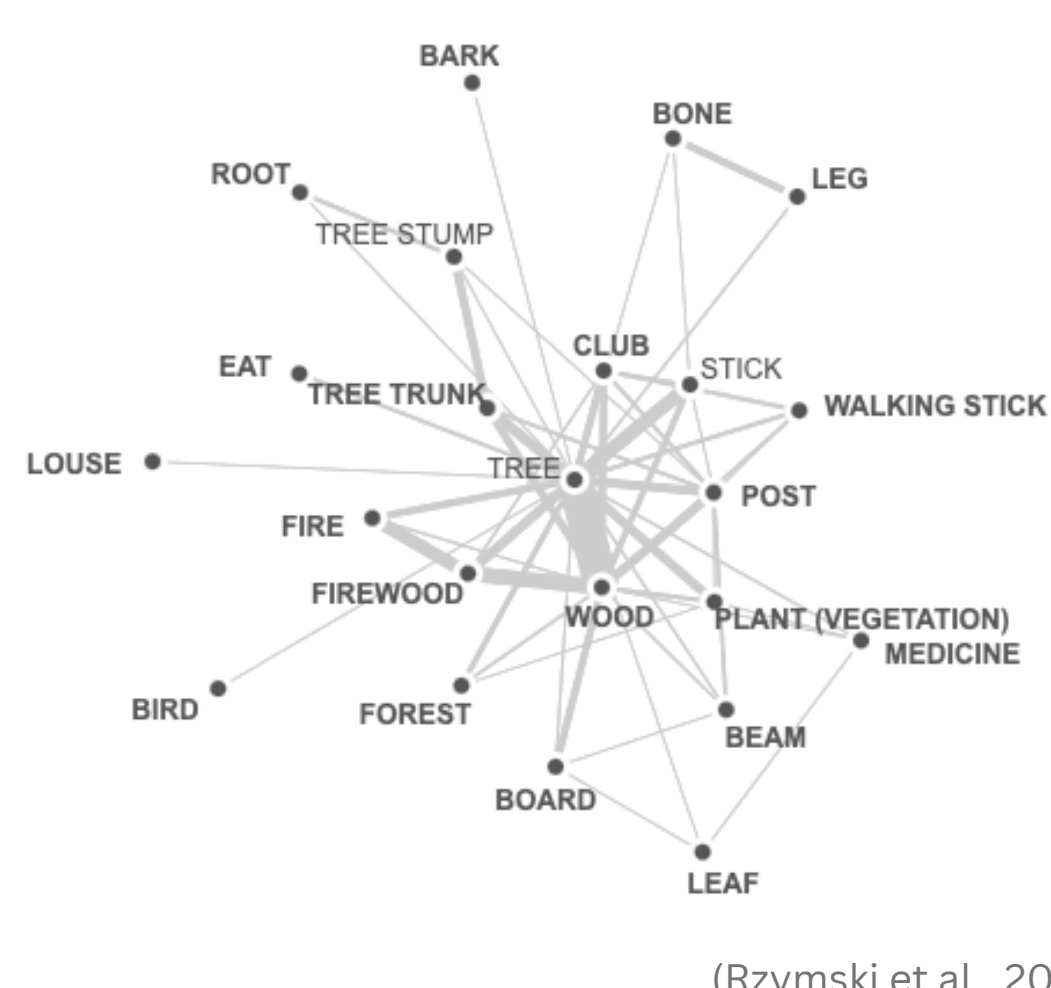
Background

- Word embeddings** have shown that various semantic relationships can be efficiently modeled in a low-dimensional vector space.
 - ...but since they embed words, they are not fit for wide cross-linguistic applications
- Colexification networks** offer a cross-linguistic, concept-based approach on semantics, encoding common pathways of semantic shift
 - ...but their network structure can not readily be processed by downstream applications

Using graph embedding techniques, we can learn embedded representations for concepts in colexification networks.



(Mikolov et al., 2013)



(Rzysmski et al., 2020)

Materials and Methods

Colexification Data

- Three types of colexification (List, 2023) inferred from the *Intercontinental Dictionary Series* (Key and Comrie, 2016):
 - full colexification
 - affix colexification
 - overlap colexification
- Colexification network for each type of colexification
- Concepts defined by the CLLD Concepticon (List et al., 2025)

Graph Embedding Techniques

- ProNe (Zhang et al., 2019), Node2Vec (Grover and Leskovec, 2016) & SDNE (Wang et al., 2016)

Training

- Train individual embeddings for each colexification network
- Combine embeddings from different colexification types as a post-processing step (concatenation + PCA)

Experiments

1) Modeling Lexical Semantic Similarity

- Multilingual similarity ratings between word pairs, obtained from MultiSimLex (Vulić et al., 2020)
- Calculate cosine similarities between corresponding concept pairs
- Calculate Spearman’s r between similarities

2) Predicting Semantic Change

- Obtain historically attested semantic changes from DatSemShift (Zalizniak et al., 2024)
- Negative sampling of random “changes”
- Train simple logistic regression classifier to tell apart true and false shifts, based on the respective concept similarities

3) Predicting Word Associations

- Obtain word association data from the Edinburgh Association Thesaurus (Kiss et al., 1973)
- Same experimental setup as previous task: Negative sampling, prediction with simple logistic regression classifier

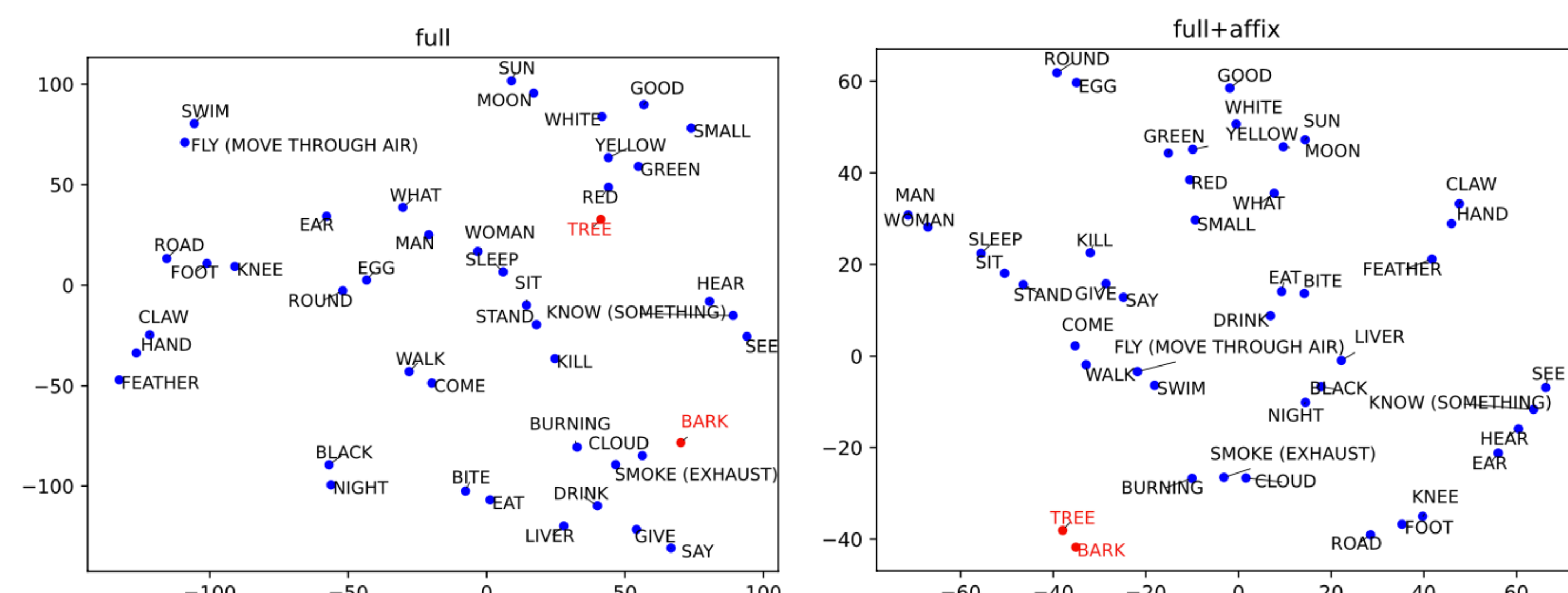
Baselines

- Various similarity metrics inferred from the graph directly
- Multilingual fastText vectors (Grave et al., 2018)

Results & Discussion

Results

- Embeddings learned with ProNE perform the best on average, closely followed by Node2Vec; SDNE does not seem viable
- Concept embeddings outperform graph-based baselines in all three tasks, and fastText embeddings in 2 of 3 tasks
- Almost identical patterns between tasks 2 and 3
- Affix colexifications lead to better embeddings on all three tasks
- Overlap colexifications are beneficial for predicting semantic change and word association, but detrimental for modeling lexical similarity (the same pattern holds for fastText embeddings!)



Discussion

- Partial colexifications can capture semantic relations that are rarely expressed by full colexification
- Affix colexification seems to capture a more direct relationship between concepts than overlap colexification

Enriching embeddings with partial colexification data always leads to better results than relying on full colexification data alone!

