

Tackling fairness, change and polysemy in word embeddings

Felipe Bravo-Marquez

Department of Computer Science, University of Chile
National Center for Artificial Intelligence Research
Millennium Institute Foundational Research on Data

November 29, 2021



dcc
CIENCIAS DE LA COMPUTACIÓN
UNIVERSIDAD DE CHILE

CEN
CENTRO NACIONAL DE INTELIGENCIA ARTIFICIAL



Millennium Institute
Foundational
Research on Data

RELELA
Representations for
Learning and Language

Word Vectors

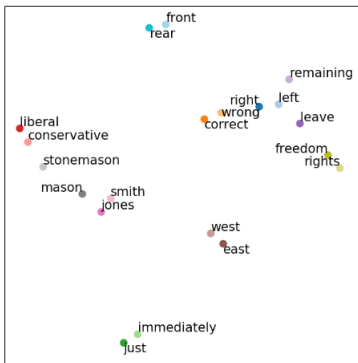
- A major component in neural networks for language is the use of an embedding layer.
- A mapping of discrete symbols to continuous vectors.
- When embedding words, they transform from being isolated distinct symbols into mathematical objects that can be operated on.
- Distance between vectors can be equated to distance between words.
- This makes easier to generalize the behavior from one word to another.

Distributional Vectors

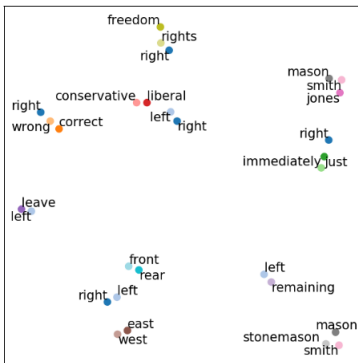
- **Distributional Hypothesis:** words occurring in the same **contexts** tend to have similar meanings.
- Or equivalently: “a word is characterized by the **company** it keeps”.
- In this talk we summarize our research addressing three limitations of static word embeddings: 1) fairness, 2) semantic change, and 3) polysemy.

PolyLM: a polysemous language model

- A language model capable of automatically learning multiple meanings of a word (e.g. apple:apple, apple:company) [Ansell et al., 2021].

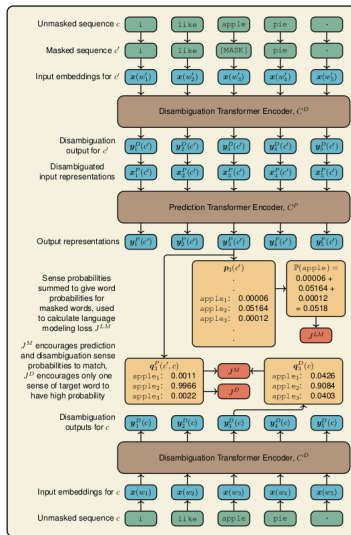


(a) Word embeddings



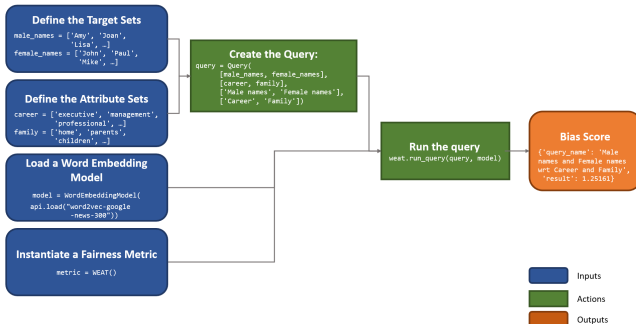
(b) Sense embeddings

PolyLM: a polysemous language model



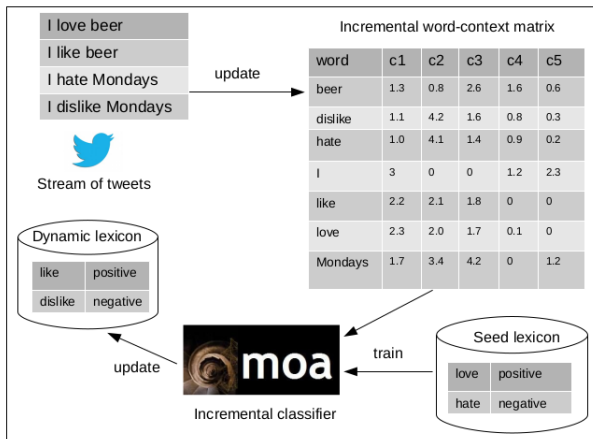
WEFE: The Word Embeddings Fairness Evaluation Framework

- The Word Embeddings Fairness Evaluation (WEFE) is a framework for measuring and mitigating bias in word embeddings (e.g. man is to programmer as woman is to housewife). [Badilla et al., 2020].



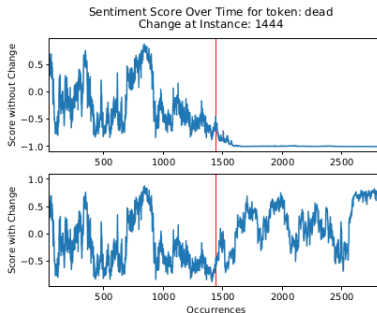
Incremental Word Vectors

- An algorithm capable of continuously learning word vectors and thus understanding how the meaning evolves over time (e.g., monitoring the word “estallido” in social networks during the Chilean social unrest).
[Bravo-Marquez et al., 2021].



Incremental Word Vectors

- We simulate sentiment change by randomly picking some words and swapping their context with the context of words exhibiting the opposite sentiment.



1. (a) dead

- Our approach allows for successfully tracking of the sentiment of words over time even when drastic change is induced.

Thanks for your Attention!

References I



Ansell, A., Bravo-Marquez, F., and Pfahringer, B. (2021).
Polylm: Learning about polysemy through language modeling.
In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 563–574.



Badilla, P., Bravo-Marquez, F., and Pérez, J. (2020).
Wefe: The word embeddings fairness evaluation framework.
In IJCAI, pages 430–436.



Bravo-Marquez, F., Khanchandani, A., and Pfahringer, B. (2021).
Incremental word vectors for time-evolving sentiment lexicon induction.
Cognitive Computation, pages 1–17.