# Affective meaning in colexification networks and applications to automatic lexica expansion

Anna Di Natale[1,2,3,*], Max Pellert[4] & David Garcia[5,1,2,3]

[1] Graz University of Technology [2] Medical University of Vienna [3] Complexity Science Hub Vienna
[*] dinatale@csh.ac.at [4] University of Mannheim [5] University of Konstanz

**Keywords**: Cross-linguistic colexifications, affective meaning, semantic networks

Machine learning methods have proven to be very effective in the analysis of human language and in understanding emotion expression. However, in order to improve AI and connect it to social sciences, we need to explore the meaning structure of human language and the applications of linguistic properties to NLP. With this aim, we analyse the relationship between colexification occurrences and meaning similarity in the domain of affective meaning, and we develop a method to automatically infer the valence, arousal and dominance of words. We find robust evidence that affective meaning is encoded in colexification networks and propose a tool that improves state of the art methods to infer the affective ratings of words. This constitutes a first step towards explainable and theory-based methods for text analysis and for the automatic expansion of affective science resources.

Colexification is a linguistic phenomenon that occurs when multiple concepts are expressed by the same word in a language (François, 2008). The collection of colexification occurrences can be shaped in the form of a network, where nodes represent concepts and edges track colexification occurrences between pairs of concepts. Edges in the network are weighted according to the number of languages and of language families that present a colexification between the same pair of concepts. Colexification patterns are believed to be determined by semantic relationships between concepts (François, 2008), thus colexification networks should also encode the semantics of concepts. In this work, we test this hypothesis in the field of affective meaning and explore applications to the automatic expansion of preexisting affective lexica in an unsupervised, theory-driven and explainable way.

To this end, we consider two English affective norms lexica (Mohammad, 2018; Warriner et al., 2013) and three colexification networks: CLICS[3] (N=1,647) (Rzymski et al., 2020) and two networks built from crowdsourced translations (OmegaWiki[1], N=10,323) and from open source bilingual dictionaries (FreeDict[2], N=27,939). We deploy the network structure to estimate the affective meaning of words according to the three dimensions of valence, arousal, and dominance in the following way. First, we map words of the affective lexicon to nodes in the network, as shown in Figure 1 left panel in the case of the OmegaWiki network. We then estimate the affective meaning of a node as the weighted mean of the affective meaning of its neighbors. We do so for the three dimensions of valence, arousal and dominance. Although this method is a simple, unsupervised computation, it reaches high correlation between the true affective rating and the predicted one, as represented in the case of OmegaWiki and the affective dataset (Mohammad, 2018) in Figure 1 right panel. Indeed, in this case the correlation between computed and true valence is significative and high ($\rho$ = 0.839). Outliers highlighted in Figure 1 right panel give further insights into the annotation procedure for ground truth data and into the cultural component of such annotations and of affective meaning in general.

We evaluate how colexification networks predict the affective norms of words that do not belong to the affective lexica with 10 repetitions of a 75/25 split cross validation as in (Mandera et al., 2015). We find high correlation coefficients between our estimates and the empirical values, which are comparable with and in some cases outperform machine learning methods on large corpora (Mandera et al., 2015). Our results also present higher coverage of the semantic space than state-of-the-art methods, that is our method can estimate the affective ratings of a higher number of words than what could be previously achieved with word embeddings (Mandera et al., 2015).

The results of this work provide strong support to the hypothesis that colexification occurrences captures meaning similarity between words and that this property is also embedded in colexification networks. Furthermore, our analysis shows that word semantics can be interpolated with colexification networks and that the unsupervised expansion of already existing lexica is possible. This practice has the potential to lower the costs of lexica creation, which usually requires a study to be designed, and a group of non-expert participants to be recruited. Moreover, the resulting algorithm is fully explainable, thus its results can be analysed with respect to, for example, cultural differences in the understanding of emotions.

---

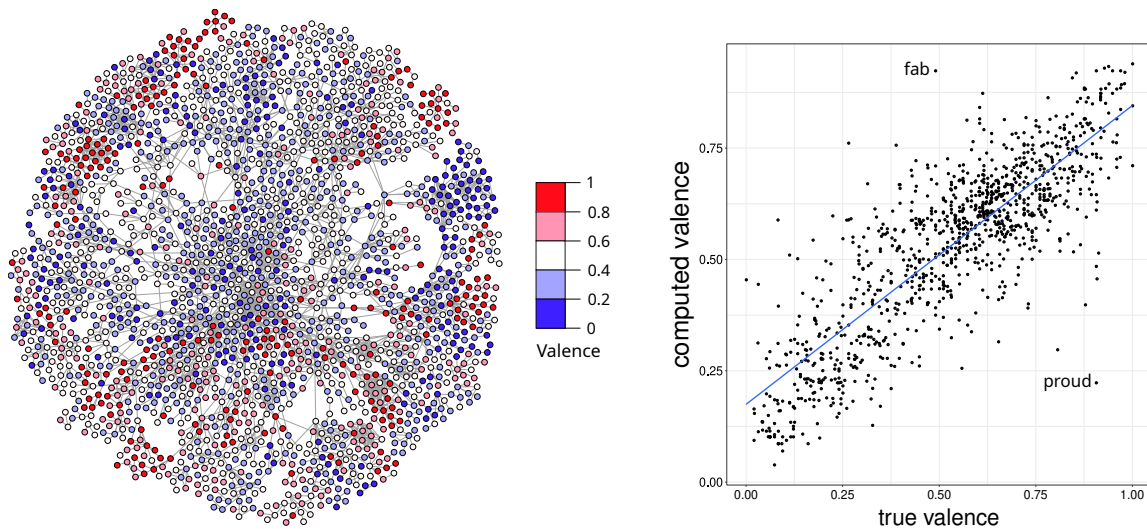[1]    http://www.omegawiki.org
[2]    http://www.freedict.org

*Fig. 1: Left panel: Words from (Mohammad, 2018) mapped in the OmegaWiki colexification network colored according to their valence (low valence, i.e. positive words are colored in blue; high valence, negative words in red). Words cluster according to their valence in the network. Right panel: Correlation of the computed and true valence ratings on the database (Mohammad, 2018) ($\rho = 0.839$, c.i.=$[0.82, 0.856]$, $p < 0.001$). Outliers are labelled.*

## References

François, Alexandre. 2008. Semantic maps and the typology of colexification. *From polysemy to semantic change: Towards a typology of lexical semantic associations* 106. 163.

Mandera, Paweł, Emmanuel Keuleers & Marc Brysbaert. 2015. How useful are corpus-based methods for extrapolating psycholinguistic variables? *The Quarterly Journal of Experimental Psychology* 68(8).

Mohammad, Saif. 2018. Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 english words. In *Proceedings of the 56th annual meeting of the association for computational linguistics (volume 1: Long papers)*, 174–184.

Rzymski, Christoph, Tiago Tresoldi, Simon J Greenhill, Mei-Shin Wu, Nathanael E Schweikhard, Maria Koptjevskaja-Tamm, Volker Gast, Timotheus A Bodt, Abbie Hantgan, Gereon A Kaiping et al. 2020. The database of cross-linguistic colexifications, reproducible analysis of cross-linguistic polysemies. *Scientific data* 7(1). 1–12.

Warriner, Amy Beth, Victor Kuperman & Marc Brysbaert. 2013. Norms of valence, arousal, and dominance for 13,915 english lemmas. *Behavior research methods* 45(4). 1191–1207.