

Automatic Detection of Copredication using Contextualized Word Embeddings

Deniz Ekin Yavaş¹, Marta Ricchiardi², Elisabetta Ježek³, Laura Kallmeyer¹,
and Rainer Osswald¹

¹Heinrich Heine University Düsseldorf

^{2,3}University of Pavia

¹{deniz.yavas, laura.kallmeyer, rainer.osswald}@hhu.de

²marta.ricchiardi01@universitadipavia.it

³jezek@unipv.it

Introduction. Copredication is a phenomenon that has been explored in formal linguistics in detail, as it is a construction important for the investigation of complex types and compositionality in natural languages. We believe that this phenomenon requires more attention from an empirical point of view, which can also eventually help the development of theoretical frameworks. For this purpose, we propose a method for automatically extracting sentences of complex type nouns with copredication from corpus data.

Copredication in general can be described as a “grammatical construction in which two predicates jointly apply to the same argument” (Asher, 2011, p. 11). A particularly interesting type of copredication concerns constructions where the two predicates require different semantic types for their arguments but apply to the same noun. This is possible if the noun is a complex type noun, whose different meaning facets can be accessed simultaneously (Pustejovsky, 1995; Asher and Pustejovsky, 2006; Pustejovsky and Ježek, 2008; Asher, 2011).¹ Copredications in this more narrow sense can involve different syntactic constructions, for instance verbs taking a noun as direct object or adjectives modifying a noun. In our study, we focus on cases of copredication that combine the two constructions just mentioned. We refer to this combination as the *verb+adj pattern*.

Ježek and Vieu (2014) adopt a semi-automatic approach for extracting verb+adj copredications for Information•Physical_Object complex type nouns. They manually construct copredication contexts with different predicate combinations, and then extract examples by searching the corpus for these contexts

Methodology. In this study, we pursue an approach that is based on exploiting the implicit knowledge in contextual language models. It only requires labeled data for the individual predications. As a proof of concept, we focus on verb+adj copredications with Event•Food nouns in Italian, as in (1). The complex type Event•Food has been frequently discussed in the literature (Pustejovsky, 1995, 1998; Asher and Pustejovsky, 2006; Pustejovsky and Ježek,

¹Different terms are used in the literature for complex types: “dot objects” (Pustejovsky, 1995, 1998), “nouns with facets” (Cruse, 1995), “dual aspect nouns” (Asher, 2011). The different meanings of a complex type are also called “facets” (Cruse, 1995) or “aspects” (Asher, 2011).

2008; Asher, 2011). Standard examples are “lunch” and “dinner”, which can refer both to an event and to food. This complex type has been selected because of the availability of data for the simple types Food and Event in the resource we used (see below).

- (1) a. Concludemmo la sostanziosa **colazione** con delle fette di dolce cocomero.
 Event Verb Food Adj
‘We concluded the substantial breakfast with slices of sweet watermelon.’
- b. Consumiamo una veloce **colazione** e poi via verso l’aeroporto.
 Food Verb Event Adj
‘We have a quick breakfast and then off to the airport.’

Our method comprises two steps: 1) training classifiers that label individual predications and 2) combining them to detect copredication.

Individual classifiers. For the first step, for each predication type (verb/adj) and for each meaning facet (event/food), we trained a binary classifier that labels the relation between a predicate of this type and its argument with respect to whether the specific facet is addressed in this predication or not. For instance, the classifier for verbs (where the noun is the direct object) and event would label (2a) as true, and the classifier for adjectives and food would label (2b) as true.

- (2) a. Arrestò la sua **corsa** davanti al portone del carcere e bussò. [EVENT]
 Event Verb
‘He stopped his run in front of the prison door and knocked.’
- b. Soffriggere le **cipolle** fresche fino a che non risultano dorate. [FOOD]
 Food Adj
‘Fry the fresh onions until golden brown.’

The contextualized embeddings for the predicate (verb or adjective) and the noun are extracted from one of the pre-trained dbmdz BERT models for Italian available at Hugging Face.² These embeddings are then concatenated and the resulting vector is classified by means of SVM classifiers.

Training data. For training the verbal classifiers, we used T-PAS (Typed Predicate Argument Structures; Ježek et al. 2014), a resource of argument structures of Italian verbs. T-PAS provides corpus-derived argument structures with manually annotated semantic argument types, e.g., [Human] mangiare [Food] (*Eng.: [Human] eat [Food]*), together with matching corpus instances. From these data, we used sentences whose verbs take direct objects of type Food (resp. Event) as training data. The sentences were parsed with the spacy-udpipe python library in order to identify the verb and the direct object.³

The training data for the adjectival classifiers is generated by employing masked language modeling with BERT. Starting from the verbal predication data we let BERT insert an adjective modifying the direct object. The assumption is that in sentences where the verbal predication over the objects targets a certain facet, the inserted adjective will do so as well. These adjective-direct object pairs are used for training the adjectival predication classifiers. In order to test the performance of adjectival predication classifiers, we used a test set that is not automatically generated by the model. The test set is created by extracting corpus instances from the ItWac corpus (Baroni and Kilgarriff, 2006), identified through a concordance search for the most typical 5-10 lexical items that express each type in corpus instances and their respective most

²<https://huggingface.co/dbmdz/bert-base-italian-cased>

³Although typed predicate argument structures are induced from corpora by manual clustering corpus instances in which the verb has the same sense, sentence-level annotations are not provided in T-PAS.

frequent adjective modifiers.

For each data set, negatives are selected from the semantic types other than the target type, and the negatives are downsized using clustering-based undersampling to make their size equal to the positive samples in order to create a balanced data set. See Table 1 for some statistics of the training and test data.

Evaluation. All 4 classifiers perform very well on the test data (see Table 1), with accuracy scores of over 90%, despite the low number of training data for some classifiers and despite the fact that the adjectival training data is automatically generated.

Detecting copredication. For the second step, we apply the classifiers for both facets on the verb/adj predications and check whether they predict food for the verb (resp. the adj) and event for the adj. (resp. the verb). The performance of the method is measured on a manually created test suite. The test suite contains 30 positive and 45 negative examples for copredication concerning Event•Food. In the positive examples, verb and adjective target different facets of a complex noun (as in (1)), while in the negative examples, they target the same facet of a complex type noun (as in (3a)) or just a simple type noun (as in (3b)).

(3) a. Sono qui pronta che attendo il **pranzo** imminente, sono abbastanza agitata.

Event Verb Event Adj

'I am here ready waiting for the upcoming lunch, I am quite nervous.'

b. Sono qui pronta che attendo il **colloquio** imminente, sono abbastanza agitata.

Event Verb Event Adj

'I am here ready waiting for the upcoming interview, I am quite nervous.'

Evaluation. Our method achieves a score of around 80% for both positives (*recall*) and negatives (*specificity*), as shown in Table 2. However, we have to distinguish between negative cases with a simple type noun and negative cases with a complex type noun, since our goal is to automatically extract copredication examples with complex type nouns. The results show that our method is more problematic in detecting negative cases with complex type nouns: Even though it detects negative cases with simple type nouns with 100% accuracy, the score drops significantly to 66% for the negative cases with complex type nouns, as seen in Table 2.

Note that for some examples in our negative set, there are two copies of the same sentence, one with a simple type noun, and the other with a complex type noun and only the noun is different in these cases, as in (3a) compared to (3b). These examples allow a more direct comparison between simple type and complex type nouns in terms of their effects on the classification. This comparison shows that the presence of a complex type noun is sometimes enough to lead classifiers to incorrectly detect copredication. For example, in (3a) and (3b), even though the sentences have two predications of the same type (Event), copredication is detected in the sentence with the complex type noun (3a), but not in the sentence with the simple type noun (3b).

Discussion. This study demonstrates that the semantic classification of different syntactic types of predications is possible with contextualized embeddings even with a limited amount of training data. However the detection of copredication, even though overall results are promising, is less straightforward, because of the tendency of false positives with complex type nouns. This might be due to a bias in the training data, and to a general distinction between simple type and complex type nouns made by BERT, independent from the specific context. We will investigate this issue further in the future. Furthermore, we will extend the approach to other complex types, other predication constructions and other language models. Finally, we plan to use this method for building a collection of corpus-based copredication instances that will provide a broad empirical basis for the qualitative and formal analysis of copredication phenomena.

Appendix: Tables

	Classifiers			
	Food Verb	Food Adj.	Event Verb	Event Adj.
Scores				
Accuracy	0.984	0.997	0.971	0.902
Data Size				
Pos/Neg (Training)	261/261	134/134	1407/1407	1044/1044
Pos/Neg (Test)	129/129	660/660	693/693	838/838

Table 1: Training and evaluation information of individual classifiers

Recall	Specificity		
	overall	over simple nouns only	over complex nouns only
0.8	0.82	1	0.66

Table 2: Evaluation results of copredication detection on test suite

References

- Asher, N. (2011). *Lexical Meaning in Context. A Web of Words*. Cambridge University Press, Cambridge.
- Asher, N. and Pustejovsky, J. (2006). A type composition logic for generative lexicon. *Journal of Cognitive Science*, 6:1–38.
- Baroni, M. and Kilgarriff, A. (2006). Large linguistically-processed web corpora for multiple languages. In *EACL'06: Proceedings of the Eleventh Conference of the European Chapter of the Association for Computational Linguistics: Posters & Demonstrations; 2006 Apr 5-6; Trento, Italy. Stroudsburg (PA): Association for Computational Linguistics; 2006. p. 87-90*. ACL (Association for Computational Linguistics).
- Cruse, D. A. (1995). Polysemy and related phenomena from a cognitive linguistic viewpoint. *Computational lexical semantics*, pages 33–49.
- Ježek, E. and Vieu, L. (2014). Distributional analysis of copredication: towards distinguishing systematic polysemy from coercion. In *Proceedings of the First Italian Conference on Computational Linguistics (CLiC-it)*, pages 219–223, Pisa.
- Ježek, E., Magnini, B., Feltracco, A., Bianchini, A., and Popescu, O. (2014). T-PAS: A resource of corpus-derived Typed Predicate Argument Structures for linguistic analysis and semantic processing. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 890–895.
- Pustejovsky, J. (1995). *The Generative Lexicon*. MIT Press, Cambridge, MA.
- Pustejovsky, J. (1998). The semantics of lexical underspecification. *Folia Linguistica*, 32(3/4):323–348.
- Pustejovsky, J. and Ježek, E. (2008). Semantic coercion in language: Beyond distributional analysis. *Italian Journal of Linguistics*, 20(1):181–214.