

OP-IMS @ DIACR-Ita: Back to the Roots: SGNS+OP+CD still Rocks Semantic Change Detection

Jens Kaiser, Dominik Schlechtweg, Sabine Schulte im Walde
Institute for Natural Language Processing, University of Stuttgart
{jens.kaiser, schlecdk, schulte}@ims.uni-stuttgart.de

Abstract

We present the results of our participation in the DIACR-Ita shared task on lexical semantic change detection for Italian. We exploit one of the earliest and most influential semantic change detection models based on Skip-Gram with Negative Sampling, Orthogonal Procrustes alignment and Cosine Distance and obtain the winning submission of the shared task with near to perfect accuracy (.94). Our results once more indicate that, within the present task setup in lexical semantic change detection, the traditional type-based approaches yield excellent performance.

1 Introduction

Lexical Semantic Change (LSC) Detection has drawn increasing attention in recent years (Kutuzov et al., 2018; Tahmasebi et al., 2018). Recently, SemEval-2020 Task 1 provided a multilingual evaluation framework to compare the variety of proposed model architectures (Schlechtweg et al., 2020). The DIACR-Ita shared task extends parts of this framework to Italian by providing an Italian data set for SemEval’s binary subtask (Basile et al., 2020a; Basile et al., 2020b).

We present the results of our participation in the DIACR-Ita shared task exploiting one of the earliest and most established semantic change detection models based on Skip-Gram with Negative Sampling, Orthogonal Procrustes alignment and Cosine Distance (Hamilton et al., 2016a). Based on our previous research (Schlechtweg et al., 2019; Kaiser et al., 2020) we optimize the dimensionality parameter assuming that high dimensionalities reduce alignment error. With our

setting win the shared task with near to perfect accuracy (.94). Our results once more demonstrate that, within the present task setup in lexical semantic change detection, the traditional type-based approaches yield excellent performance.

2 Related Work

As evident in Schlechtweg et al. (2020) the field of LSCD is currently dominated by Vector Space Models (VSMs), which can be divided into type-based (Turney and Pantel, 2010) and token-based (Schütze, 1998) models. Prominent type-based models include low-dimensional embeddings such as the Global Vectors (Pennington et al., 2014, GloVe) the Continuous Bag-of-Words (CBOW), the Continuous Skip-gram as well as a slight modification of the latter, the Skip-gram with Negative Sampling model (Mikolov et al., 2013a; Mikolov et al., 2013b, SGNS). However, as these models come with the deficiency that they aggregate all senses of a word into a single representation, token-based embeddings have been proposed (Peters et al., 2018; Devlin et al., 2019). According to Hu et al. (2019) these models can ideally capture complex characteristics of word use, and how they vary across linguistic contexts. The results of SemEval-2020 Task 1 (Schlechtweg et al., 2020), however, show that contrary to this, the token-based embedding models (Beck, 2020; Kutuzov and Giulianelli, 2020) are heavily outperformed by the type-based ones (Pražák et al., 2020; Asgari et al., 2020). The SGNS model was not only widely used, but also performed best among the participants in the task. Its fast implementation and combination possibilities with different alignment types further solidify SGNS as the standard in LSCD. A common and surprisingly robust (Schlechtweg et al., 2019; Kaiser et al., 2020) practice is to align the time-specific SGNS embeddings with Orthogonal Procrustes (OP) and measure change with Cosine Distance (CD) (Kulka-

“Copyright © 2020 for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).”

rni et al., 2015; Hamilton et al., 2016b). This has been shown in several small but independent experiments (Hamilton et al., 2016b; Schlechtweg et al., 2019; Kaiser et al., 2020; Shoemark et al., 2019) and SGNS+OP+CD has produced two of three top-performing submissions in Subtask 2 in SemEval-2020 Task 1 including the winning submission (Pömsl and Lyapin, 2020; Arefyev and Zhikov, 2020).

3 System overview

Most VSMs in LSC detection combine three subsystems: (i) creating semantic word representations, (ii) aligning them across corpora, and (iii) measuring differences between the aligned representations (Schlechtweg et al., 2019). Alignment is needed as columns from different vector spaces may not correspond to the same coordinate axes, due to the stochastic nature of many low-dimensional word representations (Hamilton et al., 2016b). Following the above-described success, we use SGNS to create word representations in combination with Orthogonal Procrustes (OP) for vector space alignment and Cosine Distance (CD) (Salton and McGill, 1983) to measure differences between word vectors. From the resulting graded change predictions we infer binary change values by comparing the target word distribution to the full distribution of change predictions between the target corpora. For our experiments we use the code provided by Schlechtweg et al. (2019).¹

3.1 Semantic Representation

SGNS is a shallow neural network trained on pairs of word co-occurrences extracted from a corpus with a symmetric window. It represents each word w and each context c as a d -dimensional vector to solve

$$\arg \max_{\theta} \sum_{(w,c) \in D} \log \sigma(v_c \cdot v_w) + \sum_{(w,c) \in D'} \log \sigma(-v_c \cdot v_w),$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$, D is the set of all observed word-context pairs and D' is the set of randomly generated negative samples (Mikolov et al., 2013a; Mikolov et al., 2013b; Goldberg and Levy, 2014). The optimized parameters θ are v_{w_i} and v_{c_i} for $i \in 1, \dots, d$. D' is obtained by drawing k contexts from the empirical unigram distribution

$P(c) = \frac{\#(c)}{|D|}$ for each observation of (w, c) , cf. Levy et al. (2015). After training, each word w is represented by its word vector v_w .

Previous research on the influence of parameter settings on SGNS+OP+CD lays the foundation for our parameter choices (Schlechtweg et al., 2019; Kaiser et al., 2020). Although this subsystem combination is extremely stable regardless of parameter settings, subtle improvements can be achieved by modifying the window size and dimensionality. A common hurdle in LSC detection is the small corpus size, increasing the standard setting for window size from 5 to 10 leads to the creation of more word-context pairs used for training the model. In addition, we also experiment with dimensionalities of 300 and 500. Higher dimensionalities alleviate the introduction of noise during the alignment process (Kaiser et al., 2020). We keep the rest of the parameter settings at their default values (learning rate $\alpha=0.025$, #negative samples $k=5$ and sub-sampling $t=0.001$).

3.2 Alignment

SGNS is trained on each corpus separately, resulting in matrices A and B . To align them we follow Hamilton et al. (2016b) and calculate an orthogonally-constrained matrix W^* :

$$W^* = \arg \min_{W \in O(d)} \|BW - A\|_F$$

where the i -th row in matrices A and B correspond to the same word. Using W^* we get the aligned matrices $A^{OP} = A$ and $B^{OP} = BW^*$. Prior to this alignment step we length-normalize and mean-center both matrices (Artetxe et al., 2017; Schlechtweg et al., 2019).

3.3 Threshold

The DIACR-Ita shared task requires a binary label for each of the target words. However, CD produces graded values between 0.0 and 2.0 when measuring differences in word vectors between the two time periods. We tackle this problem by defining a threshold parameter, similar to many approaches applied in SemEval-2020 Task 1 (Schlechtweg et al., 2020). All words with a CD greater or equal than the threshold are labeled ‘1’, indicating change. Words with a CD less than the threshold are assigned ‘0’, indicating no change.

A simplified approach is to set the threshold such that the number of words is equal in both groups. This has many disadvantages: Mainly, it

¹<https://github.com/Garrafao/LSCDetection>

relies on the assumption that the two groups are of equal size. This is rarely given in real world applications, especially if the focus is in one word at a time. Thus a more sophisticated approach is needed. In SemEval-2020’s Subtask 1 many participants faced the same problem and developed various methods to solve it. Similar to the simplified approach, Zhou and Li (2020) only look at target words, and after fitting the histogram of CDs to a gamma distribution, set the threshold at the 75% density quantile. This approach resulted in good performance but is not always applicable due to its dependence on underlying properties of the test set. Amar and Liebeskind (2020) avoid the dependence on target words by randomly selecting 200 words and setting the threshold such that 90% of the 200 words have a lower distance than the threshold. A more careful selection of words is taken by Martinc et al. (2020), they look at the CD of semantically stable stop words, accumulate them in different bins and set the threshold to the upper limit of the bin containing fewer than $\frac{\#stopwords}{\#bins}$ words. Pražák et al. (2020) propose several methods. One of them is setting the threshold at the mean of the distances of all words in the corpus vocabulary. Our method for determining a threshold is very similar to Pražák et al. (2020), but instead of taking the mean, we use the mean + one standard deviation ($\mu + \sigma$) of all words in the corpus vocabulary.

4 Experimental setup

The DIACR-Ita task definition is taken from SemEval-2020 Task 1 Subtask 1 (binary change detection): Given a list of target words and a diachronic corpus pair C_1 and C_2 , the task is to identify the respective target words which have changed their meaning between the time periods t_1 and t_2 (Basile et al., 2020a; Schlechtweg et al., 2020).² C_1 and C_2 have been extracted from Italian newspapers and books. Target words which have changed their meaning are labeled with the value ‘1’, the remaining target words are labeled with ‘0’. Gold data for the 18 target words is semi-automatically generated from Italian online dictionaries. According to the gold data, 6 of the 18 target words are subject to semantic change between t_1 and t_2 . This gold data was only made public after the evaluation phase. During the evaluation

²The time periods t_1 and t_2 were not disclosed to participants.

entry	dim	threshold	ACC	AP	
#2	300	$(\mu+\sigma)$.76	.944	.915
#4	500	$(\mu+\sigma)$.78	.889	.915
#1	300	(50:50)	.57	.833	.915
#3	500	(50:50)	.64	.833	.915
major. baseline			-	.667	.333
freq. baseline		unk.		.611	.418
colloc. baseline		unk.		.500	unk.

Table 1: Accuracy (ACC) and Average Precision (AP) for various parameter settings and thresholds and baselines; *freq. baseline*: Absolute frequency difference between the words in C_1 and C_2 and an unknown threshold; *colloc. baseline*: Bag of Words + CD and an unknown threshold; *major. baseline*: Every word labeled with ‘0’.

phase each team was allowed to submit 4 predictions for the full list of target words, which were scored using classification accuracy between the predicted labels and the gold data. The final competition ranking compares only the highest of the 4 scores achieved by each team.

5 Results

We created target word rankings using SGNS+OP+CD with a dimensionality of 300 and 500 as described above. From these rankings our predictions are calculated using two different thresholding methods: (i) Splitting the targets into two equally-sized groups (50:50) and (ii) using the mean + one standard deviation ($\mu+\sigma$) as threshold, refer to Section 3.3. The accuracy scores achieved in this way are listed in Table 1, alongside the official baselines *freq.* and *colloc.* and an additional *major.* baseline. Submission #2 is our highest scoring submission and won the DIACR-Ita task together with one other undisclosed submission. For both of our rankings the 50:50 threshold yielded lower accuracy than the $\mu+\sigma$ threshold. This is due to the imbalance of changed to unchanged target words in the test set. Using $\mu+\sigma$ as threshold resulted in an optimal split for the ranking created with $d=300$. For $d=500$ this threshold was slightly too high with a value of 0.78. The target word *palmare* which, according to the gold data, has undergone semantic change (label ‘1’) has CD of 0.76 and was thus incorrectly labeled by our system. Figure 1 shows the histogram of CD values for all words of the corpus dictionary in gray. The green and

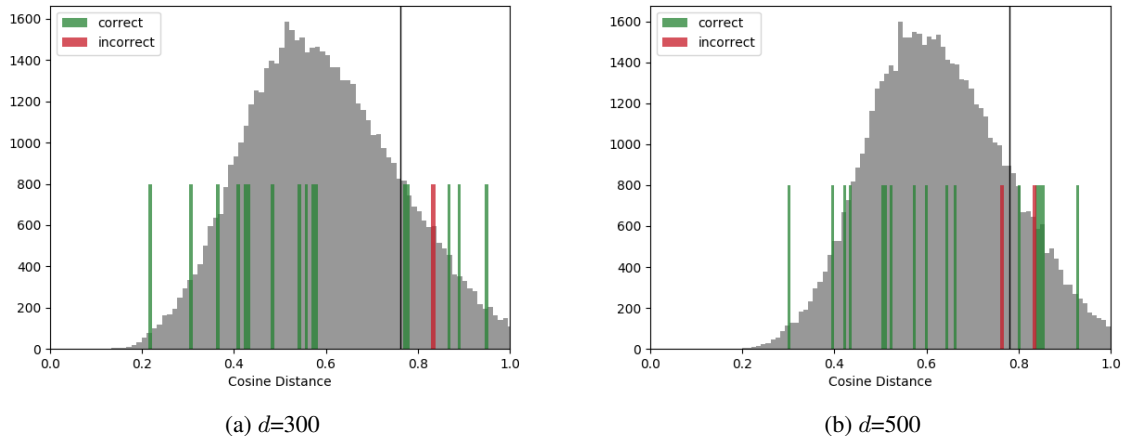


Figure 1: Background shows histogram (in gray) of CDs for all words in the corpus vocabulary. The colored bars show the CDs of target words, green indicates that the target word was correctly labeled, red indicates incorrect labeling. Vertical line marks threshold value (mean + standard deviation).

red colored bars correspond target words. If the target word was correctly labeled the bar is green, incorrect labeled target words have red bars. From this visualisation we can see that there is a pronounced gap between the CDs of target words which have changed and those which have not. Our proposed threshold method of $\mu + \sigma$ tends to slightly overshoot this gap. This has led to the lower accuracy of submission #4, despite the ranking allowing for a higher accuracy. In order to measure the quality of the rankings independent from the threshold we also report AP (Shwartz et al., 2017) in Table 1, confirming the potential equal performance.

The method of using the mean + one standard deviation of the CDs of all words in the corpus dictionary resulted in good accuracy, but leaves room for improvement. It tends to over-shoot the gap between unchanged and changed words slightly. Only using the mean shifts the tendency towards under-shooting the gap. The optimal threshold seems to lie somewhere in between. Though, this needs to be confirmed on other, larger, data sets. Furthermore, not all binary classification tasks are suitable for the approach of first creating a ranked list of graded change predictions and then choosing a threshold. The data set of SemEval-2020 Task 1 comprises two tasks, a binary and a ranked task for the same target words. It is not possible to achieve an accuracy of 1 on the binary task even if all the ranks are predicted correctly for the graded task, i.e., binary change is not just high graded change (Schlechtweg et al., 2020).

The one target word which our model labels incorrectly, across a variety of parameter settings, is *piovra*. According to the gold data this word has not undergone semantic change between t_1 and t_2 , while our system labels it as changed. A possible explanation for the error may be differences in frequency: In C_1 *piovra* appears 35 times and in C_2 it appears 643 times. SGNS often struggles to create reliable embeddings for low frequency words (Kaiser et al., 2020). Alternatively, the error could be caused by discrepancies between gold labels and corpora. Basile et al. (2020a) state that the gold data is initially based on Italian online dictionaries such as ‘Sabatini Coletti’. In a manual annotation process the gold data is further refined by providing human judges with up to 100 occurrences of each target word, for which they have to identify the used meaning according to the meanings listed in the dictionaries. A target word is labeled as changed if a meaning is observed in C_2 which has not been observed in C_1 . Although not very likely, it is possible that this annotation method fails to detect novel senses in C_2 . Sabatini Coletti reports that in addition to the sense “squid” *piovra* acquired a new sense “a secret criminal organisation deeply rooted in society” in 1983. This might explain why we detect *piovra* as a word which has undergone semantic change given that C_1 comprises texts from 1948 to 1970 and C_2 comprises texts from 1990 to 2014 (Basile et al., 2020a).

The DIACR-Ita task dataset is a very valuable contribution to the research field of LSC detec-

tion and extends the variety of available data sets to the Italian language. Nonetheless, two points are important when interpreting or results this data set: (i) it contains a small number of target words in combination with binary classification. This makes the data set vulnerable to randomness. (ii) The nature of the gold labels, in addition to possibly not being directly related to the corpus, it is unclear if they reflect semantic change as sense gain and sense loss as in SemEval’s Subtask 1. The online dictionaries which create the basis for the gold data only state sense gains. Thus, it might possible for a word to completely lose a sense but still be labeled as unchanged.

6 Conclusion

We participated in the DIACR-Ita shared task using well-established type-based methods for diachronic semantic representations in combination with a carefully calculated threshold. We were able to reach the first place with a nearly perfect accuracy of .94 confirming once more the reliability of the type-based embeddings created by SGNS, OP as an alignment method and CD to measure differences between word vectors. The presented approach is very suitable for similar tasks as no fine-tuning of parameters is needed. Yet, the system relies on the assumption that graded change is indicative of binary classes.

Acknowledgments

Dominik Schlechtweg was supported by the Konrad Adenauer Foundation and the CRETA center funded by the German Ministry for Education and Research (BMBF) during the conduct of this study. We thank the task organizers and reviewers for their efforts.

References

Efrat Amar and Chaya Liebeskind. 2020. JCT at SemEval-2020 Task 1: Combined Semantic Vector Spaces Models for Unsupervised Lexical Semantic Change Detection. In *Proceedings of the 14th International Workshop on Semantic Evaluation*, Barcelona, Spain. Association for Computational Linguistics.

Nikolay Arefyev and Vasily Zhikov. 2020. BOS at SemEval-2020 Task 1: Word Sense Induction via Lexical Substitution for Lexical Semantic Change Detection. In *Proceedings of the 14th International Workshop on Semantic Evaluation*,

Barcelona, Spain. Association for Computational Linguistics.

- Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2017. Learning bilingual word embeddings with (almost) no bilingual data. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pages 451–462. Association for Computational Linguistics.
- Ehsaneddin Asgari, Christoph Ringlstetter, and Hinrich Schütze. 2020. EmbLexChange at SemEval-2020 Task 1: Unsupervised Embedding-based Detection of Lexical Semantic Changes. In *Proceedings of the 14th International Workshop on Semantic Evaluation*, Barcelona, Spain. Association for Computational Linguistics.
- Pierpaolo Basile, Annalina Caputo, Tommaso Caselli, Pierluigi Cassotti, and Rossella Varvara. 2020a. DIACR-Ita @ EVALITA2020: Overview of the EVALITA2020 Diachronic Lexical Semantics (DIACR-Ita) Task. In Valerio Basile, Danilo Croce, Maria Di Maro, and Lucia C. Passaro, editors, *Proceedings of the 7th evaluation campaign of Natural Language Processing and Speech tools for Italian (EVALITA 2020)*, Online. CEUR.org.
- Valerio Basile, Danilo Croce, Maria Di Maro, and Lucia C. Passaro. 2020b. Evalita 2020: Overview of the 7th evaluation campaign of natural language processing and speech tools for italian. In Valerio Basile, Danilo Croce, Maria Di Maro, and Lucia C. Passaro, editors, *Proceedings of Seventh Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop (EVALITA 2020)*, Online. CEUR.org.
- Christin Beck. 2020. DiaSense at SemEval-2020 Task 1: Modeling sense change via pre-trained BERT embeddings. In *Proceedings of the 14th International Workshop on Semantic Evaluation*, Barcelona, Spain. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Yoav Goldberg and Omer Levy. 2014. Word2vec explained: Deriving Mikolov et al.’s negative-sampling word-embedding method. *arXiv:1402.3722*.
- William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016a. Cultural shift or linguistic drift? Comparing two computational measures of semantic change. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages

- 2116–2121, Austin, Texas. Association for Computational Linguistics.
- William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016b. Diachronic word embeddings reveal statistical laws of semantic change. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 1489–1501, Berlin, Germany. Association for Computational Linguistics.
- Renfen Hu, Shen Li, and Shichen Liang. 2019. Diachronic sense modeling with deep contextualized word embeddings: An ecological view. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3899–3908, Florence, Italy. Association for Computational Linguistics.
- Jens Kaiser, Dominik Schlechtweg, Sean Papay, and Sabine Schulte im Walde. 2020. IMS at SemEval-2020 Task 1: How low can you go? Dimensionality in Lexical Semantic Change Detection. In *Proceedings of the 14th International Workshop on Semantic Evaluation*, Barcelona, Spain. Association for Computational Linguistics.
- Vivek Kulkarni, Rami Al-Rfou, Bryan Perozzi, and Steven Skiena. 2015. Statistically significant detection of linguistic change. In *Proceedings of the 24th International Conference on World Wide Web*, WWW, pages 625–635, Florence, Italy.
- Andrey Kutuzov and Mario Giulianelli. 2020. UiO-UvA at SemEval-2020 Task 1: Contextualised Embeddings for Lexical Semantic Change Detection. In *Proceedings of the 14th International Workshop on Semantic Evaluation*, Barcelona, Spain. Association for Computational Linguistics.
- Andrey Kutuzov, Lilja Øvrelid, Terrence Szymanski, and Erik Velldal. 2018. Diachronic word embeddings and semantic shifts: a survey. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1384–1397, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Omer Levy, Yoav Goldberg, and Ido Dagan. 2015. Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics*, 3:211–225.
- Matej Martinc, Syrielle Montariol, Elaine Zosa, and Lidia Pivovarov. 2020. Discovery Team at SemEval-2020 Task 1: Context-sensitive Embeddings not Always Better Than Static for Semantic Change Detection. In *Proceedings of the 14th International Workshop on Semantic Evaluation*, Barcelona, Spain. Association for Computational Linguistics.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. In Yoshua Bengio and Yann LeCun, editors, *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013b. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems 26*, pages 3111–3119, Lake Tahoe, Nevada, USA.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, pages 1532–1543, Doha, Qatar.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2227–2237, New Orleans, LA, USA.
- Martin Pömsl and Roman Lyapin. 2020. CIRCE at SemEval-2020 Task 1: Ensembling Context-Free and Context-Dependent Word Representations. In *Proceedings of the 14th International Workshop on Semantic Evaluation*, Barcelona, Spain. Association for Computational Linguistics.
- Ondřej Pražák, Pavel Přibáček, Stephen Taylor, and Jakub Sido. 2020. UWB at SemEval-2020 Task 1: Lexical Semantic Change Detection. In *Proceedings of the 14th International Workshop on Semantic Evaluation*, Barcelona, Spain. Association for Computational Linguistics.
- Gerard Salton and Michael J McGill. 1983. *Introduction to Modern Information Retrieval*. McGraw-Hill Book Company, New York.
- Dominik Schlechtweg, Anna Häty, Marco del Tredici, and Sabine Schulte im Walde. 2019. A Wind of Change: Detecting and evaluating lexical semantic change across times and domains. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 732–746, Florence, Italy. Association for Computational Linguistics.
- Dominik Schlechtweg, Barbara McGillivray, Simon Hengchen, Haim Dubossarsky, and Nina Tahmasebi. 2020. SemEval-2020 Task 1: Unsupervised Lexical Semantic Change Detection. In *Proceedings of the 14th International Workshop on Semantic Evaluation*, Barcelona, Spain. Association for Computational Linguistics.
- Hinrich Schütze. 1998. Automatic word sense discrimination. *Computational Linguistics*, 24(1):97–123, March.

- Philippa Shoemark, Farhana Ferdousi Liza, Dong Nguyen, Scott Hale, and Barbara McGillivray. 2019. Room to Glo: A systematic comparison of semantic change detection approaches with word embeddings. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, pages 66–76, Hong Kong, China. Association for Computational Linguistics.
- Vered Shwartz, Enrico Santus, and Dominik Schlechtweg. 2017. Hypernyms under siege: Linguistically-motivated artillery for hypernymy detection. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, Valencia, Spain*, pages 65–75.
- Nina Tahmasebi, Lars Borin, and Adam Jatowt. 2018. Survey of computational approaches to diachronic conceptual change. *CoRR*, abs/1811.06278.
- Peter D. Turney and Patrick Pantel. 2010. From frequency to meaning: Vector space models of semantics. *J. Artif. Int. Res.*, 37(1):141–188, January.
- Jinan Zhou and Jiaxin Li. 2020. TemporalTeller at SemEval-2020 Task 1: Unsupervised Lexical Semantic Change Detection with Temporal Referencing. In *Proceedings of the 14th International Workshop on Semantic Evaluation, Barcelona, Spain*. Association for Computational Linguistics.