
AUGMENTY: A PYTHON LIBRARY FOR STRUCTURED TEXT AUGMENTATION

A PREPRINT

Kenneth Enevoldsen 

December 9, 2023

ABSTRACT

Augmnety is a Python library for structured text augmentation. It is built on top of spaCy and allows for augmentation of both the text and its annotations. Augmenty provides a wide range of augmenters which can be combined in a flexible manner to create complex augmentation pipelines. It also includes a set of primitives that can be used to create custom augmenters such as word replacement augmenters. This functionality allows for augmentations within a range of applications such as named entity recognition (NER), part-of-speech tagging, and dependency parsing.

Keywords Python • natural language processing • spacy • augmentation

1 Summary

Text augmentation is useful for tool for training (Wei and Zou 2019) and evaluating (Ribeiro et al. 2020) natural language processing models and systems. Despite its utility existing libraries for text augmentation often exhibit limitations in terms of functionality and flexibility, being confined to basic tasks such as text-classification or cater to specific downstream use-cases such as estimating robustness (Goel et al. 2021). Recognizing these constraints, Augmenty is a tool for structured text augmentation of the text along with its annotations. Augmenty integrates seamlessly with the popular NLP library spaCy (Honnibal et al. 2020) and seeks to be compatible with all models and tasks supported by spaCy. Augmenty provides a wide range of augmenters which can be combined in a flexible manner to create complex augmentation pipelines. It also includes a set of primitives that can be used to create custom augmenters such as word replacement augmenters. This functionality allows for augmentations within a range of applications such as named entity recognition (NER), part-of-speech tagging, and dependency parsing.

2 Statement of need

Augmentation is a powerful tool within disciplines such as computer vision (Wang, Perez, et al. 2017) and speech recognition (Park et al. 2019) and it used for both training more robust models and evaluating models ability to handle perturbations. Within natural language processing (NLP) augmentation has seen some uses as a tool for generating additional training data (Wei and Zou 2019), but have really shined as a tool for model evaluation, such as estimating robustness (Goel et al. 2021) and bias (Lassen et al. 2023), or for creating novel datasets (Nielsen 2023).

Despite its utility, existing libraries for text augmentation often exhibit limitations in terms of functionality and flexibility. Commonly they only provide pure string augmentation which typically leads to the annotations becoming misaligned with the text. This has limited the use of augmentation to tasks such as text classification while neglected structured prediction tasks such as named entity recognition (NER) or coreference resolution. This has limited the use of augmentation to a wide range of tasks both for training and evaluation.

Existing tools such as `textgenie` (Pandya 2023), and `textaugment` (Marivate and Sefara 2020) implements powerful techniques such as backtranslation and paraphrasing, which are useful for augmentation for text-classification tasks. However, these tools neglect a category of tasks which require that the annotations are aligned with the augmentation of the text. For instance even simple augmentation such as replacing the named entity “Jane Doe” with “John” will lead to a misalignment of the NER annotation, part-of-speech tags, etc., which if not properly handled will lead to a misinterpretation of the model performance or generation of incorrect training samples.

Augmenty seeks to remedy this by providing a flexible and easy-to-use interface for structured text augmentation. Augmenty is built to integrate well with the spaCy (Honnibal et al. 2020) and seeks to be compatible with the broad set of tasks supported by spaCy. Augmenty provides augmenters which takes in a spaCy Doc-object (but works just as well with string-objects) and returns a new Doc-object with the augmentations applied. This allows for augmentations of both the text and the annotations present in the Doc-object.

Other tools for data augmentation focus on specific downstream application such `textattack` (Morris et al. 2020) which is useful for adversarial attacks of classification systems or `robustnessgym` (Goel et al. 2021) which is useful for evaluating robustness of classification systems. Augmenty does not seek to replace any of these tools but seeks to provide a general purpose tool for augmentation of both the text and its annotations. This allows for augmentations within a range of applications such as named entity recognition, part-of-speech tagging, and dependency parsing.

3 Features & Functionality

Augmenty is a Python library that implements augmentation based on spaCy’s Doc object. spaCy’s Doc object is a container for a text and its annotations. This makes it easy to augment text and annotations simultaneously. The Doc object can easily be extended to include custom augmentation not available in spaCy by adding custom attributes to the Doc object. While Augmenty is built to augment Docs the object is easily converted into strings, lists or other formats. The annotations within a Doc can be provided either by existing annotations or by annotations provided by an existing model.

Augmenty implements a series of augmenters for token-, span- and sentence-level augmentation. These augmenters range from primitive augmentations such as word replacement which can be used to quickly construct new augmenters to language specific augmenters such as keystroke error augmentations based on a French keyboard layout. Augmenty also integrates with other libraries such as MLTK [bird2009natural] to allow for augmentations based on WordNet (Miller 1994) and allows for specification of static word vectors [pennington-etal-2014-glove] to allow for augmentations based on word similarity. Lastly, augmenty provides a set of utility functions for repeating augmentations, combining augmenters or adjust the percentage of documents that should be augmented. This allow for the flexible construction of augmentation pipelines specific to the task at hand.

Augmenty is furthermore designed to be compatible with spaCy and thus its augmenters can easily be utilized during the training of spaCy models.

4 Example Use Cases

Augmenty have already seen used in a number of projects. The code base was initially developed for evaluating the robustness and bias of DaCy (Enevoldsen, Hansen, and Nielbo 2021), a state-of-the-art Danish NLP pipeline. It is also continually used to evaluate Danish NER systems for biases and robustness on the DaCy website. Augmenty has also been used to detect intersectional biases (Lassen et al. 2023) and used within benchmark of Danish language models (Sloth and Rybner 2023).

Besides its existing use-cases Augmenty could for example also be used to a) upsample minority classes without duplicating samples, b) train less biased models by e.g. replacing names with names of minority groups c) train more robust models e.g. by augmenting with typos or d) generate pseudo historical data by augmenting with known spelling variations of words.

5 Target Audience

The package is mainly targeted at NLP researchers and practitioners who wish to augment their data for training or evaluation. The package is also targeted at researchers who wish to evaluate their models either augmentations or generating new datasets.

6 Acknowledgements

The authors thank the [contributors](#) of the package notably Lasse Hansen which provided meaningful feedback on the design of the package at early stages of development.

Enevoldsen, Kenneth, Lasse Hansen, and Kristoffer L. Nielbo. 2021. “DaCy: A Unified Framework for Danish NLP.” https://ceur-ws.org/Vol-2989/short_paper24.pdf.

Goel, Karan, Nazneen Fatema Rajani, Jesse Vig, Zachary Taschdjian, Mohit Bansal, and Christopher Ré. 2021. “Robustness Gym: Unifying the NLP Evaluation Landscape.” In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Demonstrations*, edited by Avi Sil and Xi Victoria Lin, 42–55. Online: Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.naacl-demos.6>.

- Honnibal, Matthew, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. “spaCy: Industrial-Strength Natural Language Processing in Python.” <https://doi.org/10.5281/zenodo.1212303>.
- Lassen, Ida Marie S., Mina Almasi, Kenneth Enevoldsen, and Ross Deans Kristensen-McLachlan. 2023. “Detecting Intersectionality in NER Models: A Data-Driven Approach.” In *Proceedings of the 7th Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*, edited by Stefania Degaetano-Ortlieb, Anna Kazantseva, Nils Reiter, and Stan Szpakowicz, 116–27. Dubrovnik, Croatia: Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.latechclfl-1.13>.
- Marivate, Vukosi, and Tshephisho Sefara. 2020. “Improving Short Text Classification Through Global Augmentation Methods.” In *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*, 385–99. Springer.
- Miller, George A. 1994. “WordNet: A Lexical Database for English.” In *Human Language Technology: Proceedings of a Workshop Held at Plainsboro, New Jersey, March 8-11, 1994*. <https://aclanthology.org/H94-1111>.
- Morris, John, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020. “TextAttack: A Framework for Adversarial Attacks, Data Augmentation, and Adversarial Training in NLP.” In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 119–26.
- Nielsen, Dan. 2023. “ScandEval: A Benchmark for Scandinavian Natural Language Processing.” In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, edited by Tanel Alumäe and Mark Fishel, 185–201. Tórshavn, Faroe Islands: University of Tartu Library. <https://aclanthology.org/2023.nodalida-1.20>.
- Pandya, Het. 2023. “Hetpandya/Textgenie.” <https://github.com/hetpandya/textgenie>.
- Park, Daniel S., William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin Dogus Cubuk, and Quoc V. Le. 2019. “SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition.” In *Interspeech*. <https://api.semanticscholar.org/CorpusID:121321299>.
- Ribeiro, Marco Tulio, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. “Beyond Accuracy: Behavioral Testing of NLP Models with CheckList.” In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, edited by Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, 4902–12. Online: Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.442>.
- Sloth, Thea Rolskov, and Astrid Sletten Rybner. 2023. “DaDebias/Genda-Lens.” DaDebias. <https://github.com/DaDebias/genda-lens>.
- Wang, Jason, Luis Perez, et al. 2017. “The Effectiveness of Data Augmentation in Image Classification Using Deep Learning.” *Convolutional Neural Networks Vis. Recognit* 11 (2017): 1–8.
- Wei, Jason, and Kai Zou. 2019. “EDA: Easy Data Augmentation Techniques for Boosting Performance on Text Classification Tasks.” In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, edited by Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, 6382–88. Hong Kong, China: Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-1670>.