

Small-Text: Active Learning for Text Classification in Python

Christopher Schröder[†]

CHRISTOPHER.SCHROEDER@UNI-LEIPZIG.DE

Lydia Müller^{†§}

LYDIA.MUELLER@UNI-LEIPZIG.DE

Andreas Niekler[†]

ANDREAS.NIEKLER@UNI-LEIPZIG.DE

Martin Potthast[†]

MARTIN.POTTHAST@UNI-LEIPZIG.DE

[†]*Leipzig University, Germany*

[§]*Institute for Applied Informatics (InfAI), Leipzig, Germany*

Abstract

We present `small-text`, an easy-to-use active learning library, which offers pool-based active learning for single- and multi-label text classification in Python. It features many pre-implemented state-of-the-art query strategies, including some that leverage the GPU. Standardized interfaces allow the combination of a variety of classifiers, query strategies, and stopping criteria, facilitating a quick mix and match, and enabling a rapid development of both active learning experiments and applications. To make various classifiers and query strategies accessible in a unified way, `small-text` integrates the well-known machine learning libraries `scikit-learn`, PyTorch, and `huggingface transformers`. The latter integrations are available as optionally installable extensions, making the availability of a GPU completely optional. The library is publicly available under the MIT License at <https://github.com/webis-de/small-text>, version 1.0.0b4 at the time of writing.

Keywords: active learning, text classification, query strategies, transformers

1. Introduction

Text classification—in the same way as most contemporary machine learning applications—requires large amounts of training data to achieve peak performance. However, in many real-world use cases, labeled data does not exist and is expensive to obtain—especially if domain expertise is required. *Active learning* (Lewis and Gale, 1994) solves this problem by repeatedly selecting unlabeled data deemed to be informative according to a so-called *query strategy*, which are then labeled by a human annotator. Subsequently, a new model is trained on all data labeled so far, and then this process is repeated until a *stopping criterion* has been met. Active learning aims at minimizing the amount of labeled data required, while maximizing the resulting model’s performance, e.g., in terms of classification accuracy.

An active learning solution therefore consists of corresponding components, namely a classifier, a query strategy, and an optional stopping criterion. Meanwhile, many strategies have been devised and studied for each of them. Seeing as choosing a worthwhile combination among them can only be done through extensive experimentation, and as implementing the different strategies is non-trivial while implementation details may differ dependent on the data at hand, this induces a considerable overhead. An obvious solution to this problem

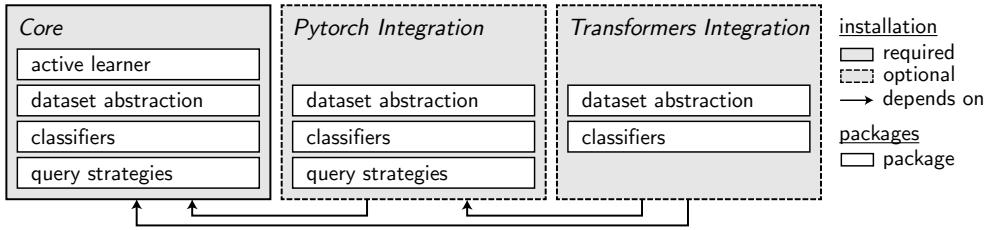


Figure 1: Architecture of `small-text`. The core installation can be optionally extended with a) the `Pytorch` integration, which allows to leverage GPU-based models, and b) the `transformers` integration, which allows to use state-of-the-art transformer-based text classification. Dependencies between packages are omitted.

is to employ open source libraries, which, among other benefits, accelerate research and ease the adoption of methods by both researchers and practitioners (Sonnenburg et al., 2007). Although there are existing solutions for general active learning, only few consider active learning for text classification, which requires functionality specific to the domain of natural language processing, such as word embeddings (Mikolov et al., 2013) or language models (Devlin et al., 2019). To fill this gap, we present a library, which provides tried-and-tested components for active learning experiments with application to text classification.

2. Overview of Small-Text

The main goal of `small-text` is to offer state-of-the-art active learning for text classification in a convenient and robust way for both researchers and practitioners. For this purpose, we implement a modular pool-based active learning mechanism, illustrated in Figure 1, which exposes interfaces for classifiers, query strategies, and stopping criteria. The core of `small-text` integrates `scikit-learn` (Pedregosa et al., 2011), enabling direct use of its classifiers. Overall, the library provides 13 query strategies, including some that are only usable on text data, and two integrations of well-known machine learning libraries, namely `PyTorch` (Paszke et al., 2019) and `transformers` (Wolf et al., 2020). The integrations ease the use of CUDA-based GPU computing and transformer models, respectively. The modular architecture renders both integrations completely optional, resulting in a slim core, which can be used without unnecessary dependencies in a CPU-only scenario.

As the query strategy, which selects the instances to be labeled, is the most salient component of an active learning setup, the range of alternative query strategies provided covers four paradigms at the time of writing: (i) confidence-based strategies: least confidence (Lewis and Gale, 1994; Culotta and McCallum, 2005), prediction entropy (Roy and McCallum, 2001), breaking ties (Luo et al., 2005), and contrastive active learning (Margatina et al., 2021); (ii) embedding-based strategies: BADGE (Ash et al., 2020), BERT k-means (Yuan et al., 2020), discriminative active learning (Gissin and Shalev-Shwartz, 2019), and SEALS (Coleman et al., 2020); (iii) gradient-based strategies: expected gradient length (EGL; Settles et al., 2007), EGL-word (Zhang et al., 2017), and EGL-sm (Zhang et al., 2017); and (iv) coreset strategies: greedy coresset (Sener and Savarese, 2018) and lightweight coresset (Bachem et al., 2018). Since there is an abundance of query strategies,

Name	Active Learning			Code					
	QS	SC	Text Focus	GPU support	Unit Tests	Language	License	Last Update	Repository
JCLAL ¹	18	2	✗	✗	✗	Java	GPL	2017	○
libact ²	19	-	✗	✗	✓	Python	BSD-2-Clause	2021	○
modAL ³	12	-	✗	✓	✓	Python	MIT	2022	○
ALiPy ⁴	22	4	✗	✗	✓	Python	BSD-3-Clause	2021	○
lrtc ⁵	7	-	✓	✓	✗	Python	Apache 2.0	2021	○
small-text	13	2	✓	✓	✓	Python	MIT	2022	○

Table 1: Comparison between `small-text` and relevant previous active learning libraries. We abbreviated the number of query strategies by “QS”, the number of stopping criteria by “SC”, and the low-resource-text-classification framework by `lrtc`. All information except “Publication Year” and “Code Repository” has been extracted from the linked Github repository of the respective library on April 22nd, 2022. Random baselines were not counted towards the total number of query strategies. Publications: ¹Reyes et al. (2016), ²(Yang et al., 2017), ³(Danka and Horvath, 2018), ⁴(Tang et al., 2019), ⁵(Ein-Dor et al., 2020).

this list will likely never be exhaustive—also because strategies from other domains, such as computer vision, are not always applicable to the text domain, e.g., when relying on the geometry of images (Konyushkova et al., 2015) and thus will be disregarded here.

The library is available via the python packaging index and can be installed with just a single command: `pip install small-text`. Similarly, the integrations can be enabled using the extra requirements argument of Python’s `setuptools`, e.g., the `transformers` integration is installed using `pip install small-text[transformers]`. The robustness of the implementation rests on extensive unit and integration tests. Detailed examples, an API documentation, and common usage patterns are available in the online documentation.¹

3. Comparison to Previous Software

Unsurprisingly, after decades of research and development on active learning, numerous other libraries are available that focus on active learning as well. In the following we present a selection of the most relevant open-source projects for which either a related publication is available or a larger user base exists: JCLAL (Reyes et al., 2016) is a generic framework for active learning which is implemented in Java and can be used either through XML configurations or directly from the code. It offers an experimental setting which includes 18 query strategies. The aim of libact (Yang et al., 2017) is to provide active learning for real-world applications. Among 19 other strategies, it includes a well-known meta-learning strategy (Hsu and Lin, 2015). The modAL library (Danka and Horvath, 2018) offers active learning including regression, multi-label classification and stream-based active learning. It offers 12 query strategies, also builds on `scikit-learn` by default, and provides instructions how to include GPU-based models using Keras and PyTorch. ALiPy (Tang et al., 2019) provides

1. <https://small-text.readthedocs.io>

an active learning framework targeted at the experimental active learning setting. Apart from providing 22 query strategies, it supports alternative active learning settings, e.g., active learning with noisy annotators. The `low-resource-text-classification-framework` (`lrtc`; Ein-Dor et al. (2020)) also focuses on text classification and has a number of built-in models, datasets, and query strategies to perform active learning experiments.

In Table 1, we compare each of those projects and `small-text` by multiple criteria related to active learning or related to the specific code: While all libraries provide a selection of query strategies, not all libraries offer stopping criteria, which are crucial to reducing the total annotation effort and thus directly influence the efficiency of the active learning process (Vlachos, 2008; Laws and Schütze, 2008; Olsson and Tomanek, 2009). We can also see a difference in the number of provided query strategies, which is however not conclusive on its own: although a higher number of query strategies is certainly not a disadvantage, it is more important to provide the most relevant strategies (either due to recency, domain-specificity, or strong performance) since active learning experiments are computationally expensive (Margatina et al., 2021; Schröder et al., 2022) and therefore not every strategy can be tested during an experiment. Likewise, for application scenarios, only the one presumed best strategy will be employed. Only `small-text` and `lrtc` focus specifically on text classification, and solely `modAL`, `lrtc` and `small-text` offer access to GPU-based deep learning frameworks, which has become indispensable for competitive text classification due to the recent success and ubiquity of transformer-based models (Vaswani et al., 2017; Devlin et al., 2019). Finally, only `small-text` provides recent strategies such as BADGE (Ash et al., 2020), BERT K-Means (Yuan et al., 2020), and contrastive active learning (Margatina et al., 2021), as well as the gradient-based strategies by Zhang et al. (2017), where the latter are unique to active learning for text classification.

The distinguishing characteristic of `small-text` is the integration of `scikit-learn`, PyTorch, and `transformers`, which makes it possible to easily combine a wide range of classifiers, query strategies and stopping criteria. It provides a broad set of features including GPU support, stopping criteria, robustness through unit tests, and most importantly, it covers concepts that are specific to text classification such as embeddings, language models, and the text-based KimCNN (Kim, 2014). In summary, `small-text` offers a wide range of components, which are specifically targeted at text classification, thereby enabling state-of-the-art active learning for natural language processing using only a few lines of code.

4. Conclusion

We introduced `small-text`, a modular Python library, which offers active learning for text classification. It integrates `scikit-learn`, PyTorch, and `transformers`, and provides robust components to quickly apply active learning in both experiments and applications, thereby making state-of-the-art active learning easily accessible to the Python ecosystem.

Acknowledgments

This research was partially funded by the Development Bank of Saxony (SAB) under project numbers 100335729 and 100400221.

References

- Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. In *Proceedings of the 8th International Conference on Learning Representations (ICLR)*, 2020.
- Olivier Bachem, Mario Lucic, and Andreas Krause. Scalable k-Means Clustering via Lightweight Coresets. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, (KDD)*, pages 1119–1127, 2018.
- Cody Coleman, Edward Chou, Sean Culatana, Peter Bailis, Alexander C. Berg, Roshan Sumbaly, Matei Zaharia, and I. Zeki Yalniz. Similarity search for efficient active learning and search of rare concepts. *arXiv preprint arXiv:2007.00077*, 2020.
- Aron Culotta and Andrew McCallum. Reducing labeling effort for structured prediction tasks. In *Proceedings of the 20th National Conference on Artificial Intelligence (AAAI)*, volume 2, pages 746–751, 2005.
- Tivadar Danka and Peter Horvath. modAL: A modular active learning framework for Python. *arXiv preprint arXiv:1805.00979*, 2018.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 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 (NAACL)*, pages 4171–4186, 2019.
- Liat Ein-Dor, Alon Halfon, Ariel Gera, Eyal Shnarch, Lena Dankin, Leshem Choshen, Marina Danilevsky, Ranit Aharonov, Yoav Katz, and Noam Slonim. Active Learning for BERT: An Empirical Study. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7949–7962, 2020.
- Daniel Gissin and Shai Shalev-Shwartz. Discriminative active learning. *arXiv preprint arXiv:1907.06347*, 2019.
- Wei-Ning Hsu and Hsuan-Tien Lin. Active learning by learning. *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI)*, 29(1), 2015.
- Yoon Kim. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1746–1751, 2014.
- Ksenia Konyushkova, Raphael Sznitman, and Pascal Fua. Introducing geometry in active learning for image segmentation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2974–2982, 2015.
- Florian Laws and Hinrich Schütze. Stopping criteria for active learning of named entity recognition. In *Proceedings of the 22nd International Conference on Computational Linguistics (COLING)*, pages 465–472, Manchester, UK, 2008.

David D. Lewis and William A. Gale. A sequential algorithm for training text classifiers. In *Proceedings of the 17th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, pages 3–12, 1994.

Tong Luo, Kurt Kramer, Dmitry B. Goldgof, Lawrence O. Hall, Scott Samson, Andrew Remsen, and Thomas Hopkins. Active Learning to Recognize Multiple Types of Plankton. *Journal of Machine Learning Research (JMLR)*, 6:589–613, 2005.

Katerina Margatina, Giorgos Vernikos, Loïc Barrault, and Nikolaos Aletras. Active learning by acquiring contrastive examples. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 650–663, 2021.

Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *Proceedings of the 1st International Conference on Learning Representations (ICLR)*, 2013.

Fredrik Olsson and Katrin Tomanek. An intrinsic stopping criterion for committee-based active learning. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning (CoNLL)*, pages 138–146, 2009.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035, 2019.

Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research (JMLR)*, 12(85):2825–2830, 2011.

Oscar Reyes, Eduardo Pérez, María del Carmen Rodríguez-Hernández, Habib M. Fardoun, and Sebastián Ventura. JCLAL: A Java Framework for Active Learning. *Journal of Machine Learning Research (JMLR)*, 17(95):1–5, 2016.

Nicholas Roy and Andrew McCallum. Toward optimal active learning through sampling estimation of error reduction. In *Proceedings of the Eighteenth International Conference on Machine Learning (ICML)*, pages 441–448, 2001.

Christopher Schröder, Andreas Niekler, and Martin Potthast. Revisiting uncertainty-based query strategies for active learning with transformers. *arXiv preprint arXiv:2107.05687v2*, 2022.

Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *Proceedings of the 6th International Conference on Learning Representations (ICLR)*, 2018.

Burr Settles, Mark Craven, and Soumya Ray. Multiple-instance active learning. In *Proceedings of the 20th International Conference on Neural Information Processing Systems (NIPS)*, pages 1289–1296, 2007.

Sören Sonnenburg, Mikio L. Braun, Cheng Soon Ong, Samy Bengio, Leon Bottou, Geoffrey Holmes, Yann LeCun, Klaus-Robert Müller, Fernando Pereira, Carl Edward Rasmussen, Gunnar Rätsch, Bernhard Schölkopf, Alexander Smola, Pascal Vincent, Jason Weston, and Robert Williamson. The Need for Open Source Software in Machine Learning. *Journal of Machine Learning Research (JMLR)*, 8(81):2443–2466, 2007.

Ying-Peng Tang, Guo-Xiang Li, and Sheng-Jun Huang. ALiPy: Active learning in python. *arXiv preprint arXiv:1901.03802*, 2019.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the Advances in Neural Information Processing Systems 30 (NeurIPS)*, pages 5998–6008, 2017.

Andreas Vlachos. A stopping criterion for active learning. *Computer Speech & Language*, 22(3):295–312, 2008.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierrick Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations (EMNLP)*, pages 38–45, 2020.

Yao-Yuan Yang, Shao-Chuan Lee, Yu-An Chung, Tung-En Wu, Si-An Chen, and Hsuan-Tien Lin. libact: Pool-based active learning in python. *arXiv preprint arXiv:1710.00379*, 2017.

Michelle Yuan, Hsuan-Tien Lin, and Jordan Boyd-Graber. Cold-start active learning through self-supervised language modeling. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7935–7948. Association for Computational Linguistics, 2020.

Ye Zhang, Matthew Lease, and Byron C. Wallace. Active discriminative text representation learning. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence (AAAI)*, pages 3386–3392, 2017.