Responsible NLP Checklist

Paper title: NLP-ADBench: NLP Anomaly Detection Benchmark
Authors: Yuangang Li, Jiaqi li, Zhuo Xiao, Tiankai Yang, Yi Nian, Xiyang Hu, Yue Zhao

How to read the checklist symbols:
the authors responded 'yes'
the authors responded 'no'
the authors indicated that the question does not apply to their work
the authors did not respond to the checkbox question
For background on the checklist and guidance provided to the authors, see the Responsible NLP Checklist page at ACL Rolling Review.

✓ A. Questions mandatory for all submissions.

- A1. Did you describe the limitations of your work? *This paper has a Limitations section.*
- A2. Did you discuss any potential risks of your work?

 This work is solely for the public release of benchmark data and algorithms, and does not involve direct risks or potential hazards. Therefore, risks are not specifically discussed. The paper focuses on promoting academic reproducibility and technological progress, without any practical application implementation or production deployment content, nor has it encountered sensitive or risky scenarios.

B. Did you use or create scientific artifacts? (e.g. code, datasets, models)

- ☑ B1. Did you cite the creators of artifacts you used?

 See References and Appendices (the dataset is detailed in Appendix A.1, and the algorithms are detailed in Appendix A.2).
- ☑ B2. Did you discuss the license or terms for use and/or distribution of any artifacts? *See Ethics Statement section.*
- ☑ B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?

 See the Ethics Statement and Appendix A.1 and A.2.
- B4. Did you discuss the steps taken to check whether the data that was collected/used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect/anonymize it?

 See Ethics Statement and A.1.
- ☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

 See 2.2, A.1 and A.2.
- ☑ B6. Did you report relevant statistics like the number of examples, details of train/test/dev splits, etc. for the data that you used/created? *See Table 1*, 2.2, 3.1, and A.1.

☑ C. Did you run computational experiments?

∠C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?

We did not report the number of parameters and specific computational budget for each model in detail, but only described the model categories and training process. Since most algorithms are standard open-source implementations, the paper focuses on algorithm evaluation and comparison. Relevant content can be obtained from the original model literature.

- ✓ C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

 See 3.1
- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

 See 3.1, 3.2, Table 2, A.3.
- ∠ C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation, such as NLTK, SpaCy, ROUGE, etc.), did you report the implementation, model, and parameter settings used?

We provide a complete Conda environment file (environment.yml) in our public GitHub repository. This file specifies all dependencies and version numbers required to reproduce our experiments. Users can create an identical environment by running 'conda env create -f environment.yml', ensuring full reproducibility.

\(\begin{aligned} \D.\) Did you use human annotators (e.g., crowdworkers) or research with human subjects?

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.? (*left blank*)
- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)? (*left blank*)
- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating (e.g., did your instructions explain how the data would be used)? (*left blank*)
- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board? (*left blank*)
- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data? (*left blank*)

E. Did you use AI assistants (e.g., ChatGPT, Copilot) in your research, coding, or writing?

☑ E1. If you used AI assistants, did you include information about their use?

See Ethics Statement. Specifically states "We used ChatGPT exclusively to improve minor grammar in the final manuscript text."