Responsible NLP Checklist

Paper title: Optimising Factual Consistency in Summarisation via Preference Learning from Multiple Imperfect Metrics

Authors: Yuxuan Ye, Raul Santos-Rodriguez, Edwin Simpson

How to read the checklist symbols:	
the authors responded 'yes'	
X 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?

 We believe our work does not pose any significant ethical or societal risks. It does not involve sensitive data, high-stakes applications, or methods that could be misused in harmful ways.
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
 - ☑ B1. Did you cite the creators of artifacts you used?

 We cited every creators of artifacts we used in the experiments in section 4.
 - B2. Did you discuss the license or terms for use and/or distribution of any artifacts? *The datasets and models used in the experiments are publicly available.*
 - 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)?

 Our use of existing artifacts is compatible with the original access conditions.
 - 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?

The datasets used in the experiments are anonymous and do not contain personal information.

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

 The documentation of artifacts used in the experiments are publicly available.
- 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?

The statistics of datasets are reported in Section 4.1.1.

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 reported the size of selected models and experimental setup in Section 4.1.

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

We described the experimental setup in Section 4.

☑ 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?

We described the experimental setup in Section 4.1, all the tables are explained in details in the caption.

✓ 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 included links for the improtant libraries used in the experiments.

\(\begin{aligned} \Bigsilon \) 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?

We use ChatGPT to evaluate our results. The details are explained in Section 4 and 5 and Appendix