

TextAttack: A Framework for Adversarial Attacks, Data Augmentation, and Adversarial Training in NLP

1. What is the primary lesson(s) you took away from this paper (avoid abstract level summary)?

The authors in this paper introduced a framework, TextAttack, a Python package for adversarial attacks, data augmentation, and adversarial training in NLP. This framework (TextAttack) provides implementations of 16 adversarial attacks from the literature and supports a variety of models and datasets, including BERT and other transformers, and all GLUE tasks.

2. Could I have done this work if I had the idea why or why not?

I think if I had the idea first, I could have done this work. The high-level idea of the framework (TextAttack) is to re-implement and organize already published text adversarial attack algorithms that have been introduced to the research community.

3. Is there anything I could do to repeat or validate? (Make notes of materials available from the paper- data sets, source code).

Yes, the authors open-sourced their framework (TextAttack) with detailed documentation here: <https://github.com/QData/TextAttack>. Docs: <https://textattack.readthedocs.io/en/master/index.html>.

4. What is my best idea for follow on work that I could personally do?

My best idea for the follow-on work that I could personally do is to experiment with this framework (TextAttack) using their English datasets to understand better how it works. Secondly, I could experiment with this framework using some Arabic datasets, examine the generated results, and test if all supported adversarial text attacks work in different languages like Arabic.

5. What is my best idea for follow on work that I'd like to see the authors or others do?

I do *not* have an idea for the follow-on work I would like to see the authors or others do.

6. Any logistical experimental lessons I learned (this is a little different than #1 here you are looking for techniques)?

I have *not* learned any logistical experimental from this paper.

7. How does this compare to the other papers we read? Most similar? How different? Other comparisons?

This paper by Morris et al. (2020) is the first paper that proposed a text adversarial attacks framework (TextAttack) that re-implements all the attack algorithms for most of the papers I have read so far in this course. Yet, this paper does *not* propose any adversarial attack algorithm.

8. What is your biggest criticism of the paper?

There is no big criticism of the paper.

9. List 3 cited references or terms/concepts that you would be most interested in reading/learning more about.

I am interested in reading this reference (Wei and Zou, 2019) about Text Data Augmentation techniques in NLP. I must read this reference to use it in one of my future projects.

Jason W. Wei and Kai Zou. 2019. EDA: easy data augmentation techniques for boosting performance on text classification tasks. CoRR, abs/1901.11196.