

1 We thank the reviewers for the insightful and helpful comments. Minor remarks will be reflected in the text. Individual
2 responses to the questions are included below.

3 **Reviewer 1**

4 The reviewer raises the very relevant issue of applicability, both theoretical and practical. The paper’s framework is
5 general and can be applied to a wide class of problems and ML models. Hence, we agree with the reviewer in that
6 the theoretical applicability is significant. Moreover, we will improve the presentation of experiments, to highlight the
7 practical applicability of the work.

8 Also, the practical applicability of our work depends on a reasoner for some relevant fragment of first order logic. The
9 same applies to a number of recent works that exploit formal methods. The improvements made to reasoners, e.g. SMT
10 solvers, in recent years offer guarantees that practical applicability will continue to improve.

11 Furthermore, there is already practical evidence to the relevance and the insights provided by the use of formal methods
12 in ML; some are referenced in our paper.

13 The duality relationship that our paper reveals will enable researchers to look at adversarial examples and explanations
14 in a new light. This will open new avenues of research, that will foster more efficient exact methods but also better
15 heuristics.

16 The reviewer is quite right that explanations find important uses in many other settings. We will expand on this.

17 **Reviewer 2**

18 We thank the reviewer for the suggestions regarding presentation.

19 We will improve readability of the paper, e.g. by clarifying the concepts that are not widely known in the community,
20 including some of the logic-related notation used.

21 The proof of Theorem 1 will be cleaned up in the revised paper, as suggested by the reviewer. All the minor issues
22 spotted will be addressed. Thank you!

23 **Reviewer 3**

24 The reviewer raises a ‘*so what*’ question.

25 The importance of having a deep understanding of the connection between adversarial examples and explanations is
26 epitomized by a number of recent works, some of which are cited in our paper.

27 Also significant is a recent keynote talk by I. Goodfellow at the AAAI 2019 conference on “*Adversarial Machine*
28 *Learning*” (please see the official video recording starting from 53m55s), which makes a very strong case for relating
29 adversarial examples and explanations. Namely, Goodfellow conjectures that there should be a connection between
30 adversarial attacks and explanations and points out that this connection has not been revealed so far. We believe that our
31 work is a significant step forward in this direction.

32 We chose an image dataset as we can visually demonstrate the discovered relationship between adversarial examples
33 and explanations to the reader. We want to convey that each explanation hits all adversarial examples and vice versa.
34 This “hitting” process can be highlighted using coloring in the image (please see Figures 1 and 2 in the paper). While
35 our method can be applied to categorical datasets, we feel that such visualization might be less visually appealing to a
36 reader.