

1 Thank you for your thoughtful feedback. We will first discuss common themes and then specific reviewer comments.

2 **Significance:** Even though ExpO is “simple” (in that it connects existing concepts, albeit in a novel way), we believe  
3 that it is highly impactful because there is no other model-agnostic and domain-knowledge free method for improving  
4 the quality of local approximation explanations such as LIME (which is a seminal method in Interpretable ML).

5 **Prior Work:** The suggested related works (which we will cite in the revision) all solve different problems than the one  
6 we consider. We will add a discussion as outlined below.

7 • “Adversarial Robustness ...” by Qin et al does not consider interpretability at all. When adapted to consider  
8 interpretability, it uses a gradient based explanation and its regularizer is quite similar to SENN’s. Consequently, it  
9 will have the same issues with flexibility, fidelity, and stability as gradient based explanations. See A.2 for details.

10 • Several methods rely on domain knowledge: “Learning credible ...” by Du et al, “Learning Deep ...” by Weinberger  
11 et al, “Interpretations are ...” by Rieger et al, and “Regional Tree ...” by Wu et al [2].

12 • “Beyond sparsity ...” by Wu et al [1] regularizes for global interpretability while ExpO regularizes for local  
13 interpretability. Despite the fact that they are globally interpretable, small decision trees are difficult to explain locally  
14 with explainers like LIME (see Figure 1 for an example). As a result, [1,2] do not solve the same problem as ExpO  
15 because making the model look more like a decision tree makes LIME less effective.

16 **Reviewer 1. Reproducibility.** The reviewer is correct that we are comparing MLPs trained using standard techniques to  
17 ones trained with ExpO. We will add a detailed discussion of the neural networks (structure, activations, widths, depths,  
18 etc), hyper-parameters (learning rate, optimizer, regularization), and selection procedures to the appendix so that the  
19 reader does not have to reference the code (which reproduces all of our results) to reproduce our results.

20 **Reviewer 2.**

21 “*computational complexity ... cube ... not usable for higher-dimensional inputs.*” We introduce ExpO-1D-Fidelity to  
22 address this concern (line 160-167). Its complexity is independent of the data dimension and we show it scales well to  
23 datasets with  $\sim 100$  features. We also note that related methods require expensive operations (FTSD and [1] both are  
24 non-differentiable; SENN and RRR both require differentiating through the model gradient).

25 “*compare to RRR in this manner.*” To the best of our knowledge, it is not technically possible to encode fidelity/stability  
26 using RRR’s regularizer.

27 **Reviewer 3.**

28 “*difference between RRR, SENN is that a neighborhood ... is introduced ... not a huge difference.*” ExpO is the only  
29 method that is differentiable and model agnostic that does not require domain knowledge; the differences are not just in  
30 whether or not a neighborhood is used. See Table 1 for details.

31 “*Algorithm 1 ... not very novel.*” Viewing the novelty of ExpO merely through the lens of Algorithm 1 sells it short; the  
32 novelty stems from its impactful connection to interpretability. It is common for algorithms designed in one area to be  
33 impactful when introduced to another area (eg, SENN/RRR are “just” regularizing the gradient which is a strategy at  
34 least as old as “Tangent prop-a formalism for specifying selected invariances in an adaptive network.” NeurIPS92.)

35 “*results ... not very surprising ... idea is to \*optimize\* those metrics during learning.*” Two small clarifications: the  
36 results are shown for points that were not regularized for during training and the results shown in the main paper were  
37 regularized only for fidelity, so the improvement in stability is not a given.

38 “*why is the accuracy of SENN explanations...measured using a post-hoc explainers.*” While the reviewer is correct  
39 that SENN’s Point-Fidelity (PF) is perfect by-design, its Neighborhood-Fidelity (NF) is not guaranteed; the setup of  
40 user study clearly motivates why NF can be preferable to PF (lines 220 - 223). Following the reviewer’s suggestion,  
41 we computed NF and Stability for SENN explaining itself. While the results are better than using LIME, they  
42 corroborate the general message that ExpO is a more flexible solution than SENN for trading off between accuracy and  
43 interpretability. Specifically, SENN explaining itself has a NF of  $3.1e-5$  and a Stability of  $2.1e-3$ ; these numbers are  
44 generally comparable to LIME explaining the appropriate ExpO model. See A.1 and Table 5 for details.

45 “*regularize neural nets ... behave similarly to decision trees, either globally or regionally ... expect tree-regularized*  
46 *models to work well together with LIME...for both linear explanations and tree-based explanations.*” As noted in the  
47 above discussion on [1], neither of these methods would improve LIME’s explanation quality for linear explanations.  
48 Although we agree that exploring non-linear local explanations is an interesting direction, ExpO focuses on the setting  
49 where the explanation is linear because this is what LIME, MAPLE, and SENN all do.

50 **Reviewer 4:** “*experimental part is somewhat not convincing ... it is not surprising to see the results in user study: the*  
51 *regularized model achieves better interpretability than the normal model.*” Fidelity/stability are the standard proxy  
52 metrics used to evaluate local approximations. However, as we emphasize in the paper (line 38-42), they are only  
53 proxies for some underlying notion of interpretability, and the goal of the user study is to directly study explanation  
54 usefulness. Consequently, it inconsistent to criticize the results on the metrics and then use those same results to criticize  
55 the results of the user study.