1 We thank the reviewers for their comments and suggestions. We will incorporate the given technical suggestions in the
2 final version of the paper. Below we address the main concerns raised in the reviews.

3 **Adversarial Robustness application.** To avoid the widespread phenomenon of breaking allegedly robust training
4 methods shortly after their publication, we decided to further stress test our method with an assortment of adversarial
5 attacks, and found some vulnerabilities of our trained models to direct decision boundary (ddb) attacks and some
6 black-box attacks. Consequently, we restricted some of the Newton projections to be in the direction of PGD-found
7 examples. We performed extensive all vs. all PGD black-box attacks using both ddb and cross-entropy (Xent). Results
8 for MNIST are shown in Table; we log test accuracy where each column represents different attack, diagonal entries
9 are white-box and off-diagonal are black box attacks; right column shows worst-case for each training method (rows).
10 Note that this procedure came with the cost of a net decrease of our performance
11 for white-box attacks, however, we still remain SoTA or comparable. We will
12 update the paper accordingly (including CIFAR10 results) and tone down some
13 of the robustness claims.

| ddb/Xent | Xent | Madry | Trades | Our | Minimum |
|---|---|---|---|---|---|
| Madry | 98.6/98.6 | 96.1/96.3 | 97.9/97.9 | 98.7/99.1 | 96.1/96.3 |
| Trades | 98.6/98.6 | 98.5/98.5 | 96.9/96.7 | 98.8/99.2 | 96.9/96.7 |
| Our | 98.6/98.5 | 98.2/98.3 | 97.8/97.7 | 96.7/98.6 | 96.7/97.7 |

14 **(R1) "How well does the model converge? Is it guaranteed to find level sets through optimizing (3)?";**
15 **(R3) "There is no guarantee that the iteration in Eq. 4 would successfully sample a point on the level set."**
16 Newton's method is not guaranteed to find zeros of non-linear functions. Although ReLU networks do not satisfy the
17 conditions required for Newton's quadratic convergence it still works well in practice. Empirically, we applied ten
18 Newton iterations and converged to the zero level set between 80-90% of the times (manifold reconstruction and early
19 robust trainings) to 20-30% (end of robust training). Note that even when the Newton projection fails we can use it with
20 non zero $c$ (see Eq. (9)), which is useful for manifold reconstruction.

21 **(R1) "What is the practical speed of training the network due to that we have to get the level sets per iteration?"**
22 When comparing training times with level set sampling phase and without we get $\times 2$ the time for manifold reconstruction
23 and $\times 8$ for adversarial training.

24 **(R2) "Doesn't the ReLU activation imply that $D_x F(p; \theta)$ is often $= 0$ at many points p?"** The last layer is not
25 followed by a ReLU activation, so for $D_x F(p; \theta)$ to be 0 you need all of the neurons from the previous fully-connected
26 layer to be on the zero-region of their respective ReLU activations. Theoretically, if all weights are i.i.d. then chances
27 this happens is $0.5$ to the power of the number of neurons in previous to last layer. Empirically, this doesn't happen.

28 **(R2) "It seems one can sample from level 0 set by instead just optimizing: $\min_x \|F(x; \theta)\|$ via gradient descent**
29 **in x. Did the authors try this procedure?"** We have tried the suggested gradient descent (GD) procedure and found it
30 required two orders of magnitude more iterations than Newton projection to converge. Intuitively, the reason Newton is
31 much faster than GD for root finding is that GD linearizes the function at a point and takes a small step toward the zero
32 set, while Newton linearizes the function and goes all the way to the root of the linear function as the next step.

33 **(R3) "A good distribution of points on the level set should also account**
34 **for local geometry, e.g., curvature, which is not addressed in the proposed**
35 **method.".** This is indeed a good point (and a true challenge). From a practical
36 point of view we quantify the quality of distribution in low dimension (where
37 ground truth dense sampling of the level set is tractable). The table logs the Cham-
38 fer and Hausdorff distances of the resulting sampling distribution and the level
39 set of a neural network trained with Xent loss in 2 dimensions where projected

| Initialization Method | Chamfer | Hausdorf |
|---|---|---|
| Uniform [-0.35,0.35] | 0.011 | 0.141 |
| Normal $\sigma = 0.01$ | 0.006 | 0.017 |
| Normal $\sigma = 0.05$ | 0.01 | 0.132 |



(a) Normal $\sigma = 0.05$  (b) Uniform

40 points (red) are initialized using a uniformly distributed points (gray, right) or normally perturbed level set samples
41 (gray, left).

42 **(R3) "A sparse set of samples may not provide adequate control over the behavior of the en-**
43 **tire level set."** Indeed in high dimensions (i.e., not for surface and curve modeling) it would
44 be impossible to densely cover the entire level set with projections since its volume is too large.
45 However, our approach does move the entire level set in the desired manner due to the effect of
46 generalization that is manifested when optimizing a neural network with SGD. This is supported
47 empirically, e.g., the inset shows the histograms of distances of MNIST *test* samples to their
48 projection on the zero-level set of model trained by our method (orange) and a baseline (blue). Note
49 that distances evaluation on the test set means sampling the level set at unseen points.



50 **(R2) Conceptual discussion (and empirical comparison) on why the proposed approach should work better than**
51 **other strategies for large-margin deep.** We will add a comparison with a popular large-margin deep model, namely
52 level set linearization methods (e.g., Elsayed et al. [2018]). Conceptually, for $\| \cdot \|_2$, this method is equivalent to working
53 in our framework with a *single* Newton iteration providing only a crude approximation to the neural level set.