1 We would like to thank the reviewers for their positive and constructive feedback. We hope that with our detailed
2 answers below we can initiate a fruitful discussion that resolves all concerns.

3 **— R1.1): Validity of Theorem for large $\epsilon$. Experiment in Sec. 7.16.** R1 is correct, the Theorem establishes *exact*
4 equivalence only for $\epsilon$ small enough s.t. $\mathcal{B}_\epsilon^p(\mathbf{x}) \subset X(\phi_\mathbf{x})$. However, we experimentally verify that the correspondence
5 holds *approximately* (to a very good degree) in a region much larger than $X(\phi_\mathbf{x})$: see Sec. 5.4 "Validity of linear
6 approx.", specific. Fig. 4 (left), as well as Sec. 5.5 "Activation Patterns" specific. Fig. 5 and Fig. 9. Note, Fig. 9 shows the
7 *average* (incl. std. errors) over $\mathbf{x}$ and corresp. adv. $\mathbf{x}^*$ from the test set (not just a single data point). These experiments
8 confirm that $\mathbf{J}_f(\mathbf{x})$ is a good approx. (negligible deviation from linearity) and that the activation pattern change is small
9 ($\sim 3\%$ in adv., $\sim 1\%$ in random directions) within $\mathcal{B}_{\epsilon^*}^p(\mathbf{x})$ for *realistic* $\epsilon^*$ commonly used during AT.

10 **Further comments:** We will clarify the power iteration notation, expand the related work and move the Frobenius
11 norm regularization to the additional ninth page for the camera-ready version. Thank you for your high-quality review.

12 **— R2.1): Still not clear how such regularization affects the training of robust classifiers and how it characterizes**
13 **the sensitivity of the model to adversarial examples.** Our Theorem confirms that a network's sensitivity to adversarial
14 examples is characterized through its spectral properties: it is the dominant singular vector (resp. the maximizer $\mathbf{v}^*$
15 in Lemma 1) corresponding to the largest singular value (resp. the $(p, q)$-operator norm) that determines the optimal
16 adversarial perturbation and hence the sensitivity of the model to adversarial examples. The effect of AT and d.d. ONR
17 on the training of robust classifiers is twofold: (i) they dampen the singular values, see Sec. 5.2 and specific. Fig. 2 and
18 (ii) they give rise to models that are significantly more linear around data than normally trained ones, see Sec. 5.4 and
19 specific. Fig. 4 (left) as well as Sec. 5.5 specific. Fig. 5 and Fig. 9. Note also that our results directly explain why input
20 gradient regularization and FGM based AT do not sufficiently protect against iterative adversarial attacks, see Sec. 5.2
21 and in particular Sec. 7.8 in the Appendix. We will add a summary at the end of the theory section to emphasize these
22 contributions more thoroughly.

23 **R2.2) $\ell_q$-norm loss between logits of the clean and perturbed inputs not consistent with practice.** It is not
24 uncommon to use $\ell_q$-norm losses for adversarial training. See for instance: [1] Harini, Kurakin, and Goodfellow,
25 "Adversarial logit pairing" or [2] Sabour et al., "Adversarial manipulation of deep representations" (both use an $\ell_2$-norm
26 loss on the logits / internal representations). We will add additional pointers to clarify this.

27 Note also that we have a Sec. 7.7 "Cross-Entropy based Adversarial Loss Function" in the Appendix, where we discuss
28 the effect of the loss function on the directional derivative in the power-method like formulation of AT. In particular, the
29 gradient of the loss w.r.t. the logits of the classifier takes the same "prediction - target" form for both the sum-of-squares
30 error as well as the softmax cross-entropy loss, see "A note on canonical link functions". We will move some of this
31 discussion to the main part in the camera-ready version to emphasize this connection.

32 **R2.3) Theorem requires small $\epsilon$.** See reply to **R1.1)**.

33 **R2.4) Plot local linearity and activation patterns during training. Explore whether AT / d.d. ONR can help**
34 **stabilize activation pattern changes against adversarial examples.** It is clear from Figs. 5 & 9 that AT and d.d. ONR
35 do improve the stability of activation patterns against adversarial examples, as both d.d. ONR and AT significantly
36 increase the size of the ReLU cells (in both random and adv. directions) compared to the normally trained model.
37 During training, we see a gradual progression towards the end-state. We will add the plots to the camera-ready version.

38 **R2.5) Generalize Theorem 1 to allow for activation pattern changes.** Proving such an extension for the "approxi-
39 mate correspondence" between AT and d.d. ONR is *highly non-trivial* and thus out-of-scope of the current paper: one
40 would have to take into account how much "nearby" Jacobians can change based on the crossings of ReLU boundaries,
41 which is complicated by the fact that the impact of such crossings depends heavily on (i) the specific activation pattern
42 at input $\mathbf{x}$, (ii) the precise values of the weights and biases in the network, and (iii) where in the network the units that
43 change their state are. We will add a note to highlight these challenges.

44 **R2.6) What if we only consider the target logit? Then the Jacobian is a vector, will the theoretical result still**
45 **hold? How will the spectrum look?** This can be formulated as a special case of our analysis and our results would
46 still hold (the optimal perturbation would align with the "Jacobian vector" and there would only be one singular value).
47 Note, however, that in general, the effectiveness of adversarial perturbations depends on the relative increase of one
48 logit over the decrease of another.

49 **— R3.1) Paper only proves correspondence for adversarial training with an $l_q$-norm loss.** See reply to **R2.2)**

50 **R3.2) Further datasets are needed to validate the effectiveness of the theoretical analysis.** We have confirmed
51 all our experiments on SVHN and TinyImageNet. The results and conclusions hold as expected. We will add the
52 corresponding plots to the camera-ready version.

53 **R3.3) Formatting issues, e.g. in bibliography.** We will standardize the formatting. Thank you for the pointer.