

1 We thank the four reviewers for their careful reading, detailed feedback, and helpful comments. Below, we begin with
2 some clarifications on our contributions and minor extensions, and end with specific responses to each reviewer.

3 **Why SDP?** The adversarial examples problem is relatively unique in ML in that *global optimality* matters a lot.
4 (In training a model, this would just overfit the model.) A central weakness of attack algorithms based on local
5 optimization (e.g. PGD of Madry et al.) is that they cannot *prove* that an adversarial example they’ve found is globally
6 optimal—even if it really is globally optimal. On the other hand, provable guarantees of global optimality are certainly
7 obtainable if we accept exp-time (e.g. Katz et al.) or exp-time in the worst-case (e.g. Tjeng et al.). Within this context,
8 *SDP is interesting because it is the best tool for proving global optimality in poly-time*. Goemans and Williamson
9 revolutionized combinatorial optimization when they used SDP to prove near-global optimality bounds right at the
10 cliff edge of poly-time (assuming $P \neq NP$). Candes, Tao and their coauthors revolutionized compressed sensing
11 when they used SDP to solve problems previous thought NP-hard to guaranteed global optimality in poly-time.

12 **Our contributions.** Amongst the adversarial examples literature, we are the first to guarantee a globally optimal
13 certificate in polynomial time. Within this context, others have used SDP, and have obtained good empirical results,
14 but we are the first to give an end-to-end theoretical proof (of any kind) for SDP. *Our paper is a first step towards global
15 optimality in poly-time; our contribution is the proof technique to get there*. The standard technique of analyzing the
16 SDP dual (e.g. Candes and Tao) immediately runs into painful, possibly insurmountable issues. Instead, we developed
17 a nonconvex technique (Appendix A) that reduces the SDP into a sequence of possibly nonconvex projections. *Viewing
18 our primary contribution as the proof technique, we have taken meticulous care in communicating the technique in
19 a clear, clean, pedagogical way, by proving clear-cut results on simple examples*. We did this because we wanted
20 to make it as easy as possible for future researchers to build off of our work. This is especially important because
21 the naive SDP relaxation doesn’t work well on its own, but it has the potential to be made to work well by future
22 researchers. (Much like how Lasserre built on top of Sherali-Adams and Lovasz-Schrijver)

23 **Problem (A) vs Problem (B).** Our proof technique works equally well for both (A) and its *convex restriction* (B). But
24 (A) is almost always loose, so we cannot prove anything on (A) other than to state reasons for why the SDP is loose. In
25 comparison, problem (B) is tight for a sufficiently small ρ (it is essentially a trace regularization to induce a low-rank
26 solution) so we’re able to prove predictive bounds and verify them numerically. This is a fairly common route in SDP.

27 **Multiple layers.** Our proof technique easily extends to the $\ell > 1$ case, as we discussed in Appendix E. The resulting
28 problem (E.1) can be “unrolled” by recursively applying the one-layer argument. But the issue here is that after the
29 first layer, we begin projecting onto hyperbolas. Conditions for hyperbola-on-hyperbola to be collinear can be derived
30 but are difficult to interpret and verify (they are themselves LMIs). Nevertheless, we believe this is a direction future
31 researchers can build off, because LMIs can always be simplified by assuming structure.

32 **Reviewer 1.** We would like to emphasize the fact that our result is the *first proof of tightness* within this context.
33 Global optimality is exceedingly important within the context of adversarial examples for obvious reasons, and we
34 give the *first provable method* that attains global optimality in poly-time. Our contribution is in the proof technique
35 used to achieve this; we had to diverge substantially from the compressed sensing SDP literature to get here. Choice
36 of ℓ_2 oracle over ℓ_∞ oracle. Both oracles are common in the literature, but generate essentially the same adversarial
37 examples (see e.g. Carlini and Wagner). Lipschitz constants techniques are ℓ_2 methods that become considerably more
38 conservative on ℓ_∞ . SDP easily accommodate ℓ_∞ ; it is the only the theoretical analysis that becomes complicated.
39 Large ρ regime. In the regime of radius $\rho \rightarrow \infty$, we view (A) as (B) with $\hat{\mathbf{z}} = -\rho\mathbf{w}/\|\mathbf{w}\|$, but this means the
40 center of the ball $\hat{\mathbf{z}} \rightarrow \infty$ as well. Our tightness guarantees for (B) require $\hat{\mathbf{z}}$ to remain bounded. Quantative measure
41 of tightness. On the one-neuron example, the relaxation is tight if $|\langle \mathbf{e}, \mathbf{x} \rangle| = \|\mathbf{x}\|$. But if $|\langle \mathbf{e}, \mathbf{x} \rangle| < \|\mathbf{e}\|$, then the
42 incidence angle $\theta = \arccos(|\langle \mathbf{e}, \mathbf{x} \rangle|/\|\mathbf{e}\|)$ gives essentially the condition number of the high-rank solution. This is an
43 insight that only becomes clear through our proof technique; we promise to add this point to the paper.

44 **Reviewer 2.** We regret the cluttering noted by the reviewer. We endeavor to reduce clutter in a future revision.

45 **Reviewer 3.** We thank the reviewer for a number of key insights, and for catching several bugs in the Appendix. We
46 hope our common response written above adequately addresses the reviewer’s concerns regarding our contributions.
47 Proof of Lemma A.3: agreed. Line 706 should read $\hat{\mathbf{z}} = \mathbf{u} - \rho\mathbf{w}/\|\mathbf{w}\|$. The next line should read $\text{tr}(\mathbf{X}_\ell) - 2\langle \hat{\mathbf{z}}, \mathbf{x}_\ell \rangle +$
48 $\|\hat{\mathbf{z}}\|^2 - \rho^2 = \text{tr}(\mathbf{X}_\ell) - 2\langle \mathbf{u}, \mathbf{x}_\ell \rangle + \|\mathbf{u}\|^2 + 2\rho[\langle \mathbf{w}, \mathbf{z} \rangle + b]$ (the term $-\rho^2$ was lost). For nonnegativity, note that
49 $\text{tr}(\mathbf{X}) - 2\langle \mathbf{u}, \mathbf{x} \rangle + \|\mathbf{u}\|^2 = \text{tr}(\mathbf{X} - \mathbf{u}\mathbf{x}^T - \mathbf{x}\mathbf{u}^T + \mathbf{u}\mathbf{u}^T) = \text{tr}(\mathbf{X} - \mathbf{x}\mathbf{x}^T + (\mathbf{x} - \mathbf{u})(\mathbf{x} - \mathbf{u})^T) = \text{tr}(\mathbf{X} - \mathbf{x}\mathbf{x}^T) + \|\mathbf{x} - \mathbf{u}\|^2$ but
50 we have $\mathbf{X} - \mathbf{x}\mathbf{x}^T \succeq 0$ and therefore $\text{tr}(\mathbf{X} - \mathbf{x}\mathbf{x}^T) \geq 0$. The claim that (A-lb) is a relaxation of (B-lb) follows from the
51 corrected version of the equation above, which shows that a feasible point $\mathbf{X}_\ell, \mathbf{x}_\ell$ satisfying $\text{tr}(\mathbf{X}_\ell) - 2\langle \hat{\mathbf{z}}, \mathbf{x}_\ell \rangle + \|\hat{\mathbf{z}}\|^2 \leq$
52 ρ^2 would then immediately satisfy $2\rho[\langle \mathbf{w}, \mathbf{z} \rangle + b] \leq 0$. These points will be clarified and fixed.

53 **Reviewer 4.** We agree with the reviewer and promise to completely rewrite the introduction to better reflect our
54 contributions. Comparison to Raghunathan et al. This previous paper was almost entirely empirical. Global optimality
55 was not at all their focus; most of their paper was spent comparing SDP vs LP.