

1 We deeply appreciate constructive and insightful comments from the reviewers.

2 **Reviewer #1**

3 **Computing OCE minimizers.** On the one hand, we clarify that Theorems 3 and 6 apply to any approximate OCE
4 minimizer, including those given by [40], and Lemma 2 applies to an arbitrary algorithm, even including those not
5 minimizing OCE risks. Our results provide generalization guarantees with the same or tighter dependency for algorithms
6 given by [40]. We will clarify this point in the revised version. On the other hand, we agree to the reviewer’s point
7 that clearly stating the optimization-relevant properties of (inverted) OCE risks may facilitate future works in this
8 direction. To this end, we will survey existing results (in the appendix) that are relevant to the optimization of standard
9 risk-sensitive measures, e.g., smoothness/convexity properties of CVaR, entropic risk, and mean-variance.

10 **Prior work on stability & convergence.** We believe that R1 refers to (Line 42), where we point to [30]. The work
11 [30] considers an MDP setting that where the standard Bellman’s optimality principle cannot handle risk-sensitive
12 cases. By contrast, we focus on nonsequential scenarios where such concerns do not arise. We will clarify this point in
13 a separated “Related Works” section, along with detailed discussions on relevant literature.

14 **Reviewer #2**

15 **Apparent lack of tailored analysis for risk-sensitivity.** We begin by noting that the risk-sensitivity is indeed distilled
16 into the smoothness characteristics of the disutility $\phi(\cdot)$ in our results, instead of completely disappearing from the
17 analysis. The proposed OCE framework enables a deceptively simple treatment of risk-sensitivity by relating the
18 problem to risk-neutral learning via contraction principles. Such strategy (similar to what [D⁺19] did for adversarially
19 robust learning) allows us to fully utilize accumulated insights on risk-neutral learning to understand risk-sensitive
20 learning, without having to resort to overly complicated machineries (e.g., adapting Dvoretzky-Kiefer-Wolfowitz
21 inequality to consider data-dependent weights). Given the significance of the considered problem, we view this as an
22 *advantage* brought by our framework, instead of a weakness. Nevertheless, we also provide several “tailored” proof
23 tools, including *product hypothesis space analysis* (Lines 430–438) and *two-sided variance-based characterization of*
24 *OCE* (Appendix A.5), may be of readers’ technical interest. We will highlight these points clearly in the main text.

25 **Vacuity of Rademacher-style analysis for DNNs, and realizability.** We thank R2 for pointing this out. While recent
26 progress [S⁺20, N⁺20] shed new light on the power of Rademacher-style uniform convergence analysis for DNNs,
27 we agree to the point that the current tone of the manuscript may be too bold. We will revise the manuscript to avoid
28 exemplifying DNNs to make sense of realizability conditions.

29 **Related works and connection to fair ML.** Thank you for this valuable suggestion. We will add a standalone “Related
30 Works” section to provide in-depth comparisons with existing statistical learning literature. Also, we will establish
31 explicit connections to the fairness risk measures axiomatically defined in [49].

32 **Reviewer #3**

33 **Experiments other than CVaR (points 1&2).** Although our main scope is on a theoretical side than an algorithmic
34 side, we agree to the reviewer’s point that clarifying optimization properties of OCE risks may help readers grasp the
35 potentials of our framework. For this purpose, we will make the following two revisions: (1) We will give a pointer
36 to the concurrent work of Li et al. [L⁺20] (which appeared on arXiv after the submission deadline). Via large-scale
37 experiments [L⁺20], the authors empirically observe that *both entropic risk and its inverse* can be optimized efficiently
38 via standard mini-batch gradient descent, as R3 correctly anticipated. (2) For completeness, we will additionally survey
39 existing theoretical results (in the appendix) that are relevant to the optimization for standard risk-sensitive measures.

40 **Vacuity of Rademacher complexities.** We resonate with your concern that any argument regarding the generalization
41 properties of deep neural networks requires a delicate care. We will revise the manuscript to avoid exemplifying neural
42 networks to support the applicability of realizability conditions.

43 **Technical benefits.** As R3 keenly points out, our primary focus is to establish an effective theoretical framework to
44 formalize risk-sensitive learning, instead of pursuing technical exquisiteness. Nevertheless, we believe that several
45 proof techniques may be of theoreticians’ interest; for details, see Lines 22–24 of this response.

46 **Reviewer #4**

47 **Suggestions on discussions and proofs.** We express our deepest gratitude for the detailed comments, especially on the
48 ideas to refine our result and the pointer to L-statistics/estimators. All these valuable comments will be incorporated.

49 **Additional References**

50 [D⁺19] Y. Dong et al. Rademacher complexity for adversarially robust generalization. In *ICML*, 2019.

51 [L⁺20] T. Li et al. Tilted empirical risk minimization. arXiv preprint 2007.01162, 2020.

52 [N⁺20] J. Negrea et al. In defense of uniform convergence: Generalization via derandomization with an application to
53 interpolating predictors. In *ICML*, 2020.

54 [S⁺20] T. Suzuki et al. Compression based bound for non-compressed network: unified generalization error analysis
55 of large compressible deep neural network. In *ICLR*, 2020.