

1 We would like to thank all the four reviewers for their comments. In this document, we try to briefly respond to the  
2 concerns and questions raised by Reviewers 2 and 4.

3 **Reviewer 2: –(major comments)** We have tried to avoid making any claims about “explaining the help from unlabeled  
4 data”. Theorem 3 only provides a generalization bound. We revise the paper again to remove any remaining claims  
5 about this issue. However, using our framework, one can (at least theoretically) characterize cases where unlabeled  
6 data can provably help. First, it should be noted that by using only the labeled data for learning (as suggested by the  
7 reviewer), the residual generalization error in the classical learning framework would be  $O(n^{-1/2}\eta^{-1/2})$ . But residual  
8 error terms of Theorem 3 are both  $O(n^{-1/2})$  (note that  $\sqrt{\eta} + \sqrt{1-\eta} \leq \sqrt{2}$ ). Therefore, we can guarantee a much  
9 smaller residual error when supervision ratio is very small, i.e.  $\eta \ll 1$ . Second, for a highly compatible pair of model  
10 set  $\Phi$  and data distribution  $P_0$ , the condition  $\text{MSR}_{\Phi, P_0}(\lambda, \zeta) < \eta$  can be satisfied even for very small (and generally  
11 negative) values of  $\lambda$ . For a sufficiently small  $\lambda$ , our  $\hat{R}_{\text{SSAR}}$  becomes smaller than the average risk computed over only  
12 the labeled data. Let us discuss this matter, mathematically: For simplicity, assume the asymptotic case of  $n \rightarrow +\infty$   
13 (similar arguments hold for  $n < +\infty$ ). Then, with a little abuse of notation and for any  $\phi \in \Phi$ , we have:

$$\lim_{n \rightarrow +\infty} \hat{R}_{\text{SSAR}}(\phi; \mathbf{D}) \stackrel{a.s.}{=} \mathbb{E}_{\mathbf{X} \sim P_{0|\mathbf{X}}} \left\{ \eta \mathbb{E}_{y \sim P_{0|\mathbf{X}}} \{ \phi(\mathbf{X}, y) \} + (1 - \eta) \underset{y \in \mathcal{Y}}{\text{softmin}}^{(\lambda)} \{ \phi(\mathbf{X}, y) \} \right\}^* \leq \mathbb{E}_{\mathbf{X}, y \sim P_0} \{ \phi(\mathbf{X}, y) \},$$

14 where  $*$  holds for sufficiently small values of  $\lambda$ , since  $\mathbb{E}_{y \sim P_{0|\mathbf{X}}}$  is an expectation operator but  $\text{softmin}^{(\lambda)}$  can go as far  
15 as being the min operator. Therefore, one can establish a set of theoretical conditions under which unlabeled data is  
16 guaranteed to be helpful, since all the three terms in the r.h.s. of the bound in Theorem 3 become smaller than their  
17 traditional counterparts. The above-mentioned conditions are very general, but at the same time very implicit. In any  
18 case, we will add a lemma to our appendix to highlight this issue for interested readers.

19 **–(minor comments)** 1. Yes, our SSM measure can also be used when  $\epsilon = 0$  (i.e. no distributional robustness). To the  
20 best of our knowledge, there are no similar theoretical treatments of this problem in the existing works. 2. Please refer  
21 to our response to Reviewer 4.

22 **Reviewer 4:–(major comments)** Considering reviewer’s comments, first let us emphasize on some of our contributions  
23 that might have been missed during the review: we have tested our method on three different datasets and outperformed  
24 state-of-the-art in at least one of them. Also, we theoretically showed that a model set with a bounded VC-dimension is  
25 also *adversarially-learnable* (Lemma E.3), even in a *semi-supervised* scenario, where a corresponding generalization  
26 bound is given in Theorem 3. We agree with both Reviewers 2 and 4 that MSR in (C.11) is very implicit and hard  
27 to evaluate. However, please note that our framework is completely general, and thus providing a way to evaluate  
28 MSR in a general scenario might lead to solving several open problems in statistics (similar to providing a general  
29 way to evaluate VC-dimension or Rademacher complexity for any model set). For example, consider the loss function  
30 set  $\Phi = \{-\log P_\theta(\cdot, \cdot) \mid \theta \in \Theta\}$ , where  $P_\theta$  can be any parametric distribution family over  $\mathcal{X} \times \mathcal{Y}$ . Also, assume  
31 dataset is sampled from  $P_{\theta_0}$ , where  $\theta_0 \in \Theta$ . Then, it can be easily seen that the proposed risk in Theorem 1 when  
32  $\lambda = -1$ , is in fact the ML estimator (which is also the optimal estimator). Characterizing MSR in this case can shed  
33 light on the sample complexity of ML in a general semi-supervised setting which is still an open problem. However,  
34 let us give a quick example of how fast MSR can be computed in some very specific and simple cases: Assume the  
35 *cluster assumption*, where data distribution  $P_0$  is a mixture of two distributions whose supports do not overlap over  
36  $\mathcal{X}$ , and correspond to only  $y = -1$  and  $+1$  over  $\mathcal{Y}$ , respectively. Consider the loss function set  $\Phi$  which is associated  
37 with a family of arbitrary binary classifiers, where for each  $\phi \in \Phi$  we have  $\phi(\mathbf{X}, y) = \infty \cdot \phi_{\text{acc}}(\mathbf{X}, y) + \phi_{\text{mar}}(\mathbf{X})$ .  
38 Here,  $\phi_{\text{acc}} \in \{0, 1\}$  checks whether the label  $y$  matches with the positioning of  $\mathbf{X}$  w.r.t. the classifier of  $\phi$ , and  
39  $\phi_{\text{mar}}(\mathbf{X}) \in \mathbb{R}$  penalizes the margin, i.e. distance of  $\mathbf{X}$  from the classifier’s border. Now, let  $\psi \subseteq \Phi$  correspond to a  
40 subset of classifiers that classify all the data correctly ( $\mathbb{E}_{P_0} \phi_{\text{acc}} = 0$ ), but have different expected margins. Also, assume  
41  $\phi^*$  (the minimum loss associated with the optimal classifier) is also inside  $\psi$ . Then, some simple calculations reveal that  
42 for every  $\phi \in \psi$  and any  $\lambda$  we have  $\rho_\lambda(\phi) = 0$  (C.6) and thus  $\Lambda(\psi) = -\infty$  (C.10). Also, we have  $\Gamma(\psi; \lambda) \geq 0$  (C.9),  
43 again for any  $\lambda$ , while  $\text{GAP}(\psi)$  (C.9) is strictly positive for any non-trivial  $\Phi$  (recall that  $\phi^* \in \psi$ ). Considering the fact  
44 that we can have  $\zeta = O(n^{-1/2})$  according to Theorem 3, then for a sufficiently large  $n$ ,  $\text{MSR}_{\Phi, P_0}(\lambda, O(n^{-1/2}))$   
45 becomes zero for any  $\lambda \in \mathbb{R} \cup \pm\infty$ . This result is in full agreement with the previous bounds that are specifically  
46 derived for generic learnability of statistical models when non-overlapping cluster assumption holds (For absolute  
47 learnability, at least one data point with a label is needed to decide which cluster is which).

48 **–(minor comments)** 1. We have rephrased the sentences to avoid any possible confusions. 2. Yes,  $n = n_l + n_{\text{ul}}$ .  
49 Reviewer is correct and notations w.r.t.  $\mathbf{D}$  will be corrected. 5. The model used in our experiments is a deep neural  
50 network whose structure is completely explained in the supplementary document. Unfortunately, we cannot give more  
51 info in the main text due to the page limit. 3.4.6. We will correct all the grammatical mistakes, and also update the  
52 references.