

1 We thank the reviewers for their time and for consideration. They will find responses to their specific points below.

2 **Reviewer 2.** While we use fairness and robustness to illustrate the use of constraints in learning, the main results of the
3 paper (Thm. 1–3) are independent of these applications or the experimental settings. Still, we believe they are relevant.
4 Models invariant to gender, e.g., can be used for fair decision-making and, in a more data analysis context, to seek
5 patterns masked by gender or identify, through the dual variables, the features (phenotypes) more sensitive to gender.
6 We illustrate such analyses in Section D of the appendices. As the reviewer notes, this framework has applications
7 beyond model structure and preliminary results suggest that pointwise, per data point constraints can be used to improve
8 fit. However, this is beyond the scope of this work that focuses on understanding what can be achieved when learning
9 under constraints. Due to space limitations we could not address these points in the main body of the paper, but we
10 intend to use the extra page in the final manuscript to bring experimental details from the appendices. If this does not
11 fully address the reviewer concerns about the broader impact of the manuscript, we can work on expanding that section
12 for the camera-ready.

13 **Reviewer 3.** We agree that, being a new area of research, constrained learning theory still has many open questions,
14 such as understanding the effect of constraints on other learning models (e.g., structured complexity and PAC-Bayes)
15 and learning forms (e.g., reinforcement learning). We will include these discussions in the camera-ready. Still, this
16 work already shows that constrained learning is fundamentally different than PAC learning, especially dual constrained
17 learning where the dependence on μ and λ appears. Nevertheless, they affect only the *parametrization error* ϵ_0 ,
18 which is independent of the sample size (Def. 3). Hence, ϵ_0 is an intrinsic limitation of the parametrized learning
19 problem and sample sizes can always be estimated relative to this lower error bound. Alternatively, there is a well-
20 known upper bound from optimization theory, which for the case of $[0, B]$ -bounded objective yields $\|\mu^*\|_1 \leq Bs^{-1}$,
21 where $s = \min_i \mathbb{E}[\ell_0(f_{\theta'}(\mathbf{x}), y)] - c_i$ is the smallest slack for any strictly feasible solution θ' (i.e., such that $s > 0$).
22 If the output $\theta^{(T)}$ of Alg. 1 is sufficiently in the interior of the feasible set, bounds for the proof of Thm. 2 yield $s \leq$
23 $\min_i \left[c_i - M\nu - B\zeta(N_i) - \frac{1}{N_i} \sum_{n_i=1}^{N_i} \ell_i(f_{\theta^{(T)}}(\mathbf{x}_{n_i}), y_{n_i}) \right]_+$. While more tractable, we do not think that this result
24 is convenient as it requires the constraints to be quite loose for the bound not to vanish. We therefore did not include
25 these details in the manuscript, but will discuss this limitation more clearly in the camera-ready. Also, we use CSL
26 to mean “Constrained Statistical Learning” and write $\|\cdot\|_1$ to denote the ℓ_1 -norm (i.e., on a space of sequences) and $\|\cdot\|_{L_1}$
27 to denote the L_1 -norm (i.e., on a space of functions). This distinction occurs because the dual variable of the j -th
28 pointwise constraint is actually a continuous (as opposed to discrete) function.

29 **Reviewer 4.** The reviewer raises good points to improve the clarity of the paper. We intend to use the extra page of the
30 camera-ready to address them as well as bring part of the numerical experiments discussion from the appendices. In
31 particular, we will emphasize how the contributions of the paper provide generalization guarantees and therefore ensure
32 that an accurate, robust classifier trained using our constrained method will also perform well in testing (as evidenced
33 by Figures 4 and 5). Addressing the reviewer’s questions more specifically, the indices in (3) refer to the constraints
34 in (P-CSL), i.e., (3a) relates to the i -th average constraint of the learning problem and (3b) relates to the j -th pointwise
35 constraint. The reason we write “PACC Learning is as Hard as...” is because Thm. 1 is both necessary and sufficient.
36 That PAC is not harder than PACC is trivial since PACC includes PAC in (2). But Thm. 1 also proves the converse. We
37 then say these forms of learning are equivalent in the sense that a hypothesis class is PACC learnable if and only if it is
38 PAC learnable (Thm. 1). As for the relation between (PIV) and its dual (\hat{D} -CSL), recall that (PIV) is non-convex so that
39 its dual only provides a lower bound on its value (weak duality). However, Thm. 2 shows that (\hat{D} -CSL) provides both
40 lower and upper bounds for the original (P-CSL). In Def. 3, ϵ and ϵ_0 are not related: ϵ is the *estimation error* while ϵ_0 is
41 the *parametrization error*. Indeed, ϵ_0 is independent of the sample size and depends only on the learning problem (M
42 and dual variables) and the richness of the parametrization (ν). Although the parametrization error is treated separately
43 in unconstrained learning, it is not possible to untangle it for (dual) PACC learning (see proof of Thm. 2). In terms of
44 Alg. 1, its input is a sample set and its output is a parameter vector $\theta^{(T)}$ that describes the classifier. Dual variables can
45 also be extracted for analysis (see Section D). In these analyses, we leverage the sensitivity interpretation of the dual
46 variables, which is formalized by the fact that the optimal dual variables are *subgradients* of the empirical Lagrangian
47 at the optimal primal solution. In fact, we show that this is (approximately) true even at the approximate oracle from
48 Assumption 4 (see Lemma 1 in Section C). To be more specific, it means that the larger the dual variable, the more the
49 optimal value would increase if we were to tighten that constraint, indicating that the constraint is “harder” to satisfy.
50 We use a mini-batch SGD algorithm for step 3 of Alg. 1 (details can be found in the appendices). Note that the results
51 of Thm. 3 do not depend on this method and there may be better suited approaches in certain problems, but this was a
52 natural choice for the logistic and CNN classifiers used in the experiments. While we will certainly expand Section 2,
53 we note that, to the best of our knowledge, this is the first paper to explicitly analyze constraints in the context of (PAC)
54 learning theory. Previous work focus on specific constraints, e.g., rates, and provide guarantees for specific algorithms
55 that typically yield randomized solutions. In contrast, this paper provides algorithm-independent results (Thm. 1 and 2)
56 for deterministic learners. It also considers pointwise constraints that are important in domains such as fairness.