

Table 5. Key notation.

Symbol	Description
$m$	Number of constraints in Equation 1
$\ell_0$	Objective loss function in Equation 1
$\ell_i$	$i$ th constraint loss function in Equation 1
$\tilde{\ell}_i$	$i$ th proxy-constraint loss function (corresponding to $\ell_i$ ) in Definition 1
$x$	Feature vector for an example
$\mathcal{X}$	Space of feature vectors $x \in \mathcal{X}$
$\mathcal{D}$	Data distribution over feature vectors $x \in \mathcal{X}$
$S^{(\text{trn})}$	An <i>i.i.d.</i> training sample from $\mathcal{D}$
$S^{(\text{val})}$	An <i>i.i.d.</i> “validation” sample from $\mathcal{D}$
$\theta$	Parameter vector defining a model
$\Theta$	Space of parameter vectors ( $\theta \in \Theta$ )
$\hat{\Theta}$	Subset of $\Theta$ consisting of the sequence of algorithm outputs $\theta^{(1)}, \dots, \theta^{(T)}$
$\bar{\theta}$	Random variable over $\hat{\Theta}$ defining a stochastic model (Theorems 1 and 2)
$\Delta^m$	The $m$ -probability-simplex, i.e. the set of all $p \in \mathbb{R}_+^{m+1}$ for which $\sum_{i=1}^{m+1} p_i = 1$
$\lambda$	Vector of Lagrange-multiplier-like “hyperparameters” (Definition 1)
$\Lambda$	Space of Lagrange-multiplier-like “hyperparameter” vectors ( $\lambda \in \Lambda := \Delta^m \subseteq \mathbb{R}_+^{m+1}$ )
$\bar{\lambda}$	Average “hyperparameters” $\bar{\lambda} := (\sum_{t=1}^T \lambda^{(t)})/T$ (Theorems 1 and 2)
$\bar{\lambda}_1$	First coordinate of $\bar{\lambda}$ , measuring our “belief” that $\bar{\theta}$ is feasible
$\mathcal{L}_\theta$	In-expectation proxy-Lagrangian function minimized by $\theta$ -player (Definition 4)
$\mathcal{L}_\lambda$	In-expectation proxy-Lagrangian function maximized by $\lambda$ -player (Definition 4)
$\hat{\mathcal{L}}_\theta$	Empirical proxy-Lagrangian function minimized by $\theta$ -player (Definition 1)
$\hat{\mathcal{L}}_\lambda$	Empirical proxy-Lagrangian function maximized by $\lambda$ -player (Definition 1)
$\tilde{G}^{(\text{trn})}(\Theta)$	Generalization bound for $\ell_0, \tilde{\ell}_1, \dots, \tilde{\ell}_m$ on $S^{(\text{trn})}$ for $\theta \in \Theta$ (Definition 2)
$G^{(\text{val})}(\hat{\Theta})$	Generalization bound for $\ell_1, \dots, \ell_m$ on $S^{(\text{val})}$ for $\theta \in \hat{\Theta}$ (Definition 2)
$\mathcal{O}_\rho$	Bayesian oracle of Definition 3
$\rho$	Additive approximation of the Bayesian oracle of Definition 3
$C_r$	A radius- $r$ external covering of $\Lambda := \Delta^m \subseteq \mathbb{R}_+^{m+1}$ w.r.t. the 1-norm
$r$	Radius of the covering $C_r$
$\Theta_{C_r}$	Oracle evaluations at the covering centers ( $\Theta_{C_r} := \{\mathcal{O}_\rho(\hat{\mathcal{L}}_\theta(\cdot, \tilde{\lambda})) : \tilde{\lambda} \in C_r\}$ )
$T$	Number of iterations performed in Algorithms 1, 3 and 4
$T_\lambda$	Number of iterations performed in the outer loop of Algorithm 2
$T_\theta$	Number of iterations performed in the inner loop of Algorithm 2
$\eta_\lambda$	Step size associated with $\lambda$ -player in Algorithms 1, 2, 3 and 4
$\eta_\theta$	Step size associated with $\theta$ -player in Algorithms 3 and 4
$M$	A left-stochastic $(m+1) \times (m+1)$ matrix
$\mathcal{M}$	Space of all left-stochastic $(m+1) \times (m+1)$ matrices ( $M \in \mathcal{M}$ )
$\gamma$	Maximum margin by which the proxy-constraints can be satisfied in Theorems 1 and 2
$\mu$	Strong convexity parameter of $\ell_0, \tilde{\ell}_1, \dots, \tilde{\ell}_m$ in Theorem 2
$L$	Lipschitz constant of $\ell_1, \dots, \ell_m$ in Theorem 2
$B_{\ell_0}$	Upper bound on $b_0 - a_0$ , where $\text{range}(\ell_0) = [a_0, b_0]$
$B_\ell$	Upper bound on $b_i - a_i$ for all $i \in [m]$ , where $\text{range}(\ell_i) = [a_i, b_i]$
$B_{\tilde{\ell}}$	Upper bound on $ \ell(x; \theta) $ for all $\ell \in \{\ell_0, \tilde{\ell}_1, \dots, \tilde{\ell}_m\}$
$B_{\tilde{\Delta}}$	Upper bound on the 2-norms of subgradients of $\hat{\mathcal{L}}_\theta$ w.r.t. $\theta$
$B_{\Delta}$	Upper bound on the $\infty$ -norms of gradients of $\hat{\mathcal{L}}_\lambda$ w.r.t. $\lambda$

**Algorithm 3** “Practical” algorithm for optimizing the empirical proxy-Lagrangian game (Definition 1). This is essentially Algorithm 2 of Cotter et al. (2019), differing only in that it is applied to the two-dataset formulation of Definition 1.

---

PracticalTwoDataset  $(\hat{\mathcal{L}}_\theta, \hat{\mathcal{L}}_\lambda : \Theta \times \Delta^m \rightarrow \mathbb{R}, T \in \mathbb{N}, \eta_\theta, \eta_\lambda \in \mathbb{R}_+)$ :

- 1 Initialize  $\theta^{(1)} = 0$  // Assumes  $0 \in \Theta$
- 2 Initialize  $M^{(1)} \in \mathbb{R}^{(m+1) \times (m+1)}$  with  $M_{i,j} = 1/(m+1)$
- 3 For  $t \in [T]$ :
- 4     Let  $\lambda^{(t)} = \text{fix } M^{(t)}$  // fixed point of  $M^{(t)}$ , i.e. a stationary distribution
- 5     Let  $\check{\Delta}_\theta^{(t)}$  be a stochastic subgradient of  $\hat{\mathcal{L}}_\theta(\theta^{(t)}, \lambda^{(t)})$  w.r.t.  $\theta$
- 6     Let  $\check{\Delta}_\lambda^{(t)}$  be a stochastic gradient of  $\hat{\mathcal{L}}_\lambda(\theta^{(t)}, \lambda^{(t)})$  w.r.t.  $\lambda$
- 7     Update  $\theta^{(t+1)} = \Pi_\Theta(\theta^{(t)} - \eta_\theta \check{\Delta}_\theta^{(t)})$
- 8     Update  $\tilde{M}^{(t+1)} = M^{(t)} \odot \cdot \exp(\eta_\lambda \Delta_\lambda^{(t)} (\lambda^{(t)})^T)$  //  $\Delta \lambda^T$  is an outer product;  $\odot$  and  $\cdot \exp$  are element-wise
- 9     Project  $M_{:,i}^{(t+1)} = \tilde{M}_{:,i}^{(t+1)} / \|\tilde{M}_{:,i}^{(t+1)}\|_1$  for  $i \in [m+1]$  // Column-wise projection w.r.t. KL divergence
- 10 Return  $\theta^{(1)}, \dots, \theta^{(T)}$  and  $\lambda^{(1)}, \dots, \lambda^{(T)}$

---

**Algorithm 4** “Practical” algorithm for optimizing a variant of the standard Lagrangian game, modified to support proxy constraints and two datasets. Here, instead of using the proxy-Lagrangian formulation of Definition 1, we take  $\hat{\mathcal{L}}_\theta := \mathbb{E}_{x \sim S(\text{trn})} [\ell_0(x; \theta) + \sum_{i=1}^m \lambda_i \tilde{\ell}_i(x; \theta)]$  and  $\hat{\mathcal{L}}_\lambda := \mathbb{E}_{x \sim S(\text{val})} [\ell_0(x; \theta) + \sum_{i=1}^m \lambda_i \ell_i(x; \theta)]$ , with  $\lambda \in \Lambda := \mathbb{R}_+^m$ . Compared to Algorithm 3, this algorithm is further from those for which we can prove theoretical results (Algorithms 1 and 2), but is much closer to the Lagrangian-based approach of Agarwal et al. (2018). We include it to demonstrate that our two-dataset proposal works *as a heuristic*.

---

LagrangianTwoDataset  $(\hat{\mathcal{L}}_\theta, \hat{\mathcal{L}}_\lambda : \Theta \times \mathbb{R}_+^m \rightarrow \mathbb{R}, T \in \mathbb{N}, \eta_\theta, \eta_\lambda \in \mathbb{R}_+)$ :

- 1 Initialize  $\theta^{(1)} = 0, \lambda^{(1)} = 0$  // Assumes  $0 \in \Theta$
- 2 For  $t \in [T]$ :
- 3     Let  $\check{\Delta}_\theta^{(t)}$  be a stochastic subgradient of  $\hat{\mathcal{L}}_\theta(\theta^{(t)}, \lambda^{(t)})$  w.r.t.  $\theta$
- 4     Let  $\Delta_\lambda^{(t)}$  be a stochastic gradient of  $\hat{\mathcal{L}}_\lambda(\theta^{(t)}, \lambda^{(t)})$  w.r.t.  $\lambda$
- 5     Update  $\theta^{(t+1)} = \Pi_\Theta(\theta^{(t)} - \eta_\theta \check{\Delta}_\theta^{(t)})$  // Projected SGD updates (w.r.t. the Euclidean norm)...
- 6     Update  $\lambda^{(t+1)} = \Pi_\Lambda(\lambda^{(t)} + \eta_\lambda \Delta_\lambda^{(t)})$  // ...
- 7 Return  $\theta^{(1)}, \dots, \theta^{(T)}$  and  $\lambda^{(1)}, \dots, \lambda^{(T)}$

---

## A. Glossary

**(External) Regret:** Suppose we have a sequence of losses  $\ell_t : \Theta \rightarrow \mathbb{R}$  that we wish to minimize, and an algorithm that produces a sequence  $\theta^{(1)}, \dots, \theta^{(T)} \in \Theta$ . We can measure the performance of this algorithm via its *external regret*:

$$R = \sum_{t=1}^T \ell_t(\theta^{(t)}) - \min_{\theta^* \in \Theta} \sum_{t=1}^T \ell_t(\theta^*)$$

External regret is the difference between the total loss of the sequence generated by the algorithm, and that of the best element of  $\Theta$ , chosen with the benefit of hindsight.

**Swap Regret:** To define swap regret, we must adjust the above setting slightly. In particular, we will assume that at each step we may select from a *finite* number of choices—we’ll call this number  $m$ —and define each  $\ell_t \in \mathbb{R}^m$  as a vector for which its  $i$ th element is the loss associated with the  $i$ th choice. The algorithm we are evaluating will be assumed to produce a sequence of probability vectors  $p^{(1)}, \dots, p^{(T)} \in \Delta^{m-1} \subseteq \mathbb{R}_+^m$ , for which the external regret would be:

$$R = \sum_{t=1}^T \langle p^{(t)}, \ell_t \rangle - \min_{p^* \in \Delta^{m-1}} \sum_{t=1}^T \langle p^*, \ell_t \rangle$$

Unlike external regret, which compares to the best constant strategy, *swap regret* compares to a strategy that varies in a particular way: we can retroactively decide that the probability mass that the  $p^{(t)}$  sequence assigned to the  $i$ th element should instead be assigned to the  $i'$ th, and the  $j$ th to the  $j'$ th, and so on. Mathematically, swap regret can be defined as:

$$R = \sum_{t=1}^T \langle p^{(t)}, \ell_t \rangle - \min_{M \in \mathcal{M}} \sum_{t=1}^T \langle Mp^{(t)}, \ell_t \rangle$$

here,  $\mathcal{M} \subseteq \mathbb{R}^{m \times m}$  is the space of left-stochastic matrices.

**Pure Equilibrium:** Consider a two-player game with payoff functions  $\mathcal{L}_\theta : \Theta \times \Lambda \rightarrow \mathbb{R}$  and  $\mathcal{L}_\lambda : \Theta \times \Lambda \rightarrow \mathbb{R}$ . We'll adopt the convention that the first player wishes to choose  $\theta$  to minimize  $\mathcal{L}_\theta$ , while the second wishes to choose  $\lambda$  to maximize  $\mathcal{L}_\lambda$ . If  $\mathcal{L}_\theta = \mathcal{L}_\lambda$ , then the game is said to be *zero-sum*. Otherwise, it is non-zero-sum. In either case, a *pure equilibrium* is a pair  $\theta^*, \lambda^* \in \Theta \times \Lambda$  satisfying some properties (depending on the type of equilibrium). For example, for a *pure Nash equilibrium*, both players suffer no external regret:

$$\begin{aligned} \mathcal{L}_\theta(\theta^*, \lambda^*) &\leq \mathcal{L}_\theta(\theta, \lambda^*) \\ \mathcal{L}_\lambda(\theta^*, \lambda^*) &\geq \mathcal{L}_\lambda(\theta^*, \lambda) \end{aligned}$$

the above holding for all  $\theta \in \Theta$  and  $\lambda \in \Lambda$ . In words, the  $\theta$ -player cannot improve upon  $\theta^*$ , and the  $\lambda$ -player cannot improve upon  $\lambda^*$ .

**Mixed Equilibrium:** A *mixed equilibrium* differs from a pure equilibrium in that, instead of being a pair of elements taken from  $\Theta \times \Lambda$ , it is defined either as a joint distribution  $P$  supported on  $\Theta \times \Lambda$ , or—when possible—as a pair of marginal distributions  $P_\theta$  and  $P_\lambda$  supported on  $\Theta$  and  $\Lambda$ , respectively (with  $P := P_\theta \times P_\lambda$ ). For example, a *mixed Nash equilibrium* is a pair of marginal distributions for which both players suffer no external regret:

$$\begin{aligned} \mathbb{E}_{\theta^* \sim P_\theta, \lambda^* \sim P_\lambda} [\mathcal{L}_\theta(\theta^*, \lambda^*)] &\leq \mathbb{E}_{\lambda^* \sim P_\lambda} [\mathcal{L}_\theta(\theta, \lambda^*)] \\ \mathbb{E}_{\theta^* \sim P_\theta, \lambda^* \sim P_\lambda} [\mathcal{L}_\theta(\theta^*, \lambda^*)] &\geq \mathbb{E}_{\theta^* \sim P_\theta} [\mathcal{L}_\theta(\theta^*, \lambda)] \end{aligned}$$

the above holding, as before, for all  $\theta \in \Theta$  and  $\lambda \in \Lambda$ .

**Coarse Correlated Equilibrium:** A *coarse correlated equilibrium* is a joint distribution  $P$  for which both players suffer no external regret:

$$\begin{aligned} \mathbb{E}_{\theta^*, \lambda^* \sim P} [\mathcal{L}_\theta(\theta^*, \lambda^*)] &\leq \mathbb{E}_{\theta^*, \lambda^* \sim P} [\mathcal{L}_\theta(\theta, \lambda^*)] \\ \mathbb{E}_{\theta^*, \lambda^* \sim P} [\mathcal{L}_\theta(\theta^*, \lambda^*)] &\geq \mathbb{E}_{\theta^*, \lambda^* \sim P} [\mathcal{L}_\theta(\theta^*, \lambda)] \end{aligned}$$

the above holding, once again, for all  $\theta \in \Theta$  and  $\lambda \in \Lambda$ . A coarse correlated equilibrium is *weaker* than a mixed Nash equilibrium: the former is a joint distribution  $P$  supported on  $\Theta \times \Lambda$ , while the latter can be decomposed as the product of two marginal distributions ( $P := P_\theta \times P_\lambda$ ).

In Theorems 1 and 2, the relevant type of mixed equilibrium—which does not have an established name<sup>¶</sup>—is slightly stronger than a coarse correlated equilibrium, but is still easier to find than a mixed Nash equilibrium: it is a joint distribution  $P$  for which the  $\theta$ -player suffers no external regret, while the  $\lambda$ -player suffers no *swap* regret. For further details, see Theorem 3 in Appendix E, or Theorem 8 of Cotter et al. (2019).

## B. Examples of Constraints

In this appendix, we'll provide some examples of constrained optimization problems in the form of Equation 1.

### B.1. Neyman-Pearson

The first example we'll consider is Neyman-Pearson classification (Davenport et al., 2010; Gasso et al., 2011). Imagine that we wish to learn a classification function  $f : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$  parameterized by  $\theta$ , with the goal being to minimize the false

<sup>¶</sup>We sometimes call it a *semi-coarse correlated equilibrium*, since a *correlated equilibrium* is a joint distribution for which all players suffer no swap regret.

positive rate, subject to the constraint that the false negative rate be at most 10%:

$$\begin{aligned} & \underset{\theta \in \Theta}{\text{minimize}} \mathbb{E}_{x,y|y=-1} [\mathbf{1}\{f(x;\theta) \geq 0\}] \\ & \text{s.t. } \mathbb{E}_{x,y|y=1} [\mathbf{1}\{f(x;\theta) \leq 0\}] \leq 0.1 \end{aligned}$$

One way to convert this problem into the form of Equation 1 is to define  $\mathcal{D}_+$  and  $\mathcal{D}_-$  as the marginal distributions over  $x$ s for which  $y = +1$  and  $y = -1$  (respectively), and take  $\mathcal{D} := \mathcal{D}_+ \times \mathcal{D}_-$  so that  $\mathcal{D}$  is a distribution over *pairs* of feature vectors, the first having a positive label, and the second a negative label. Defining  $\ell_0(x_+, x_-; \theta) := \mathbf{1}\{f(x_-) \geq 0\}$  and  $\ell_1(x_+, x_-; \theta) := \mathbf{1}\{f(x_+) \leq 0\} - 0.1$  puts the original Neyman-Pearson problem in the form of Equation 1.

In practice, the fact that  $\ell_0$  and  $\ell_1$  are defined in terms of indicator functions, and are therefore discontinuous, will be problematic. To fix this, using the formulation of Definition 1, one could instead define  $\ell_0(x_+, x_-; \theta) := \max\{0, 1 + f(x_-)\}$  as a hinge upper bound on the false positive rate, and leave  $\ell_1$  as-is while defining the corresponding proxy-constraint to be  $\tilde{\ell}_1(x_+, x_-; \theta) := \max\{0, 1 - f(x_+)\}$ .

## B.2. Equal Opportunity

The second example we'll consider is a fairness-constrained problem. As before, we'll take  $f : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$  to be a classification function, but we'll imagine that each  $x \in \mathcal{X}$  contains a feature  $x_k \in \{1, 2, 3\}$  indicating to which of three protected classes the corresponding example belongs. We will seek to minimize the overall error rate, subject to the constraint that, for each of the three protected classes, the false negative rate is at most 110% of the false negative rate across all three classes (this is essentially an equal opportunity constraint (Hardt et al., 2016)):

$$\begin{aligned} & \underset{\theta \in \Theta}{\text{minimize}} \mathbb{E}_{x,y} [\mathbf{1}\{yf(x;\theta) \leq 0\}] \\ & \text{s.t. } \forall i \in \{1, 2, 3\} \mathbb{E}_{x,y|y=1 \wedge x_k=i} [\mathbf{1}\{f(x;\theta) \leq 0\}] \leq 1.1 \cdot \mathbb{E}_{x,y|y=1} [\mathbf{1}\{f(x;\theta) \leq 0\}] \end{aligned}$$

While we could use the same approach as in the Neyman-Pearson example, i.e. taking marginals and crossing them to define a data distribution over tuples of examples, we'll instead take  $\mathcal{D}$  to be the data distribution over  $\mathcal{X} \times \{\pm 1\}$  pairs, and use the indicator feature  $x_k$  to define:

$$\begin{aligned} \ell_0(x, y; \theta) & := \mathbf{1}\{yf(x;\theta) \leq 0\} \\ \ell_i(x, y; \theta) & := \frac{\mathbf{1}\{y = 1 \wedge x_k = i\} \mathbf{1}\{f(x;\theta) \leq 0\}}{\Pr\{y = 1 \wedge x_k = i \mid x, y \sim \mathcal{D}\}} - 1.1 \cdot \frac{\mathbf{1}\{y = 1\} \mathbf{1}\{f(x;\theta) \leq 0\}}{\Pr\{y = 1 \mid x, y \sim \mathcal{D}\}} \end{aligned}$$

for all  $i \in \{1, 2, 3\}$ , where we assume that the probabilities in the denominators of the ratios defining  $\ell_i$  are constants known a priori.

As in the Neyman-Pearson example, in practice the indicators in the objective function could be replaced with differentiable upper bounds, and a differentiable proxy-constraint  $\tilde{\ell}_i$  could be introduced for each  $\ell_i$ .

## C. Shrinking

Cotter et al. (2019) introduced a procedure for “shrinking” the support size of a  $\bar{\theta}$ . When adapted to our setting, the first step is to evaluate the objective and constraints for every iterate (in practice, this is overkill; one should subsample the iterates):

$$\begin{aligned} \vec{\ell}_0^{(t)} & = \frac{1}{|S^{(\text{trn})}|} \sum_{x \in S^{(\text{trn})}} \ell_0(x; \theta^{(t)}) \\ \vec{\ell}_i^{(t)} & = \frac{1}{|S^{(\text{val})}|} \sum_{x \in S^{(\text{val})}} \ell_i(x; \theta^{(t)}) \end{aligned}$$

Next, we optimize a linear program (LP) that seeks a distribution  $p$  over  $\hat{\Theta}$  that minimizes the objective while violating no constraint by more than  $\epsilon$ :

$$\min_{p \in \Delta^{T-1} \subseteq \mathbb{R}_+^T} \langle p, \vec{\ell}_0 \rangle \quad \text{s.t.} \quad \langle p, \vec{\ell}_i \rangle \leq \epsilon \quad \forall i \in [m]$$

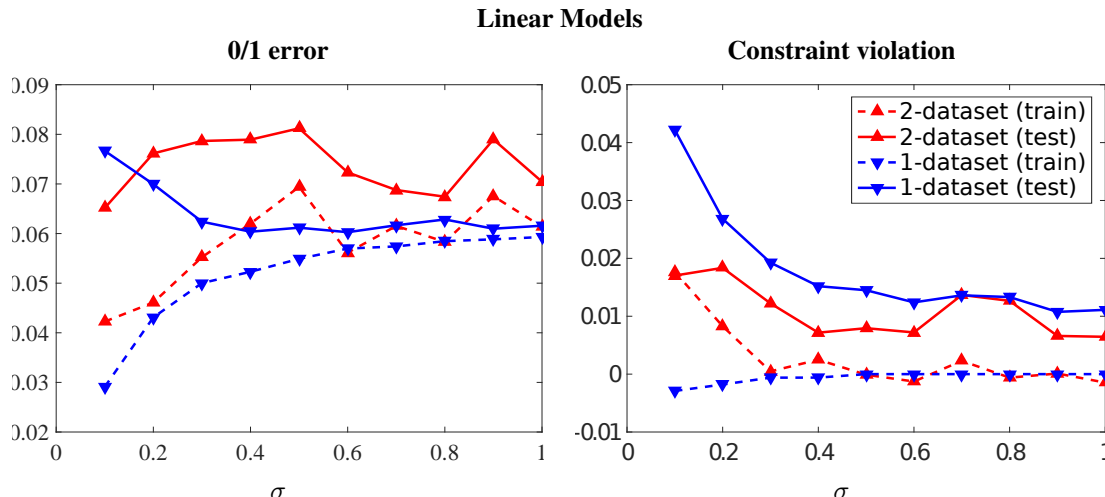


Figure 1. Results of the experiment of Appendix D, using Algorithm 3 on a linear model. The left-hand plot demonstrates both that overfitting increases as  $\sigma$  decreases, as expected, and that our model performs worse than the baseline, in terms of accuracy. In the right-hand plot, however, we can see that our approach generalizes better on the constraints (the two red curves are closer together than the two blue ones), and also does a better job of satisfying the constraints on the testing set (the solid red curve is below the solid blue one).

Notice that the  $p \in \Delta^{T-1}$  condition adds several implicit simplex constraints to this LP. The key to this procedure is that, as Cotter et al. (2019) show, every vertex  $p$  of this linear program has at most  $m + 1$  nonzero elements, where  $m$  is the number of constraints.

In particular, if  $\epsilon$  is chosen to be the maximum validation constraint violation of the “original” stochastic classifier  $\bar{\theta}$ , then an *optimal* vertex  $p^*$  will have a training objective function value and maximum validation constraint violation that are no larger than those of  $\bar{\theta}$ , and  $p^*$  will be supported on only  $m + 1$   $\theta^{(t)}$ s. Furthermore, since the resulting stochastic classifier is still supported on a subset of  $\bar{\Theta}$ , the generalization definitions of Section 3.1 apply to it just as well as they did to  $\bar{\theta}$ .

While it would be possible to provide optimality and feasibility guarantees for the result of this “shrinking” procedure, we will only use it in our experiments (Section 5). There, instead of taking the LP’s  $\epsilon$  parameter to be the maximum validation constraint violation of  $\bar{\theta}$ , we follow Cotter et al. (2019)’s suggestion to use an outer bisection search to find the smallest  $\epsilon \geq 0$  for which the LP is feasible.

## D. Simulated-data Experiments

These experiments are performed on a simulated binary classification problem designed to be especially prone to overfitting, and is intended to explore the relationship between generalization and model complexity in the context of our proposed approach.

To generate the dataset, we first draw  $n = 1000$  points  $z_1, \dots, z_n$  from two overlapping Gaussians in  $\mathbb{R}^2$ , and another  $n$  points  $w_1, \dots, w_n$  from the same distribution. For each  $i$ , we let the classification label  $y_i$  indicate which of the two Gaussians  $z_i$  was drawn from, and generate a feature vector  $x_i \in \mathcal{X} := \mathbb{R}^n$  such that the  $j$ th feature satisfies  $x_{i,j} := \exp(-\|z_i - w_j\|^2 / 2\sigma^2)$ . Our results are averaged over ten runs, with different random splits of the data into equally-sized training, validation and testing datasets.

The classification task is learn a classifier on  $\mathcal{X}$  that determines which of the two Gaussian distributions generated the example, with the model’s recall constrained to be at least 97%. The  $\sigma$  parameter partly controls the amount of overfitting: as  $\sigma \rightarrow 0$ , a linear classifier on  $\mathcal{X}$  approaches a 1-nearest neighbor classifier over  $w_1, \dots, w_n$ , which one would expect to overfit badly.

We trained four sets of models using Algorithm 3: linear, and one-hidden-layer neural networks with 5, 10 and 100 hidden ReLU units. We also varied  $\sigma$  between 0 and 1. Figures 1 and 2 show that our approach consistently comes closer to satisfying the constraints on the testing set—and the training and testing constraint violations are much closer to each other for our approach than for the baseline—but that, as one would expect, this comes at a slight cost in testing accuracy. In light

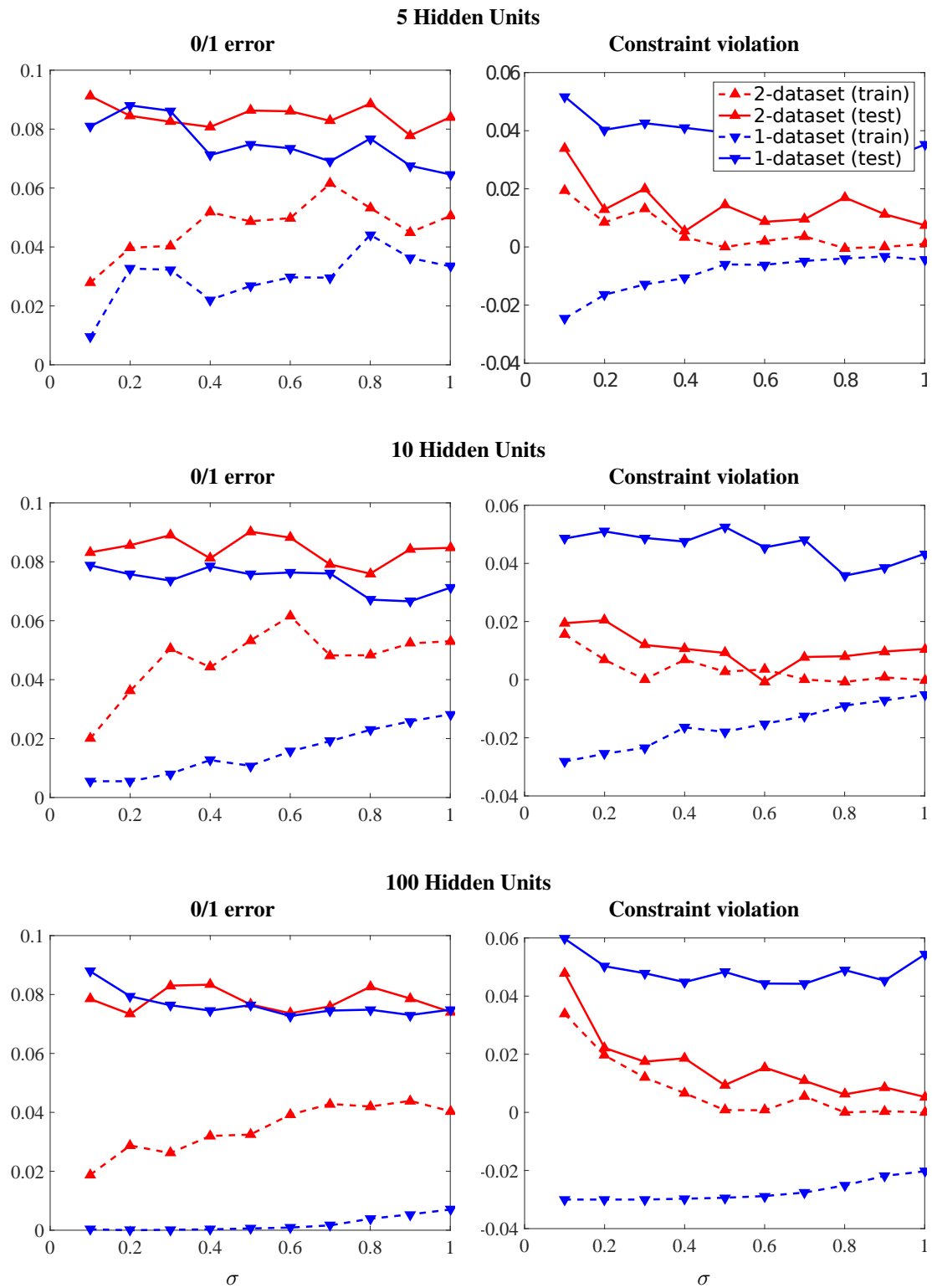


Figure 2. Same as Figure 1, but for one-hidden-layer neural networks with 5, 10 and 100 hidden units. The results follow the same general trend as those of Figure 1, but, just as one would expect, the benefits of our proposed two-dataset approach become more pronounced as the model complexity increases.

of the fact that our theoretical results (Section 4) bound constraint generalization independently of model complexity, it's unsurprising that our approach is most advantageous for the most complex models (100-hidden unit), and less so for the simplest (linear).

## E. Proofs

We'll begin by reproducing a definition from Cotter et al. (2019), which introduced the idea of the in-expectation proxy-Lagrangian game:

**Definition 4. (Definition 2 of Cotter et al. (2019))** Given proxy loss functions  $\tilde{\ell}_i(x; \theta) \geq \ell_i(x; \theta)$  for all  $x \in \mathcal{X}$  and  $i \in [m]$ , the proxy-Lagrangians  $\mathcal{L}_\theta, \mathcal{L}_\lambda : \Theta \times \Lambda \rightarrow \mathbb{R}$  of Equation 1 are:

$$\begin{aligned}\mathcal{L}_\theta(\theta, \lambda) &:= \mathbb{E}_{x \sim \mathcal{D}} \left[ \lambda_1 \ell_0(x; \theta) + \sum_{i=1}^m \lambda_{i+1} \tilde{\ell}_i(x; \theta) \right] \\ \mathcal{L}_\lambda(\theta, \lambda) &:= \mathbb{E}_{x \sim \mathcal{D}} \left[ \sum_{i=1}^m \lambda_{i+1} \ell_i(x; \theta) \right]\end{aligned}$$

where  $\Lambda := \Delta^m \subseteq \mathbb{R}_+^{m+1}$  is the  $m$ -probability-simplex.

This definition differs from Definition 1 (in Section 3) in that the former is the in-expectation version of the latter, which additionally is written in terms of separate *i.i.d.* training and validation sets.

Theorem 2 of Cotter et al. (2019) characterizes the optimality and feasibility properties of a particular type of  $\Phi$ -correlated equilibrium of Definition 4. We adapt it to our setting, giving the analogous result for such an equilibrium of Definition 1:

**Theorem 3.** Define  $\mathcal{M}$  as the set of all left-stochastic  $(m+1) \times (m+1)$  matrices,  $\Lambda := \Delta^m \subseteq \mathbb{R}_+^{m+1}$  as the  $m$ -probability-simplex, and assume that each  $\tilde{\ell}_i$  upper bounds the corresponding  $\ell_i$ . Let  $\theta^{(1)}, \dots, \theta^{(T)} \in \Theta$  and  $\lambda^{(1)}, \dots, \lambda^{(T)} \in \Lambda$  be sequences satisfying:

$$\begin{aligned}\frac{1}{T} \sum_{t=1}^T \hat{\mathcal{L}}_\theta(\theta^{(t)}, \lambda^{(t)}) - \inf_{\theta^* \in \Theta} \frac{1}{T} \sum_{t=1}^T \hat{\mathcal{L}}_\theta(\theta^*, \lambda^{(t)}) &\leq \epsilon_\theta \\ \max_{M^* \in \mathcal{M}} \frac{1}{T} \sum_{t=1}^T \hat{\mathcal{L}}_\lambda(\theta^{(t)}, M^* \lambda^{(t)}) - \frac{1}{T} \sum_{t=1}^T \hat{\mathcal{L}}_\lambda(\theta^{(t)}, \lambda^{(t)}) &\leq \epsilon_\lambda\end{aligned}\tag{8}$$

Define  $\hat{\Theta} := \{\theta^{(1)}, \dots, \theta^{(T)}\}$ . Let  $\bar{\theta}$  be a random variable taking values from  $\hat{\Theta}$ , defined such that  $\bar{\theta} = \theta^{(t)}$  with probability  $\lambda_1^{(t)} / \sum_{s=1}^T \lambda_1^{(s)}$ , and let  $\bar{\lambda} := (\sum_{t=1}^T \lambda^{(t)}) / T$ . Then  $\bar{\theta}$  is nearly-optimal in expectation:

$$\begin{aligned}\mathbb{E}_{\bar{\theta}, x \sim \mathcal{D}} [\ell_0(x; \bar{\theta})] &\leq \inf_{\theta^* \in \Theta: \forall i, \mathbb{E}_{x \sim \mathcal{D}}[\tilde{\ell}_i(x; \theta^*)] \leq 0} \mathbb{E}_{x \sim \mathcal{D}} [\ell_0(x; \theta^*)] \\ &\quad + \frac{1}{\lambda_1} (\epsilon_\theta + \epsilon_\lambda + 2\tilde{G}^{(\text{trn})}(\Theta) + G^{(\text{val})}(\hat{\Theta}))\end{aligned}\tag{9}$$

and nearly-feasible:

$$\max_{i \in [m]} \mathbb{E}_{\bar{\theta}, x \sim \mathcal{D}} [\ell_i(x; \bar{\theta})] \leq \frac{\epsilon_\lambda}{\lambda_1} + G^{(\text{val})}(\hat{\Theta})\tag{10}$$

Additionally, if there exists a  $\theta^l \in \Theta$  that satisfies all of the constraints with margin  $\gamma$  (i.e.  $\mathbb{E}_{x \sim \mathcal{D}} [\ell_i(x; \theta^l)] \leq -\gamma$  for all  $i \in [m]$ ), then:

$$\bar{\lambda}_1 \geq \frac{\gamma - \epsilon_\theta - \epsilon_\lambda - 2\tilde{G}^{(\text{trn})}(\Theta) - G^{(\text{val})}(\hat{\Theta})}{\gamma + B_{\ell_0}}\tag{11}$$

where  $B_{\ell_0} \geq \sup_{\theta \in \Theta} \mathbb{E}_{x \sim \mathcal{D}} [\ell_0(x; \theta)] - \inf_{\theta \in \Theta} \mathbb{E}_{x \sim \mathcal{D}} [\ell_0(x; \theta)]$  is a bound on the range of the objective loss.

*Proof.* This proof closely follows those of Theorem 4 and Lemma 7 in Cotter et al. (2019)—the only addition is that, in this proof, we use the empirical proxy-Lagrangian formulation with separate training and validation datasets (Definition 1), instead of the in-expectation proxy-Lagrangian (Definition 4), and therefore need to account for generalization.

**Optimality:** If we choose  $M^*$  to be the matrix with its first row being all-one, and all other rows being all-zero, then  $\hat{\mathcal{L}}_\lambda(\theta, M^*\lambda) = 0$ , which shows that the first term in the LHS of the second line of Equation 8 is nonnegative. Hence:

$$-\mathbb{E}_{t \sim [T]} \left[ \hat{\mathcal{L}}_\lambda \left( \theta^{(t)}, \lambda^{(t)} \right) \right] \leq \epsilon_\lambda$$

so by the definitions of  $\hat{\mathcal{L}}_\lambda$  (Definition 1) and  $G^{(\text{val})}(\hat{\Theta})$  (Definition 2), and the facts that  $\tilde{\ell}_i \geq \ell_i$  and  $\lambda^{(t)} \in \Delta^m$ :

$$\mathbb{E}_{t \sim [T], x \sim \mathcal{D}} \left[ \sum_{i=1}^m \lambda_{i+1}^{(t)} \tilde{\ell}_i \left( x; \theta^{(t)} \right) \right] \geq -\epsilon_\lambda - G^{(\text{val})}(\hat{\Theta})$$

Notice that  $\hat{\mathcal{L}}_\theta$  is linear in  $\lambda$ , so the first line of Equation 8, combined with the definition of  $\tilde{G}^{(\text{trn})}(\Theta)$ , becomes:

$$\begin{aligned} & \mathbb{E}_{t \sim [T], x \sim \mathcal{D}} \left[ \lambda_1^{(t)} \ell_0 \left( x; \theta^{(t)} \right) + \sum_{i=1}^m \lambda_{i+1}^{(t)} \tilde{\ell}_i \left( x; \theta^{(t)} \right) \right] - \inf_{\theta^* \in \Theta} \mathbb{E}_{x \sim \mathcal{D}} \left[ \bar{\lambda}_1 \ell_0 \left( x; \theta^* \right) + \sum_{i=1}^m \bar{\lambda}_{i+1} \tilde{\ell}_i \left( x; \theta^* \right) \right] \\ & \leq \epsilon_\theta + 2\tilde{G}^{(\text{trn})}(\Theta) \end{aligned}$$

Combining the above two results:

$$\begin{aligned} & \mathbb{E}_{t \sim [T], x \sim \mathcal{D}} \left[ \lambda_1^{(t)} \ell_0 \left( x; \theta^{(t)} \right) \right] - \inf_{\theta^* \in \Theta} \mathbb{E}_{x \sim \mathcal{D}} \left[ \bar{\lambda}_1 \ell_0 \left( x; \theta^* \right) + \sum_{i=1}^m \bar{\lambda}_{i+1} \tilde{\ell}_i \left( x; \theta^* \right) \right] \\ & \leq \epsilon_\theta + \epsilon_\lambda + 2\tilde{G}^{(\text{trn})}(\Theta) + G^{(\text{val})}(\hat{\Theta}) \end{aligned} \quad (12)$$

Choose  $\theta^*$  to be the optimal solution that satisfies the proxy constraints, so that  $\mathbb{E}_{x \sim \mathcal{D}} \left[ \tilde{\ell}_i \left( x; \theta^* \right) \right] \leq 0$  for all  $i \in [m]$ . Then:

$$\mathbb{E}_{t \sim [T], x \sim \mathcal{D}} \left[ \lambda_1^{(t)} \ell_0 \left( x; \theta^{(t)} \right) \right] - \mathbb{E}_{x \sim \mathcal{D}} \left[ \bar{\lambda}_1 \ell_0 \left( x; \theta^* \right) \right] \leq \epsilon_\theta + \epsilon_\lambda + 2\tilde{G}^{(\text{trn})}(\Theta) + G^{(\text{val})}(\hat{\Theta})$$

which is the optimality claim.

**Feasibility:** We'll begin by simplifying our notation: define  $g_1(\theta) := 0$  and  $g_{i+1}(\theta) := \mathbb{E}_{x \sim \mathcal{S}^{(\text{val})}} [\ell_i(x; \theta)]$  for  $i \in [m]$ , so that  $\hat{\mathcal{L}}_\lambda(\theta, \lambda) = \langle \lambda, g \cdot (\theta) \rangle$ . Consider the first term in the LHS of the second line of Equation 8:

$$\begin{aligned} \max_{M^* \in \mathcal{M}} \mathbb{E}_{t \sim [T]} \left[ \hat{\mathcal{L}}_\lambda \left( \theta^{(t)}, M^* \lambda^{(t)} \right) \right] &= \max_{M^* \in \mathcal{M}} \mathbb{E}_{t \sim [T]} \left[ \left\langle M^* \lambda^{(t)}, g \cdot \left( \theta^{(t)} \right) \right\rangle \right] \\ &= \max_{M^* \in \mathcal{M}} \mathbb{E}_{t \sim [T]} \left[ \sum_{i=1}^{m+1} \sum_{j=1}^{m+1} M_{j,i}^* \lambda_i^{(t)} g_j \left( \theta^{(t)} \right) \right] \\ &= \sum_{i=1}^{m+1} \max_{M^*, i \in \Delta^m} \sum_{j=1}^{m+1} \mathbb{E}_{t \sim [T]} \left[ M_{j,i}^* \lambda_i^{(t)} g_j \left( \theta^{(t)} \right) \right] \\ &= \sum_{i=1}^{m+1} \max_{j \in [m+1]} \mathbb{E}_{t \sim [T]} \left[ \lambda_i^{(t)} g_j \left( \theta^{(t)} \right) \right] \end{aligned}$$

where we used the fact that, since  $M^*$  is left-stochastic, each of its columns is a  $(m+1)$ -dimensional multinoulli distribution. For the second term in the LHS of the second line of Equation 8, we can use the fact that  $g_1(\theta) = 0$ :

$$\mathbb{E}_{t \sim [T]} \left[ \sum_{i=2}^{m+1} \lambda_i^{(t)} g_i \left( \theta^{(t)} \right) \right] \leq \sum_{i=2}^{m+1} \max_{j \in [m+1]} \mathbb{E}_{t \sim [T]} \left[ \lambda_i^{(t)} g_j \left( \theta^{(t)} \right) \right]$$

Plugging these two results into the second line of Equation 8, the two sums collapse, leaving:

$$\max_{i \in [m+1]} \mathbb{E}_{t \sim [T]} \left[ \lambda_1^{(t)} g_i \left( \theta^{(t)} \right) \right] \leq \epsilon_\lambda$$



Substituting the definitions of  $g_i$  and  $G^{(\text{val})}(\hat{\Theta})$  (since the  $g_i$ s are defined on  $S^{(\text{val})}$ , but we want our result to hold on  $\mathcal{D}$ ) then yields the feasibility claim.

**Bound on  $\bar{\lambda}_1$ :** Choosing  $\theta^* = \theta'$  in Equation 12 (recall that  $\theta'$  satisfies all of the proxy constraints with margin  $\gamma$ ) and substituting the definition of  $B_{\ell_0}$ :

$$\begin{aligned} \epsilon_\theta + \epsilon_\lambda + 2\tilde{G}^{(\text{trn})}(\Theta) + G^{(\text{val})}(\hat{\Theta}) &\geq \mathbb{E}_{t \sim [T], x \sim \mathcal{D}} \left[ \lambda_1^{(t)} \ell_0(x; \theta^{(t)}) - \lambda_1^{(t)} \ell_0(x; \theta') \right] + (1 - \bar{\lambda}_1) \gamma \\ &\geq -\bar{\lambda}_1 B_{\ell_0} + (1 - \bar{\lambda}_1) \gamma \end{aligned}$$

Solving for  $\bar{\lambda}_1$  yields the claim.  $\square$

Before moving on to the convergence and generalization properties of our actual algorithms, we need to state some (fairly standard) elementary results:

**Definition 5.** We say that  $C_r \subseteq \mathbb{R}^{m+1}$  is a radius- $r$  external covering of  $\Lambda := \Delta^m \subseteq \mathbb{R}_+^{m+1}$  w.r.t. the 1-norm if for every  $\lambda \in \Lambda$  there exists a  $\tilde{\lambda} \in C_r$  for which  $\|\lambda - \tilde{\lambda}\|_1 \leq r$ . Notice that we do not require  $C_r$  to be a subset of  $\Lambda$ —this is why it’s an external covering.

**Lemma 1.** Assuming that  $r \leq 1$ , there exists a radius- $r$  external covering of  $\Lambda := \Delta^m \subseteq \mathbb{R}_+^{m+1}$  w.r.t. the 1-norm of size no larger than  $(5/r)^m$ .

*Proof.* Consider the  $m$ -dimensional unit ball  $\tilde{B} := \{\tilde{\lambda} \in \mathbb{R}^m : \|\tilde{\lambda}\|_1 \leq 1\}$  w.r.t. the 1-norm (note that we could instead consider only the positive orthant, which would improve the constant in the overall result). There exists a radius- $r/2$  covering  $\tilde{C}_r \subseteq \mathbb{R}^m$  of  $\tilde{B}$  with  $|\tilde{C}_r| \leq (1 + 4/r)^m \leq (5/r)^m$  (Bartlett, 2013).

Define  $C_r \subseteq \mathbb{R}^{m+1}$  as:

$$C_r = \left\{ \left[ \begin{array}{c} \tilde{\lambda} \\ 1 - \|\tilde{\lambda}\|_1 \end{array} \right] : \tilde{\lambda} \in \tilde{C}_r \right\}$$

Notice that we do not necessarily have that  $C_r \subseteq \Delta^m$ , i.e. this will be an *external* covering.

From any  $\lambda \in \Delta^m$ , we can define  $\lambda' \in \mathbb{R}_+^m$  by dropping the last coordinate of  $\lambda$ , and we’ll have that  $\|\lambda'\|_1 \leq \|\lambda\|_1 = 1$ , so there will exist a  $\tilde{\lambda} \in \tilde{C}_r$  such that  $\|\tilde{\lambda} - \lambda'\|_1 \leq r/2$ , which implies that the corresponding element of  $C_r$  is  $r$ -far from  $\lambda$ , showing that  $C_r$  is a radius- $r$  covering of  $\Delta^m$  w.r.t. the 1-norm.  $\square$

**Lemma 2.** Let  $S$  be an i.i.d. sample from a distribution  $\mathcal{D}$  supported on  $\mathcal{X}$ , and  $\hat{\Theta} \subseteq \Theta$  the finite set of permitted model parameters, which defines a finite function class ( $\hat{\Theta}$  may be a random variable, but must be independent of  $S$ ). Suppose that  $\ell : \mathcal{X} \times \Theta \rightarrow [a, b]$  with  $B_\ell := b - a$ . Then:

$$\left| \frac{1}{|S|} \sum_{x \in S} \ell(x; \theta) - \mathbb{E}_{x \sim \mathcal{D}} [\ell(x; \theta)] \right| < B_\ell \sqrt{\frac{\ln(2|\hat{\Theta}|/\delta)}{2|S|}}$$

for all  $\theta \in \hat{\Theta}$ , with probability at least  $1 - \delta$  over the sampling of  $S$ .

*Proof.* Allowing  $\hat{\Theta}$  to be a random variable independent of  $S$ , instead of a constant set, doesn’t significantly change the standard proof (e.g. Srebro, 2016). By Hoeffding’s inequality:

$$\Pr \left\{ \left| \frac{1}{|S|} \sum_{x \in S} \ell(x; \theta) - \mathbb{E}_{x \sim \mathcal{D}} [\ell(x; \theta)] \right| \geq \epsilon \right\} \leq 2 \exp \left( -\frac{2|S| \epsilon^2}{B_\ell^2} \right)$$

the above holding for any  $\theta \in \Theta$ . Since  $\hat{\Theta}$  is independent of  $S$ , we can apply the union bound over all  $\theta \in \hat{\Theta}$ :

$$\Pr \left\{ \exists \theta \in \hat{\Theta}. \left| \frac{1}{|S|} \sum_{x \in S} \ell(x; \theta) - \mathbb{E}_{x \sim \mathcal{D}} [\ell(x; \theta)] \right| \geq \epsilon \right\} \leq 2|\hat{\Theta}| \exp \left( -\frac{2|S| \epsilon^2}{B_\ell^2} \right)$$

Rearranging terms yields the claimed result.  $\square$

**E.1. Algorithm 1**

**Lemma 3.** *If we take  $\eta_\lambda := \sqrt{(m+1) \ln(m+1) / TB_\Delta^2}$ , then the result of Algorithm 1 satisfies the conditions of Theorem 3 with:*

$$\begin{aligned}\epsilon_\theta &= \rho + 2rB_{\tilde{\ell}} \\ \epsilon_\lambda &= 2B_\Delta \sqrt{\frac{(m+1) \ln(m+1)}{T}}\end{aligned}$$

where  $B_{\tilde{\ell}} \geq |\ell(x, \theta)|$  for all  $\ell \in \{\ell_0, \tilde{\ell}_1, \dots, \tilde{\ell}_m\}$ , and  $B_\Delta \geq \max_{t \in [T]} \|\Delta_\lambda^{(t)}\|_\infty$  is a bound on the gradients.

*Proof.* Since  $C_r$  is a radius- $r$  external covering of  $\Lambda := \Delta^m \subseteq \mathbb{R}_+^{m+1}$  w.r.t. the 1-norm, we must have that  $\|\tilde{\lambda}^{(t)} - \lambda^{(t)}\|_1 \leq r$ , which implies by Definition 1 that:

$$\left| \hat{\mathcal{L}}_\theta(\theta, \tilde{\lambda}^{(t)}) - \hat{\mathcal{L}}_\lambda(\theta, \lambda^{(t)}) \right| \leq rB_{\tilde{\ell}}$$

the above holding for all  $\theta \in \Theta$ , and all  $t$ . In particular, this implies that:

$$\begin{aligned}\hat{\mathcal{L}}_\theta(\theta^{(t)}, \lambda^{(t)}) &\leq \hat{\mathcal{L}}_\theta(\theta^{(t)}, \tilde{\lambda}^{(t)}) + rB_{\tilde{\ell}} \\ &\leq \inf_{\theta^* \in \Theta} \hat{\mathcal{L}}_\theta(\theta^*, \tilde{\lambda}^{(t)}) + \rho + rB_{\tilde{\ell}} \\ &\leq \inf_{\theta^* \in \Theta} \hat{\mathcal{L}}_\theta(\theta^*, \lambda^{(t)}) + \rho + 2rB_{\tilde{\ell}}\end{aligned}$$

Therefore:

$$\begin{aligned}\frac{1}{T} \sum_{t=1}^T \hat{\mathcal{L}}_\theta(\theta^{(t)}, \lambda^{(t)}) &\leq \frac{1}{T} \sum_{t=1}^T \inf_{\theta^* \in \Theta} \hat{\mathcal{L}}_\theta(\theta^*, \lambda^{(t)}) + \rho + 2rB_{\tilde{\ell}} \\ &\leq \inf_{\theta^* \in \Theta} \frac{1}{T} \sum_{t=1}^T \hat{\mathcal{L}}_\theta(\theta^*, \lambda^{(t)}) + \rho + 2rB_{\tilde{\ell}}\end{aligned}$$

so the first condition of Theorem 3 is satisfied with the claimed  $\epsilon_\theta$ .

The second condition, on  $\epsilon_\lambda$ , follows immediately from Lemma 8 of Appendix C.1 of Cotter et al. (2019), taking  $\tilde{m} = m + 1$ .  $\square$

**Lemma 4.** *If we take  $\hat{\Theta} := \{\theta^{(1)}, \dots, \theta^{(T)}\}$  as in Theorem 3, where  $\theta^{(1)}, \dots, \theta^{(T)}$  are the result of Algorithm 1, then with probability  $1 - \delta$  over the sampling of  $S^{(\text{val})}$ :*

$$G^{(\text{val})}(\hat{\Theta}) < B_\ell \sqrt{\frac{\ln(2m|C_r|/\delta)}{2|S^{(\text{val})}|}}$$

where  $B_\ell \geq \max_{i \in [m]} (b_i - a_i)$  assuming that the range of each  $\ell_i$  is the interval  $[a_i, b_i]$ .

*Proof.* Since each  $\theta^{(t)}$  is uniquely associated with a  $\tilde{\lambda}^{(t)} \in C_r$  (Definition 3), we will have that  $\hat{\Theta} \subseteq \Theta_{C_r}$ , where:

$$\Theta_{C_r} := \left\{ \mathcal{O}_\rho \left( \hat{\mathcal{L}}_\theta(\cdot, \tilde{\lambda}) \right) : \tilde{\lambda} \in C_r \right\}$$

Because the oracle call defining  $\Theta_{C_r}$  depends only on  $\hat{\mathcal{L}}_\theta$ , which itself depends only on  $S^{(\text{trn})}$ , we can apply Lemma 2 to  $\Theta_{C_r}$ , yielding that, for each  $i \in [m]$ , the following holds with probability  $\delta/m$  for all  $\theta \in \Theta_{C_r}$ :

$$\left| \frac{1}{|S^{(\text{val})}|} \sum_{x \in S^{(\text{val})}} \ell_i(x; \theta) - \mathbb{E}_{x \sim \mathcal{D}} [\ell_i(x; \theta)] \right| < B_\ell \sqrt{\frac{\ln(2m|\Theta_{C_r}|/\delta)}{2|S^{(\text{val})}|}}$$

The claimed result on  $G^{(\text{val})}(\hat{\Theta})$  then follows from the union bound and the facts that  $\hat{\Theta} \subseteq \Theta_{C_r}$  and  $|\Theta_{C_r}| = |C_r|$ .  $\square$

**Theorem 1.** Given any  $\epsilon > 0$ , there exists a covering  $C_r$  such that, if we take  $T \geq 4B_\Delta^2 (m+1) \ln(m+1) / \epsilon^2$  and  $\eta_\lambda = \sqrt{(m+1) \ln(m+1) / TB_\Delta^2}$ , where  $B_\Delta \geq \max_{t \in [T]} \left\| \Delta_\lambda^{(t)} \right\|_\infty$  is a bound on the gradients, then the following hold, where  $\hat{\Theta} := \{\theta^{(1)}, \dots, \theta^{(T)}\}$  is the set of results of Algorithm 1.

**Optimality and Feasibility:** Let  $\bar{\theta}$  be a random variable taking values from  $\hat{\Theta}$ , defined such that  $\bar{\theta} = \theta^{(t)}$  with probability  $\lambda_1^{(t)} / \sum_{s=1}^T \lambda_1^{(s)}$ , and let  $\bar{\lambda} := (\sum_{t=1}^T \lambda^{(t)}) / T$ . Then  $\bar{\theta}$  is nearly-optimal in expectation:

$$\mathbb{E}_{\bar{\theta}, x \sim \mathcal{D}} [\ell_0(x; \bar{\theta})] \leq \mathbb{E}_{x \sim \mathcal{D}} [\ell_0(x; \theta^*)] + \frac{1}{\lambda_1} \left( \rho + 2\epsilon + 2\tilde{G}^{(\text{trn})}(\Theta) + G^{(\text{val})}(\hat{\Theta}) \right) \quad (3)$$

where  $\theta^*$  minimizes  $\mathbb{E}_{x \sim \mathcal{D}} [\ell_0(x; \cdot)]$  subject to the proxy-constraints  $\mathbb{E}_{x \sim \mathcal{D}} [\tilde{\ell}_i(x; \theta^*)] \leq 0$ . It is also nearly-feasible:

$$\max_{i \in [m]} \mathbb{E}_{\bar{\theta}, x \sim \mathcal{D}} [\ell_i(x; \bar{\theta})] \leq \frac{\epsilon}{\lambda_1} + G^{(\text{val})}(\hat{\Theta}) \quad (4)$$

Additionally, if there exists a  $\theta' \in \Theta$  that satisfies all of the constraints with margin  $\gamma$  (i.e.  $\mathbb{E}_{x \sim \mathcal{D}} [\ell_i(x; \theta')] \leq -\gamma$  for all  $i \in [m]$ ), then:

$$\bar{\lambda}_1 \geq \frac{1}{\gamma + B_{\ell_0}} \left( \gamma - \rho - 2\epsilon - 2\tilde{G}^{(\text{trn})}(\Theta) - G^{(\text{val})}(\hat{\Theta}) \right) \quad (5)$$

where  $B_{\ell_0} \geq \sup_{\theta \in \Theta} \mathbb{E}_{x \sim \mathcal{D}} [\ell_0(x; \theta)] - \inf_{\theta \in \Theta} \mathbb{E}_{x \sim \mathcal{D}} [\ell_0(x; \theta)]$  is a bound on the range of the objective loss.

**Generalization:** With probability  $1 - \delta$  over the sampling of  $S^{(\text{val})}$ :

$$G^{(\text{val})}(\hat{\Theta}) < B_\ell \sqrt{\frac{m \ln(10B_{\tilde{\ell}}/\epsilon) + \ln(2m/\delta)}{2|S^{(\text{val})}|}} \quad (6)$$

where  $B_{\tilde{\ell}} \geq |\ell(x, \theta)|$  for all  $\ell \in \{\ell_0, \tilde{\ell}_1, \dots, \tilde{\ell}_m\}$ , and  $B_\ell \geq \max_{i \in [m]} (b_i - a_i)$  assuming that the range of each  $\ell_i$  is the interval  $[a_i, b_i]$ .

*Proof.* The particular values we choose for  $T$  and  $\eta_\lambda$  come from Lemma 3, taking  $r = \epsilon/2B_{\tilde{\ell}}$ ,  $\epsilon_\theta = \rho + 2rB_{\tilde{\ell}} = \rho + \epsilon$ , and  $\epsilon_\lambda = \epsilon$ . The optimality and feasibility results then follow from Theorem 3.

For the bound on  $G^{(\text{val})}(\hat{\Theta})$ , notice that by Lemma 1, there exists a radius- $r$  covering  $C_r$  w.r.t. the 1-norm with  $|C_r| \leq (5/r)^m = (10B_{\tilde{\ell}}/\epsilon)^m$ . Substituting this, and the definition of  $r$ , into the bound of Lemma 4 yields the claimed bound.  $\square$

## E.2. Algorithm 2

**Lemma 5.** Suppose that  $\Theta$  is compact and convex, and that  $\ell(x; \theta)$  is  $\mu$ -strongly convex in  $\theta$  for all  $\ell \in \{\ell_0, \tilde{\ell}_1, \dots, \tilde{\ell}_m\}$ . If we take  $\eta_\lambda := \sqrt{(m+1) \ln(m+1) / T_\lambda B_\Delta^2}$ , then the result of Algorithm 2 satisfies the conditions of Theorem 3 with:

$$\begin{aligned} \epsilon_\theta &= \frac{B_\Delta^2 (1 + \ln T_\theta)}{2\mu T_\theta} \\ \epsilon_\lambda &= 2B_\Delta \sqrt{\frac{(m+1) \ln(m+1)}{T_\lambda}} \end{aligned}$$

where  $B_{\tilde{\Delta}} \geq \max_{s, t \in [T_\theta] \times [T_\lambda]} \left\| \tilde{\Delta}_\theta^{(t, s)} \right\|_2$  is a bound on the subgradients, and  $B_\Delta \geq \max_{t \in [T_\lambda]} \left\| \Delta_\lambda^{(t)} \right\|_\infty$  is a bound on the gradients.

*Proof.* By Lemma 1 of Shalev-Shwartz et al. (2011), the fact that  $\hat{\mathcal{L}}_\theta(\theta, \lambda)$  is  $\mu$ -strongly convex in  $\theta$  (because  $\ell(x; \theta)$  is  $\mu$ -strongly convex in  $\theta$  for  $\ell \in \{\ell_0, \tilde{\ell}_1, \dots, \tilde{\ell}_m\}$ ), and  $\lambda \in \Lambda := \Delta^m \subseteq \mathbb{R}_+^{m+1}$ , and Jensen's inequality:

$$\hat{\mathcal{L}}_\theta(\theta^{(t)}, \lambda^{(t)}) \leq \frac{1}{T_\theta} \sum_{s=1}^{T_\theta} \hat{\mathcal{L}}_\theta(\tilde{\theta}^{(t, s)}, \lambda^{(t)}) \leq \min_{\theta^* \in \Theta} \hat{\mathcal{L}}_\theta(\theta^*, \lambda^{(t)}) + \frac{B_\Delta^2 (1 + \ln T_\theta)}{2\mu T_\theta} \quad (13)$$

the above holding for all  $t$ . Therefore:

$$\frac{1}{T} \sum_{t=1}^T \hat{\mathcal{L}}_{\theta} \left( \theta^{(t)}, \lambda^{(t)} \right) \leq \min_{\theta^* \in \Theta} \frac{1}{T} \sum_{t=1}^T \hat{\mathcal{L}}_{\theta} \left( \theta^*, \lambda^{(t)} \right) + \frac{B_{\Delta}^2 (1 + \ln T_{\theta})}{2\mu T_{\theta}}$$

so the first condition of Theorem 3 is satisfied with the claimed  $\epsilon_{\theta}$ .

As in the proof of Lemma 3, the second condition, on  $\epsilon_{\lambda}$ , follows immediately from Lemma 8 of Appendix C.1 of Cotter et al. (2019), taking  $\tilde{m} = m + 1$ .  $\square$

**Lemma 6.** *In addition to the conditions of Lemma 5, suppose that  $\ell(x; \theta)$  is  $L$ -Lipschitz continuous in  $\theta$  for all  $\ell \in \{\ell_1, \dots, \ell_m\}$ . If we take  $\hat{\Theta} := \{\theta^{(1)}, \dots, \theta^{(T)}\}$  as in Theorem 3, where  $\theta^{(1)}, \dots, \theta^{(T)}$  are the result of Algorithm 2, then with probability  $1 - \delta$  over the sampling of  $S^{(\text{val})}$ :*

$$G^{(\text{val})}(\hat{\Theta}) < 2L \sqrt{\frac{4rB_{\tilde{\ell}}}{\mu}} + \frac{2LB_{\Delta}}{\mu} \sqrt{\frac{1 + \ln T_{\theta}}{T_{\theta}}} + B_{\ell} \sqrt{\frac{\ln(2m|C_r|/\delta)}{2|S^{(\text{val})}|}}$$

where  $B_{\tilde{\ell}} \geq |\ell(x, \theta)|$  for all  $\ell \in \{\ell_0, \tilde{\ell}_1, \dots, \tilde{\ell}_m\}$ ,  $B_{\ell} \geq \max_{i \in [m]} (b_i - a_i)$  assuming that the range of each  $\ell_i$  is the interval  $[a_i, b_i]$ , and  $C_r$  is a radius- $r$  covering of  $\Lambda := \Delta^m \subseteq \mathbb{R}_+^{m+1}$  w.r.t. the 1-norm.

*Proof.* Define  $\tilde{\lambda}^{(t)} := \operatorname{argmin}_{\tilde{\lambda} \in C_r} \|\lambda^{(t)} - \tilde{\lambda}\|$  for all  $t$ . Since  $C_r$  is a radius- $r$  covering of  $\Lambda := \Delta^m$  w.r.t. the 1-norm, we must have that  $\|\tilde{\lambda}^{(t)} - \lambda^{(t)}\|_1 \leq r$ , which implies by Definition 1 that:

$$\left| \hat{\mathcal{L}}_{\theta} \left( \theta, \tilde{\lambda}^{(t)} \right) - \hat{\mathcal{L}}_{\lambda} \left( \theta, \lambda^{(t)} \right) \right| \leq rB_{\tilde{\ell}}$$

the above holding for all  $\theta \in \Theta$ , and all  $t$ . Take  $\theta^{(t^*)} := \operatorname{argmin}_{\theta^* \in \Theta} \hat{\mathcal{L}}_{\theta} \left( \theta^*, \lambda^{(t)} \right)$  and  $\tilde{\theta}^{(t^*)} := \operatorname{argmin}_{\tilde{\theta}^* \in \Theta} \hat{\mathcal{L}}_{\theta} \left( \tilde{\theta}^*, \tilde{\lambda}^{(t)} \right)$ . Then, by the above result and the triangle inequality:

$$\begin{aligned} \left| \hat{\mathcal{L}}_{\theta} \left( \theta^{(t^*)}, \lambda^{(t)} \right) - \hat{\mathcal{L}}_{\theta} \left( \tilde{\theta}^{(t^*)}, \tilde{\lambda}^{(t)} \right) \right| &\leq rB_{\tilde{\ell}} \\ \left| \hat{\mathcal{L}}_{\theta} \left( \theta^{(t^*)}, \tilde{\lambda}^{(t)} \right) - \hat{\mathcal{L}}_{\theta} \left( \tilde{\theta}^{(t^*)}, \tilde{\lambda}^{(t)} \right) \right| &\leq 2rB_{\tilde{\ell}} \end{aligned}$$

so by the fact that  $\hat{\mathcal{L}}_{\theta}(\theta, \lambda)$  is  $\mu$ -strongly convex in  $\theta$  for all  $\lambda$ :

$$\left\| \theta^{(t^*)} - \tilde{\theta}^{(t^*)} \right\|_2 \leq \sqrt{\frac{4rB_{\tilde{\ell}}}{\mu}}$$

Again by strong convexity, but applied to Equation 13 in the proof of Lemma 5:

$$\left\| \theta^{(t)} - \theta^{(t^*)} \right\|_2 \leq \sqrt{\frac{B_{\Delta}^2 (1 + \ln T_{\theta})}{\mu^2 T_{\theta}}}$$

so by the triangle inequality:

$$\left\| \theta^{(t)} - \tilde{\theta}^{(t^*)} \right\|_2 \leq \sqrt{\frac{4rB_{\tilde{\ell}}}{\mu}} + \sqrt{\frac{B_{\Delta}^2 (1 + \ln T_{\theta})}{\mu^2 T_{\theta}}}$$

and by  $L$ -Lipschitz continuity:

$$\left| \ell_i \left( x; \theta^{(t)} \right) - \ell_i \left( x; \tilde{\theta}^{(t^*)} \right) \right| \leq L \left( \sqrt{\frac{4rB_{\tilde{\ell}}}{\mu}} + \sqrt{\frac{B_{\Delta}^2 (1 + \ln T_{\theta})}{\mu^2 T_{\theta}}} \right) \quad (14)$$

the above holding for all  $t$ , and all  $i \in [m]$ .

Define  $\tilde{\Theta} := \{\tilde{\theta}^{(1*)}, \dots, \tilde{\theta}^{(T*)}\}$ . Observe that since  $\tilde{\theta}^{(t*)}$  is uniquely associated with a  $\tilde{\lambda}^{(t)} \in C_r$ , we will have that  $\tilde{\Theta} \subseteq \Theta_{C_r}$ , where:

$$\Theta_{C_r} := \left\{ \underset{\tilde{\theta}^* \in \Theta}{\operatorname{argmin}} \hat{\mathcal{L}}_{\theta}(\tilde{\theta}^*, \tilde{\lambda}) : \tilde{\lambda} \in C_r \right\}$$

Because the argmins defining  $\Theta_{C_r}$  depend only on  $\hat{\mathcal{L}}_{\theta}$ , which itself depends only on  $S^{(\text{trn})}$ , we can apply Lemma 2 to  $\Theta_{C_r}$ , yielding that, for each  $i \in [m]$ , the following holds with probability  $\delta/m$  for any  $\theta \in \Theta_{C_r}$ :

$$\left| \frac{1}{|S^{(\text{val})}|} \sum_{x \in S^{(\text{val})}} \ell_i(x; \theta) - \mathbb{E}_{x \sim \mathcal{D}}[\ell_i(x; \theta)] \right| < B_{\ell} \sqrt{\frac{\ln(2m |\Theta_{C_r}| / \delta)}{2 |S^{(\text{val})}|}}$$

By the union bound, we could instead take the above to hold uniformly for all  $i \in [m]$  with probability  $1 - \delta$ . Substituting Equation 14 and using the facts that  $\tilde{\Theta} \subseteq \Theta_{C_r}$  and  $|\Theta_{C_r}| = |C_r|$  yields the claimed result.  $\square$

**Theorem 2.** *Suppose that  $\Theta$  is compact and convex, and that  $\ell(x; \theta)$  is  $\mu$ -strongly convex in  $\theta$  for all  $\ell \in \{\ell_0, \tilde{\ell}_1, \dots, \tilde{\ell}_m\}$ . Given any  $\epsilon > 0$ , if we take  $T_{\theta} \geq (B_{\Delta}^2 / \mu \epsilon) \ln(B_{\Delta}^2 / \mu \epsilon)$ ,  $T_{\lambda} \geq 4B_{\Delta}^2 (m+1) \ln(m+1) / \epsilon^2$  and  $\eta_{\lambda} = \sqrt{(m+1) \ln(m+1) / T_{\lambda} B_{\Delta}^2}$ , where  $B_{\Delta}$  is as in Theorem 1 and  $B_{\Delta} \geq \max_{s, t \in [T_{\theta}] \times [T_{\lambda}]} \|\tilde{\Delta}_{\theta}^{(t, s)}\|_2$  is a bound on the subgradients, then the following hold, where  $\hat{\Theta} := \{\theta^{(1)}, \dots, \theta^{(T_{\lambda})}\}$  is the set of results of Algorithm 1.*

**Optimality and Feasibility:** *Let  $\bar{\theta}$  be a random variable taking values from  $\hat{\Theta}$ , defined such that  $\bar{\theta} = \theta^{(t)}$  with probability  $\lambda_1^{(t)} / \sum_{s=1}^T \lambda_1^{(s)}$ , and let  $\bar{\lambda} := (\sum_{t=1}^T \lambda^{(t)}) / T_{\lambda}$ . Then  $\bar{\theta}$  is nearly-optimal in expectation:*

$$\mathbb{E}_{\bar{\theta}, x \sim \mathcal{D}}[\ell_0(x; \bar{\theta})] \leq \mathbb{E}_{x \sim \mathcal{D}}[\ell_0(x; \theta^*)] + \frac{1}{\lambda_1} \left( 2\epsilon + 2\tilde{G}^{(\text{trn})}(\Theta) + G^{(\text{val})}(\hat{\Theta}) \right)$$

where  $\theta^*$  minimizes  $\mathbb{E}_{x \sim \mathcal{D}}[\ell_0(x; \cdot)]$  subject to the proxy-constraints  $\mathbb{E}_{x \sim \mathcal{D}}[\tilde{\ell}_i(x; \theta^*)] \leq 0$ . It is also nearly-feasible:

$$\max_{i \in [m]} \mathbb{E}_{\bar{\theta}, x \sim \mathcal{D}}[\ell_i(x; \bar{\theta})] \leq \frac{\epsilon}{\lambda_1} + G^{(\text{val})}(\hat{\Theta})$$

Additionally, if there exists a  $\theta' \in \Theta$  that satisfies all of the constraints with margin  $\gamma$  (i.e.  $\mathbb{E}_{x \sim \mathcal{D}}[\ell_i(x; \theta')] \leq -\gamma$  for all  $i \in [m]$ ), then:

$$\bar{\lambda}_1 \geq \frac{1}{\gamma + B_{\ell_0}} \left( \gamma - 2\epsilon - 2\tilde{G}^{(\text{trn})}(\Theta) - G^{(\text{val})}(\hat{\Theta}) \right)$$

where  $B_{\ell_0}$  is as in Theorem 1.

**Generalization:** *If, in addition to the above requirements,  $\ell(x; \theta)$  is  $L$ -Lipschitz continuous in  $\theta$  for all  $\ell \in \{\ell_1, \dots, \ell_m\}$ , then with probability  $1 - \delta$  over the sampling of  $S^{(\text{val})}$ :*

$$G^{(\text{val})}(\hat{\Theta}) < B_{\ell} \sqrt{\frac{2m}{|S^{(\text{val})}|} \max \left\{ 1, \ln \left( \frac{160L^2 B_{\tilde{\ell}} |S^{(\text{val})}|}{m \mu B_{\tilde{\ell}}^2} \right) \right\}} + B_{\ell} \sqrt{\frac{\ln(2m/\delta)}{2 |S^{(\text{val})}|}} + 2L\epsilon \sqrt{\frac{2}{\mu}} \quad (7)$$

where  $B_{\tilde{\ell}}$  and  $B_{\ell}$  are as in Theorem 1.

*Proof.* The particular values we choose for  $T_{\theta}$ ,  $T_{\lambda}$  and  $\eta_{\lambda}$  come from Lemma 5, taking  $\epsilon_{\theta} = \epsilon_{\lambda} = \epsilon$ . The optimality and feasibility results then follow from Theorem 3.

For the bound on  $G^{(\text{val})}(\hat{\Theta})$ , notice that by Lemma 1, there exists a radius- $r$  external covering  $C_r$  w.r.t. the 1-norm with  $|C_r| \leq \max\{1, (5/r)^m\}$ . Substituting into the bound of Lemma 6:

$$\begin{aligned} G^{(\text{val})}(\hat{\Theta}) &< 2L \sqrt{\frac{4r B_{\tilde{\ell}}}{\mu}} + \frac{2LB_{\Delta}}{\mu} \sqrt{\frac{1 + \ln T_{\theta}}{T_{\theta}}} + B_{\ell} \sqrt{\frac{m \max\{0, \ln(5/r)\} + \ln(2m/\delta)}{2 |S^{(\text{val})}|}} \\ &< 2L \sqrt{\frac{4r B_{\tilde{\ell}}}{\mu}} + B_{\ell} \sqrt{\frac{m \max\{0, \ln(5/r)\}}{2 |S^{(\text{val})}|}} + \frac{2LB_{\Delta}}{\mu} \sqrt{\frac{1 + \ln T_{\theta}}{T_{\theta}}} + B_{\ell} \sqrt{\frac{\ln(2m/\delta)}{2 |S^{(\text{val})}|}} \end{aligned}$$

Taking  $r := (m\mu B_\ell^2) / (32L^2 B_{\tilde{\ell}} |S^{(\text{val})}|)$ :

$$\begin{aligned} G^{(\text{val})}(\hat{\Theta}) &< 2B_\ell \sqrt{\frac{m \max\{1, \ln(5/r)\}}{2|S^{(\text{val})}|}} + \frac{2LB_{\tilde{\Delta}}}{\mu} \sqrt{\frac{1 + \ln T_\theta}{T_\theta}} + B_\ell \sqrt{\frac{\ln(2m/\delta)}{2|S^{(\text{val})}|}} \\ &< B_\ell \sqrt{\frac{2m}{|S^{(\text{val})}|} \max\left\{1, \ln\left(\frac{160L^2 B_{\tilde{\ell}} |S^{(\text{val})}|}{m\mu B_\ell^2}\right)\right\}} \\ &\quad + B_\ell \sqrt{\frac{\ln(2m/\delta)}{2|S^{(\text{val})}|}} + \frac{2LB_{\tilde{\Delta}}}{\mu} \sqrt{\frac{1 + \ln T_\theta}{T_\theta}} \end{aligned}$$

Substituting the definition of  $T_\theta$  then yields the claimed result.  $\square$

## F. One-dataset Lagrangian Baseline Approach

In this appendix, we'll analyze the most natural theoretical baseline for our proposed approach, namely using a single training dataset, and optimizing the empirical Lagrangian:

$$\hat{\mathcal{L}}(\theta, \lambda) := \frac{1}{|S^{(\text{trn})}|} \sum_{x \in S^{(\text{trn})}} \left( \ell_0(x; \theta) + \sum_{i=1}^m \lambda_i \ell_i(x; \theta) \right) \quad (15)$$

This is essentially a minor extension of the approach proposed by Agarwal et al. (2018)—who proposed using the Lagrangian formulation in the particular case of fair classification—to the slightly more general setting of inequality constrained optimization.

Theorem 1 of Cotter et al. (2019) characterizes the optimality and feasibility properties of Nash equilibria of the in-expectation Lagrangian (Equation 2). The analogous result for the empirical Lagrangian of Equation 15 is given in the following theorem:

**Theorem 4.** Define  $\Lambda = \{\lambda \in \mathbb{R}_+^m : \|\lambda\|_1 \leq R\}$ , let  $\tilde{\ell}_i := \ell_i$  for all  $i \in [m]$ , and consider the empirical Lagrangian of Equation 15. Let  $\theta^{(1)}, \dots, \theta^{(T)} \in \Theta$  and  $\lambda^{(1)}, \dots, \lambda^{(T)} \in \Lambda$  be sequences satisfying:

$$\max_{\lambda^* \in \Lambda} \frac{1}{T} \sum_{t=1}^T \hat{\mathcal{L}}(\theta^{(t)}, \lambda^*) - \inf_{\theta^* \in \Theta} \frac{1}{T} \sum_{t=1}^T \hat{\mathcal{L}}(\theta^*, \lambda^{(t)}) \leq \epsilon \quad (16)$$

Define  $\bar{\theta}$  as a random variable for which  $\bar{\theta} = \theta^{(t)}$  with probability  $1/T$ , and let  $\bar{\lambda} := (\sum_{t=1}^T \lambda^{(t)})/T$ . Then  $\bar{\theta}$  is nearly-optimal in expectation:

$$\mathbb{E}_{\bar{\theta}, x \sim \mathcal{D}} [\ell_0(x; \bar{\theta})] \leq \inf_{\theta^* \in \Theta: \forall i, \mathbb{E}_{x \sim \mathcal{D}}[\ell_i(x; \theta^*)] \leq 0} \mathbb{E}_{x \sim \mathcal{D}} [\ell_0(x; \theta^*)] + \epsilon + 2\tilde{G}^{(\text{trn})}(\Theta)$$

and nearly-feasible:

$$\max_{i \in [m]} \mathbb{E}_{\bar{\theta}, x \sim \mathcal{D}} [\ell_i(x; \bar{\theta})] \leq \frac{\epsilon}{R - \|\bar{\lambda}\|_1} + \tilde{G}^{(\text{trn})}(\Theta)$$

Additionally, if there exists a  $\theta' \in \Theta$  that satisfies all of the constraints with margin  $\gamma$  (i.e.  $\mathbb{E}_{x \sim \mathcal{D}} [\ell_i(x; \theta')] \leq -\gamma$  for all  $i \in [m]$ ), then:

$$\|\bar{\lambda}\|_1 \leq \frac{\epsilon + B_{\ell_0}}{\gamma - \tilde{G}^{(\text{trn})}(\Theta)}$$

assuming that  $\gamma > \tilde{G}^{(\text{trn})}(\Theta)$ , where  $B_{\ell_0} \geq \sup_{\theta \in \Theta} \mathbb{E}_{x \sim \mathcal{D}} [\ell_0(x; \theta)] - \inf_{\theta \in \Theta} \mathbb{E}_{x \sim \mathcal{D}} [\ell_0(x; \theta)]$  is a bound on the range of the objective loss.

*Proof.* The empirical Lagrangian is nothing but the in-expectation Lagrangian over the finite training sample  $S^{(\text{trn})}$ , so by Theorem 1 of Cotter et al. (2019),  $\bar{\theta}$  is nearly-optimal in expectation:

$$\mathbb{E}_{\bar{\theta}, x \sim S^{(\text{trn})}} [\ell_0(x; \bar{\theta})] \leq \inf_{\theta^* \in \Theta: \forall i, \mathbb{E}_{x \sim S^{(\text{trn})}} [\ell_i(x; \theta^*)] \leq 0} \mathbb{E}_{x \sim S^{(\text{trn})}} [\ell_0(x; \theta^*)] + \epsilon$$

and nearly-feasible:

$$\max_{i \in [m]} \mathbb{E}_{\tilde{\theta}, x \sim S^{(\text{trn})}} [\ell_i(x; \tilde{\theta})] \leq \frac{\epsilon}{R - \|\tilde{\lambda}\|_1}$$

Since  $\theta'$  satisfies the constraints with margin  $\gamma$  in expectation, it will satisfy them with margin  $\gamma - \tilde{G}^{(\text{trn})}(\Theta)$  on the training dataset, so the same theorem gives the claimed upper-bound  $\|\tilde{\lambda}\|_1 \leq (\epsilon + B_{\ell_0}) / (\gamma - \tilde{G}^{(\text{trn})}(\Theta))$  when  $\gamma > \tilde{G}^{(\text{trn})}(\Theta)$ .

Notice that the above expectations are taken over the finite training sample  $S^{(\text{trn})}$ , rather than the data distribution  $\mathcal{D}$ . To fix this, we need only define  $\tilde{\ell}_i = \ell_i$ , and appeal to the definition of  $\tilde{G}^{(\text{trn})}(\Theta)$  (Definition 2), yielding the claimed results.  $\square$

Here,  $R$  is the maximum allowed 1-norm of the vector of Lagrange multipliers (such a bound is necessary for Cotter et al. (2019)'s proof to work out). Notice that we have assumed that  $\tilde{\ell}_i := \ell_i$  for all  $i$ . This is purely for notational reasons (the Lagrangian does not involve proxy constraints *at all*)—it allows us to re-use the definition of  $\tilde{G}^{(\text{trn})}(\Theta)$  (Definition 2) in the above Theorem, and causes  $\gamma$  to have the same definition here, as it did in Theorems 1 and 2.