

## CONTRIBUTIONS

1. *What does it mean to learn under constraints?*  
Define PAC constrained (PACC) learning
2. *When (if at all) is it possible to learn under constraints?*  
Constrained learning is as hard as unconstrained learning (PACC  $\Leftrightarrow$  PAC)
3. *Is there a practical constrained learning rule?*  
Dual empirical learning and primal-dual algorithms (under mild conditions)

## INTRODUCTION

Learning is ubiquitous, but has serious shortcomings



Problem (Constrained learning)

Train a CNN that is *robust to input perturbations*

$$\begin{aligned} & \underset{\theta \in \mathbb{R}^p}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \\ & \text{subject to } \mathbb{E}_{(\tilde{\mathbf{x}}, \tilde{y}) \sim \mathcal{A}} [\ell(f_{\theta}(\tilde{\mathbf{x}}), \tilde{y})] \leq c \end{aligned} \quad \approx \quad ?$$

### Unconstrained alternatives

► **Learning** (PAC learning theory)

$$\underset{\theta \in \mathbb{R}^p}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \approx \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n)$$

- May not satisfy the requirements

► **Regularized learning** (PAC learning theory)

$$\begin{aligned} & \underset{\theta \in \mathbb{R}^p}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] + \lambda \mathbb{E}_{(\tilde{\mathbf{x}}, \tilde{y}) \sim \mathcal{A}} [\ell(f_{\theta}(\tilde{\mathbf{x}}), \tilde{y})] \\ & \approx \frac{1}{N} \sum_{n=1}^N [\ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \ell(f_{\theta}(\tilde{\mathbf{x}}_n), \tilde{y}_n)] \end{aligned}$$

- For what  $\lambda$  is the solution feasible? Does it generalize or is it dataset-dependent?

## CONSTRAINED LEARNING THEORY

### 1. What is constrained learning?

Definition (PACC learnability)

$\mathcal{H}$  is PACC learnable if for every  $\epsilon, \delta$  and every distributions  $\mathcal{D}_k$ , we can obtain  $f_{\theta^\dagger} \in \mathcal{H}$  from  $N_{\mathcal{H}}(\epsilon, \delta)$  samples that is, with probability  $1 - \delta$ ,

$$\begin{aligned} & \text{near-optimal } (\Rightarrow \text{PAC learning}) \quad \text{approximately feasible} \\ & \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta^\dagger}(\mathbf{x}), y)] \leq P^* + \epsilon \quad \mathbb{E}_{(\tilde{\mathbf{x}}, \tilde{y}) \sim \mathcal{A}} [\ell(f_{\theta^\dagger}(\tilde{\mathbf{x}}), \tilde{y})] \leq c + \epsilon \end{aligned}$$

### 2. When (if at all) is it possible to learn under constraints?

Theorem 1

$\mathcal{H}$  is PACC learnable  $\Leftrightarrow$   $\mathcal{H}$  is PAC learnable

### 3. Is there a practical constrained learning rule?

Theorem 2

Let  $f$  be  $\nu$ -universal, i.e., for each  $\theta_1, \theta_2$ , and  $\gamma \in [0, 1]$  there exists  $\theta$  such that  $\mathbb{E}[\gamma f_{\theta_1}(\mathbf{x}) + (1 - \gamma)f_{\theta_2}(\mathbf{x}) - f_{\theta}(\mathbf{x})] \leq \nu$  and  $\ell$  be convex, bounded, and  $M$ -Lipschitz continuous. Then  $\hat{D}^*$  is a (near-)PACC learner, i.e., if  $\theta^\dagger$  achieves  $\hat{D}^*$ , then with probability  $1 - \delta$ ,

$$\begin{aligned} & \mathbb{E}[\ell(f_{\theta^\dagger}(\mathbf{x}), y)] \leq P^* + \epsilon_0 + \epsilon \quad \mathbb{E}[\ell(f_{\theta^\dagger}(\tilde{\mathbf{x}}), \tilde{y})] \leq c + \epsilon \\ & \epsilon_0 = (2 + \|\lambda_p^*\|_1) M \nu \quad \epsilon = B \sqrt{\frac{1}{N} \left[ 1 + \log \left( \frac{4(m+2)(2N)^{dvc}}{\delta} \right) \right]} \end{aligned}$$

A primal-dual algorithm

$$\begin{aligned} \text{Minimize the primal: } \theta^+ & \approx \underset{\theta \in \mathbb{R}^p}{\text{argmin}} \frac{1}{N} \sum_{n=1}^N [\ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \ell(f_{\theta}(\tilde{\mathbf{x}}_n), \tilde{y}_n)] \\ & = \theta - \eta [\nabla_{\theta} \ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \nabla_{\theta} \ell(f_{\theta}(\tilde{\mathbf{x}}_n), \tilde{y}_n)] \end{aligned}$$

$$\text{Update the dual: } \lambda^+ = \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta^+}(\tilde{\mathbf{x}}_n), \tilde{y}_n) - c \right) \right]_+$$

Theorem 3

$\theta^+$  is a  $\rho$ -approximate minimizer  $\Rightarrow$  converges to a  $\rho$ -neighborhood of  $\hat{D}^*$  (under mild conditions)

## MAIN RESULTS

$$\begin{aligned} P^* & = \underset{\theta \in \mathbb{R}^p}{\text{min}} \mathbb{E} [\ell_0(f_{\theta}(\mathbf{x}), y)] \\ & \text{s. to } \mathbb{E} [\ell_i(f_{\theta}(\tilde{\mathbf{x}}), \tilde{y})] \leq c \end{aligned} \quad \xleftrightarrow{\text{Thm. 1}} \quad \begin{aligned} & \underset{\theta \in \mathbb{R}^p}{\text{min}} \frac{1}{N} \sum_{n=1}^N \ell_0(f_{\theta}(\mathbf{x}_n), y_n) \\ & \text{s. to } \frac{1}{N} \sum_{n=1}^N \ell_i(f_{\theta}(\tilde{\mathbf{x}}_n), \tilde{y}_n) \leq c \end{aligned}$$

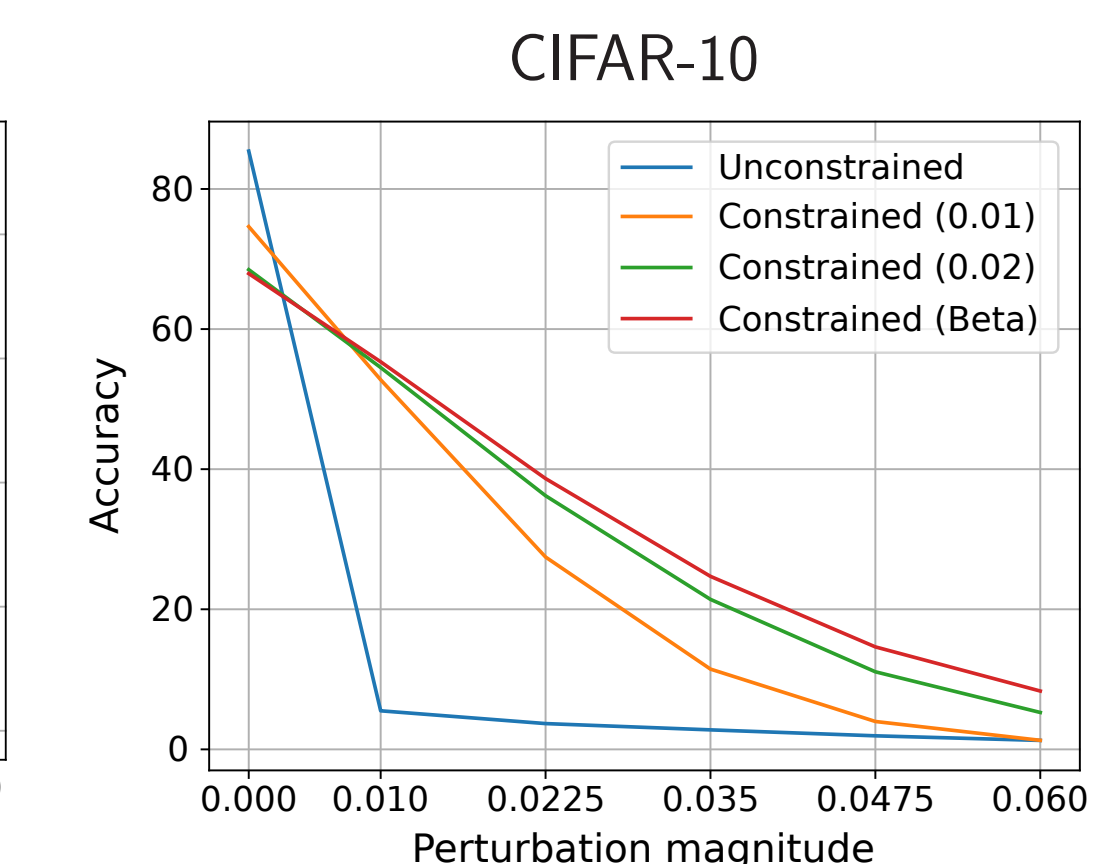
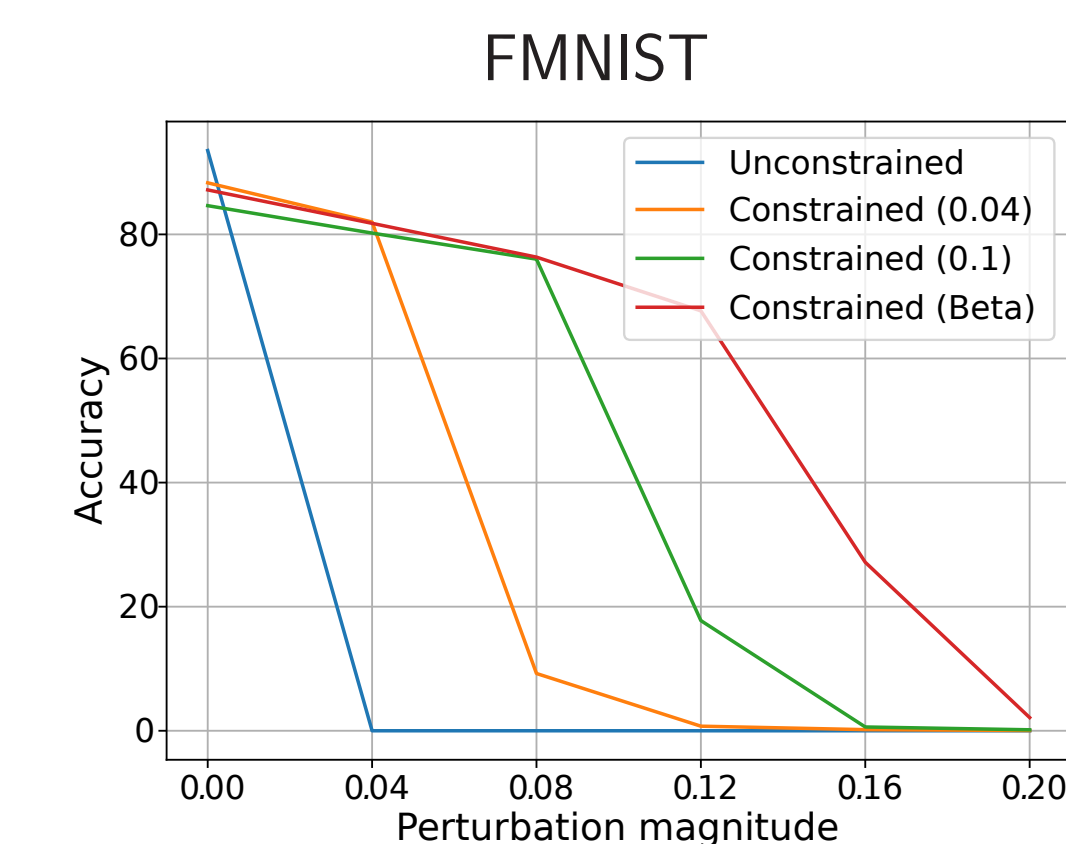
$$\xleftrightarrow{\text{Thm. 2}} \quad \hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{n=1}^N \left( \ell_0(f_{\theta}(\mathbf{x}_n), y_n) + \lambda [\ell(f_{\theta}(\tilde{\mathbf{x}}_n), \tilde{y}_n) - c] \right)$$

In the paper...

- Multiple constraints, different losses
- Pointwise constraints:  $\ell(f_{\theta}(\mathbf{x}), y) \leq c$ ,  $\mathcal{D}$ -a.e.

## APPLICATION

### Robustness



And also...

