



Luiz F. O. Chamon

Workshop USP-IPP
Mar. 10th, 2024



**learning
under
requirements**

System engineering cycle

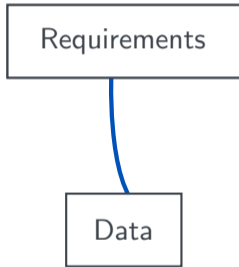


System engineering cycle

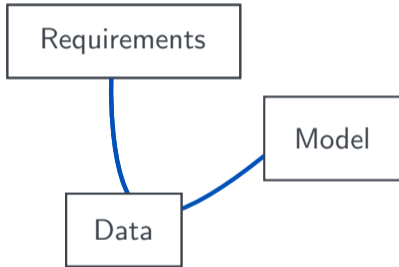
Requirements



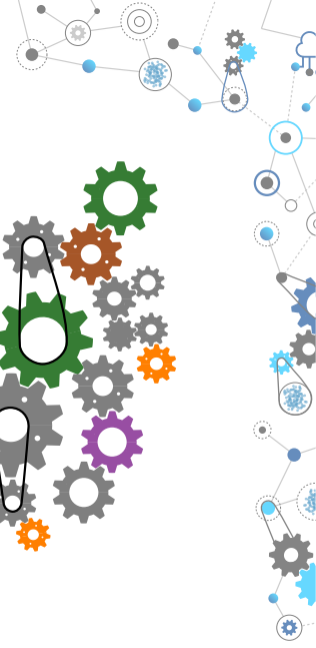
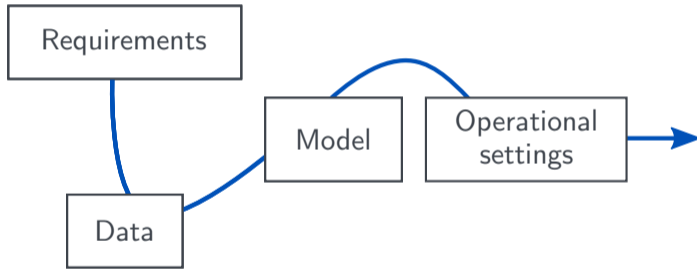
System engineering cycle



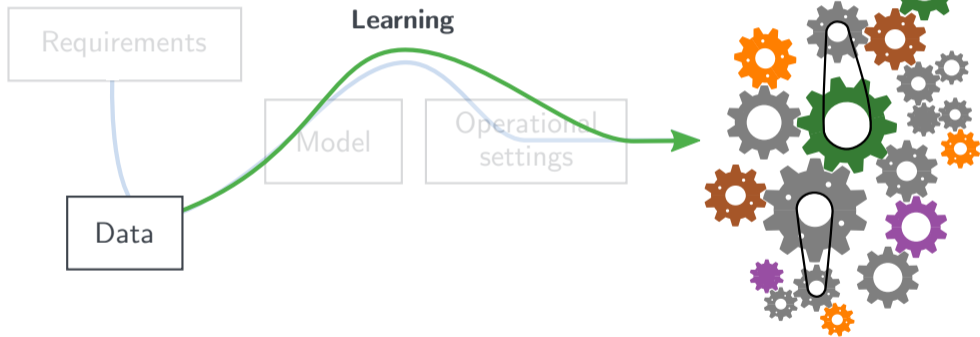
System engineering cycle



System engineering cycle



The promise of learning



The promise emerging reality of learning

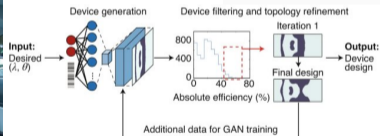
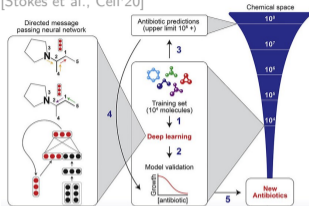


Microsoft Azure will become the preferred cloud platform for Johns Hopkins inHealth precision medicine initiative

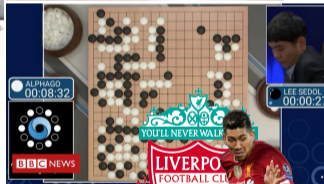
June 18, 2020 | Microsoft News Center

Five-year agreement will support Johns Hopkins Medicine inHealth in driving new medical discoveries to improve disease management and patient care

[Stokes et al., Cell'20]



[Ma et al., Nature Photonics'20]



[Piggott et al., Nature Photonics'15]

The promise emerging reality of learning

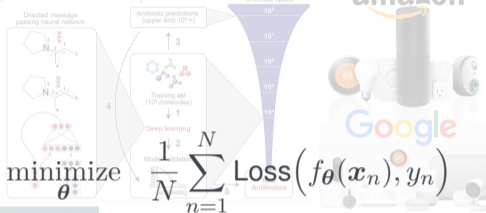


Microsoft Azure will be the preferred cloud platform for Johns Hopkins inHealth precision medicine initiative

June 18, 2020 | Microsoft News Center

Five-year agreement will support Johns Hopkins Medicine inHealth in driving new medical discoveries to improve disease management and patient care

[Stokes et al., Cell'20]



amazon



maximize θ

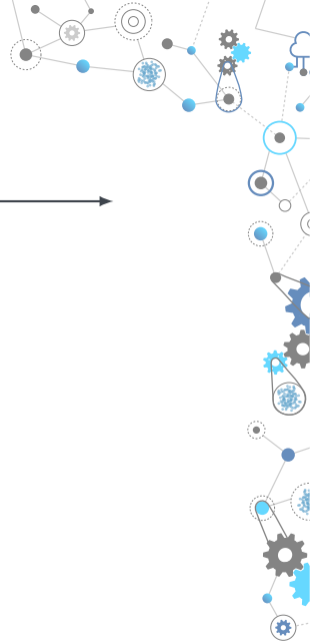


[Ma et al., Nature Photonics'20]



[Piggott et al., Nature Photonics'15]

Learning breakthroughs



Learning breakthroughs

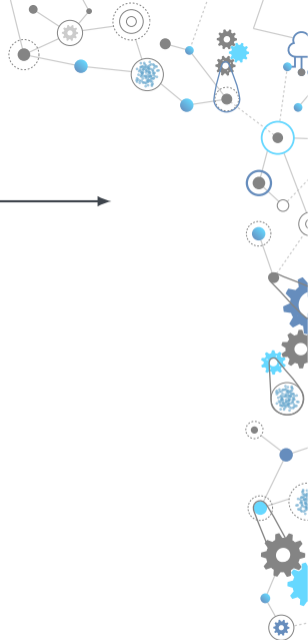


Classical learning theory [Vapnik & Chervonenkis, TP'71; Valiant, CACM'84]:

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) \xrightarrow{\text{"LLN"}} \min_{\theta} \mathbb{E} \left[\text{Loss}(f_{\theta}(\mathbf{x}), y) \right]$$

- e.g., linear functions, smooth functions (finite RKHS norm, bandlimited), NNs. . .

Learning breakthroughs



The promise emerging reality of learning

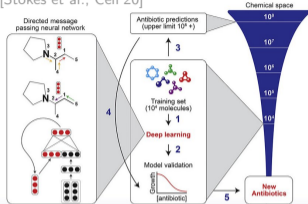


Microsoft Azure will become the preferred cloud platform for Johns Hopkins inHealth precision medicine initiative

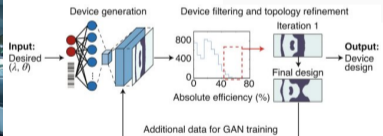
June 18, 2020 | Microsoft News Center

Five-year agreement will support Johns Hopkins Medicine inHealth in driving new medical discoveries to improve disease management and patient care

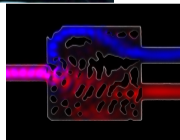
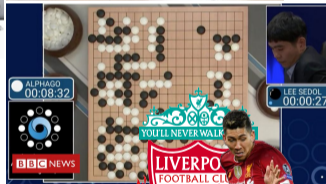
[Stokes et al., Cell'20]



amazon



[Ma et al., Nature Photonics'20]



[Piggott et al., Nature Photonics'15]

The promise limitations of learning



The New York Times

Self-Driving Uber Car Kills Pedestrian in Arizona, Where Robots Roam

cloud platform for Johns Hopkins inHealth precision medicine initiative

AAAS Become a Member

Science

Read our COVID-19 research and news.

RESEARCH ARTICLE

Dissecting racial bias in an algorithm used to manage the health of populations

How nonprofit and business leaders can equitably and responsibly use AI systems in the fight against COVID-19.

Stanford SOCIAL INNOVATION Review

The Problem With COVID-19 Artificial Intelligence Solutions and How to Fix Them

How nonprofit and business leaders can equitably and responsibly use AI systems in the fight against COVID-19.

REUTERS

Amazon scraps secret AI recruiting tool that showed bias against women

By Jeffrey Dastin

TO KNOWLEDGE

How Auto Encourag Flawed A

By CHLOE HADAV MARCH 06, 2020



MIT Technology Review

The way we train AI is fundamentally flawed

Artificial intelligence / Machine learning

The process used to build most of the machine-learning models we use today can't tell if they will work in the real world or not—and that's a problem.

THE APPEAL

ALGORITHMS OF INEQUALITY

POLITICAL REPORT

MIT Technology Review

Facebook's ad-serving algorithm discriminates by gender and race

Artificial intelligence / Machine learning

Even if an advertiser is well-intentioned, the algorithm still prefers certain groups of people over others.

PRO PUBLICA

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

THE APPEAL

Tesla Auto Found Pro in 2018 Cr

The National Transportation Safety Board called for electric-car companies to feature and cite agencies.

[Prggott et al., Nature Photonics '15]

Improving ERM



Improving ERM

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

[Xie & Yuille, ICLR'20; Guo et al., CVPR'20; Finzi et al., ICML'20; Li et al., ICRL'21;
Lu et al., Nature Mach. Intel.'21; Raissi et al., J. Comp. Phys.'19; ...]



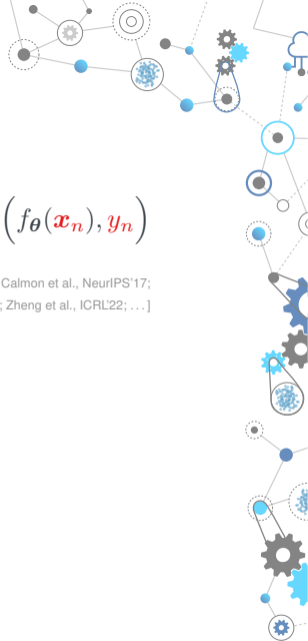
Improving ERM

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

[Xie & Yuille, ICLR'20; Guo et al., CVPR'20; Finzi et al., ICML'20; Li et al., ICRL'21;
Lu et al., Nature Mach. Intel.'21; Raissi et al., J. Comp. Phys.'19; ...]

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

[Kamiran & Calders, KIS'12; Feldman et al., SIGKDD'15; Calmon et al., NeurIPS'17;
Chen et al., ICML'20; Müller & Hutter, ICCV'21; Zheng et al., ICRL'22; ...]



Improving ERM

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

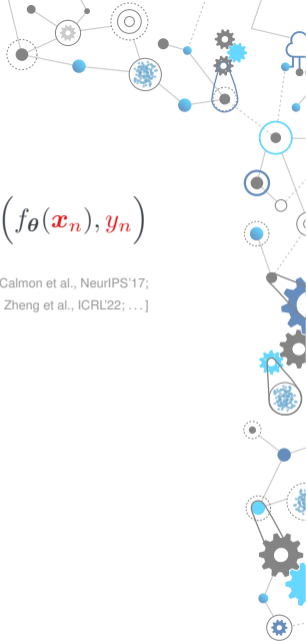
[Xie & Yuille, ICLR'20; Guo et al., CVPR'20; Finzi et al., ICML'20; Li et al., ICRL'21;
Lu et al., Nature Mach. Intel.'21; Raissi et al., J. Comp. Phys.'19; ...]

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

[Goodfellow et al., ICLR'15; Arjovsky et al., ICML'17; Madry et al., ICLR'18;
Zhang et al., ICML'19; Raissi et al., J. Comp. Phys.'19; Krishnan et al., NeurIPS'20; ...]

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

[Kamiran & Calders, KIS'12; Feldman et al., SIGKDD'15; Calmon et al., NeurIPS'17;
Chen et al., ICML'20; Müller & Hutter, ICCV'21; Zheng et al., ICRL'22; ...]



Improving ERM

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

[Xie & Yuille, ICLR'20; Guo et al., CVPR'20; Finzi et al., ICML'20; Li et al., ICRL'21;
Lu et al., Nature Mach. Intel.'21; Raissi et al., J. Comp. Phys.'19; ...]

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

[Goodfellow et al., ICLR'15; Arjovsky et al., ICML'17; Madry et al., ICLR'18;
Zhang et al., ICML'19; Raissi et al., J. Comp. Phys.'19; Krishnan et al., NeurIPS'20; ...]

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

[Kamiran & Calders, KIS'12; Feldman et al., SIGKDD'15; Calmon et al., NeurIPS'17;
Chen et al., ICML'20; Müller & Hutter, ICCV'21; Zheng et al., ICRL'22; ...]

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

[Helmbold & Long, JMLR'15; Mianjy et al., ICML'18; Tashiro et al., NeurIPS'20;
Li et al., AISTATS'20; Lin et al., ICML'20; Foret et al., ICRL'21; ...]

A different paradigm...



Learning is doing exactly
what we asked for.



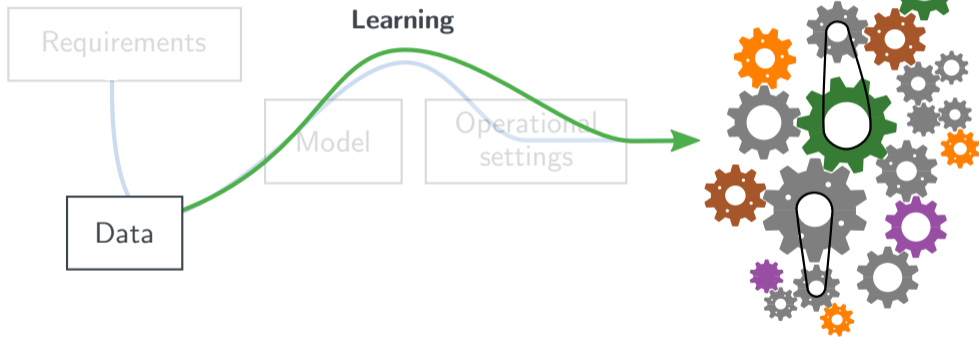
A different paradigm...



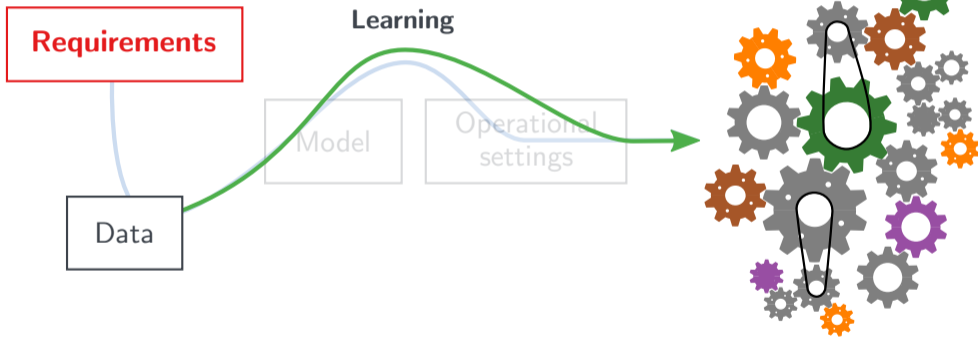
Learning is doing exactly
what we asked for.

How can AI learn to do what we want?

The promise of learning



The promise of learning



A different paradigm...



Learning is doing exactly
what we asked for.

How can AI learn to do what we want?

Constrained learning

Claims

Constrained learning is the right way to learn under requirements

Constrained learning is hard...

...but possible



Claims

Constrained learning is the right way to learn under requirements

Constrained learning is hard...

... but possible



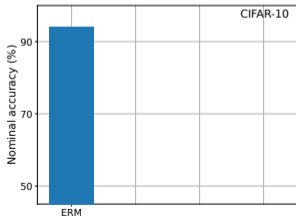
Robust image recognition

Problem

Learn an image classifier



Cello



Robust image recognition

Problem

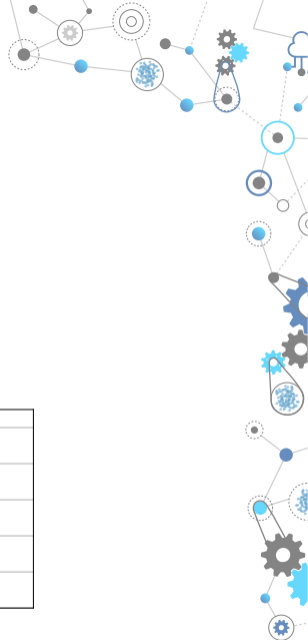
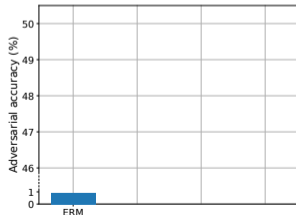
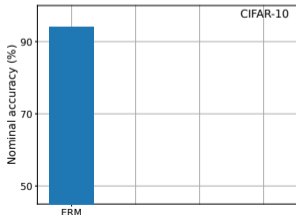
Learn an image classifier



Cello



Hammer



Robust image recognition

Problem

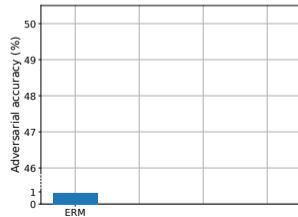
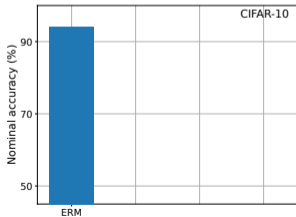
Learn an image classifier **that is robust to input perturbations**



Cello



Hammer



Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

- Adversarial training (e.g., [Szegedy et al., ICLR'14; Goodfellow et al., ICLR'15; Madry et al., ICLR'18])

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) \longrightarrow \min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(\mathbf{x}_n + \delta), y_n) \right]$$



Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

- Adversarial training (e.g., [Szegedy et al., ICLR'14; Goodfellow et al., ICLR'15; Madry et al., ICLR'18])

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) \longrightarrow \min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(\mathbf{x}_n + \delta), y_n) \right]$$



\approx gradient ascent

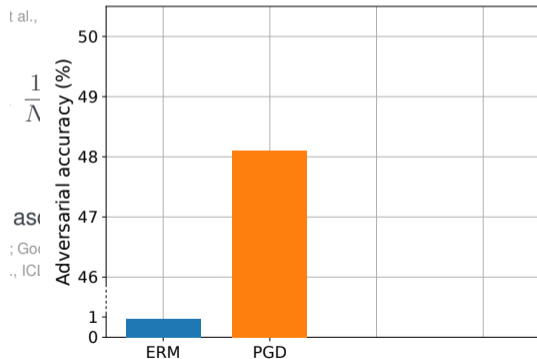
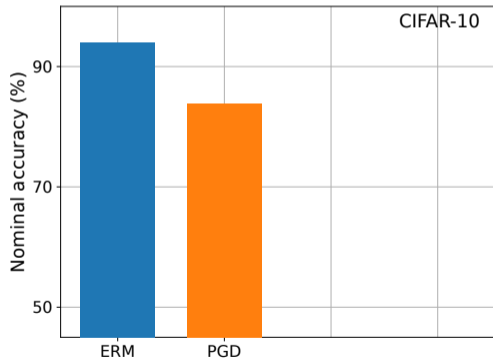
[Szegedy et al., ICLR'14; Goodfellow et al., ICLR'15; Madry et al., ICLR'18; ...]



Adversarial training

Problem

Learn an image classifier that is robust to input perturbations



Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

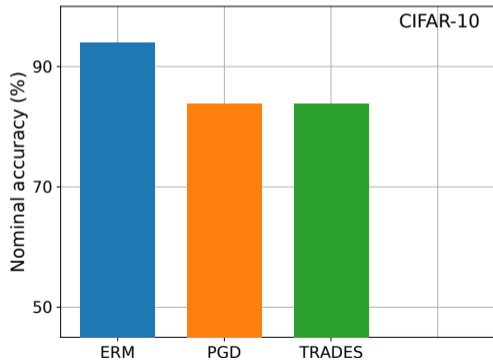
- Adversarial training (e.g., [Zhang et al., ICML'19])

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) \longrightarrow \min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(\mathbf{x}_n + \delta), y_n) \right]$$
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(\mathbf{x}_n + \delta), y_n) \right]$$

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations



Penalty-based methods

Problem

Learn an image classifier that is robust to input perturbations

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(\mathbf{x}_n + \delta), y_n) \right]$$

- ✘ No straightforward relation between λ and **adversarial loss**
- ✘ λ depends on the values of the losses (dataset, model, performance measure)
- ✘ Requirement generalization

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

$$\begin{aligned} \min_{\theta} \quad & \text{Nominal loss} \\ \text{subject to} \quad & \text{Adversarial loss} \leq c \end{aligned}$$



Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n)$$

subject to **Adversarial loss** $\leq c$



Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) \\ \text{subject to} \quad & \frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(\mathbf{x}_n + \delta), y_n) \right] \leq c \end{aligned}$$

- ✓ More natural: requirement is a constraint, not a cost
- ✓ Decouple performance and requirements



Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

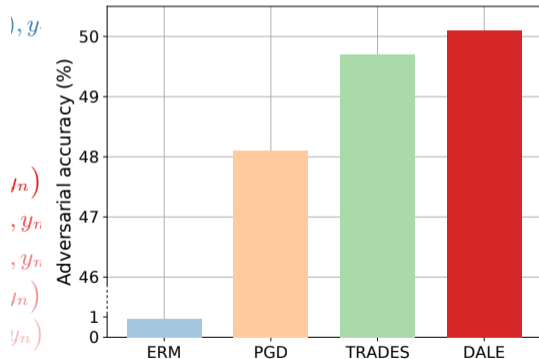
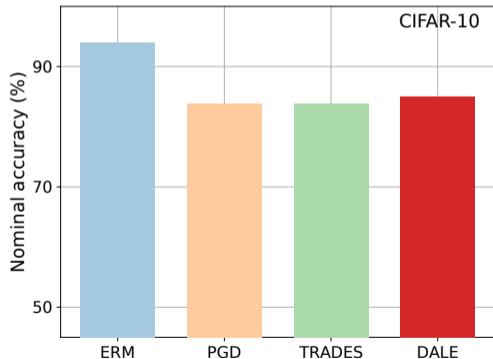
$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) \\ \text{subject to} \quad & \frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(\mathbf{x}_n + \delta), y_n) \right] \leq c \end{aligned}$$



Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

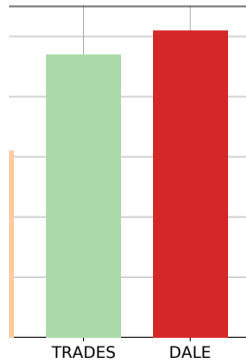
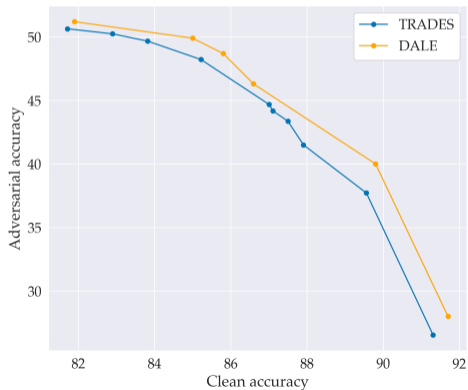
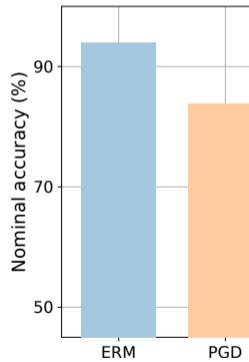


$$\mathcal{L}_{\text{ERM}}(\theta; \mathcal{D}) = \mathbb{E}_{(x, y) \sim \mathcal{D}} \ell(\theta; x, y)$$

Constrained learning for robustness

Problem

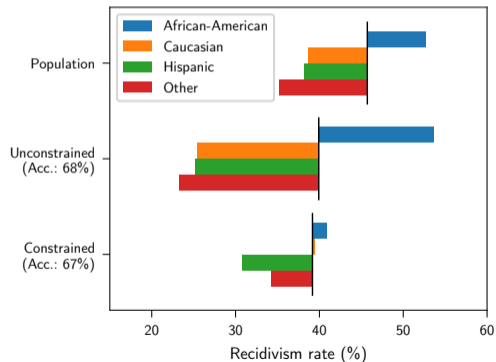
Learn an image classifier that is robust to input perturbations



Fair learning

Problem

Predict whether an individual will recidivate at the same rate across races



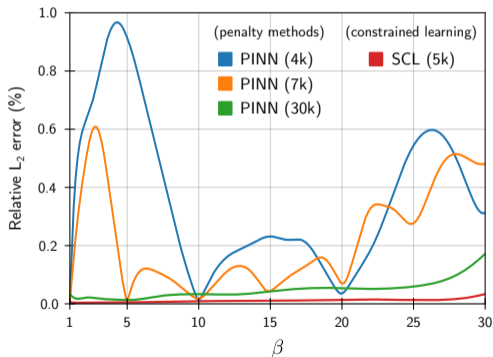
minimize Prediction error

subject to Prediction rate disparity (Race) $\leq c$

Learning to solve PDEs

Problem

Obtain (weak) solutions for a parametric family of boundary value problems

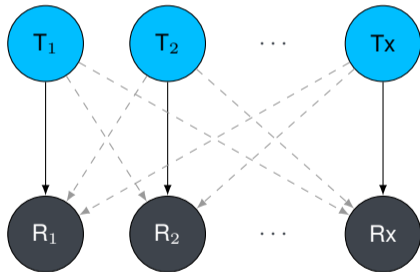


minimize Boundary condition error
subject to Weak formulation error $\leq \epsilon, \forall(x, t)$

Wireless resource allocation

Problem

Allocate the least transmit power to m device pairs to achieve a communication rate

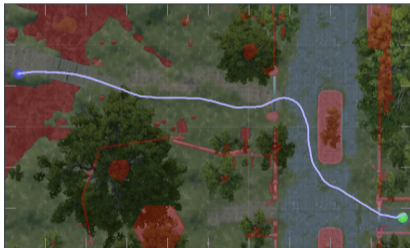


minimize Total power
subject to Communicate rate $(T_i) \geq c_i$

Safe reinforcement learning

Problem

Learn a control policy that navigates the environment effectively and safely



maximize Task reward

subject to $\mathbb{P}(\text{Colliding with obstacles}) \leq \delta$

Claims

Constrained learning is the right way to learn under requirements

Constrained learning is hard...

... but possible



(Un)constrained learning

$$\hat{P}_U^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n)$$

- ℓ, g are bounded, Lipschitz continuous (possibly non-convex) functions
- f_{θ} is a (possibly nonlinear) parametrization [e.g., logistic classifier, (G)(C)NN]
- $(\mathbf{x}_n, y_n) \sim \mathcal{D}, (\mathbf{x}_m, y_m) \sim \mathcal{A}$ (i.i.d.)



(Un)constrained learning

$$\hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n)$$

subject to $\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$

- ℓ, g are bounded, Lipschitz continuous (possibly non-convex) functions
- f_{θ} is a (possibly nonlinear) parametrization [e.g., logistic classifier, (G)(C)NN]
- $(\mathbf{x}_n, y_n) \sim \mathcal{D}, (\mathbf{x}_m, y_m) \sim \mathcal{A}$ (i.i.d.)



(Un)constrained learning

$$\hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n)$$

subject to $\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$

- ℓ, g are bounded, Lipschitz continuous (possibly non-convex) functions
- f_{θ} is a (possibly nonlinear) parametrization [e.g., logistic classifier, (G)(C)NN]
- $(\mathbf{x}_n, y_n) \sim \mathcal{D}, (\mathbf{x}_m, y_m) \sim \mathcal{A}$ (i.i.d.)



Constrained learning

$$\begin{array}{l} \hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) \\ \text{subject to } \frac{1}{N} \sum_{m=1}^M g(f_{\theta}(\mathbf{x}_m), y_m) \leq c \end{array} \quad \xrightarrow{\quad ? \quad} \quad \begin{array}{l} P^* = \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \\ \text{subject to } \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{A}} [g(f_{\theta}(\mathbf{x}), y)] \leq c \end{array}$$

Challenges

- ✘ *Statistical*: does the solution of the constrained empirical problem generalize?

Constrained learning

$$\hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) \xrightarrow{\text{"LLN"}} P^* = \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)]$$

subject to $\frac{1}{N} \sum_{m=1}^M g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$ subject to $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [g(f_{\theta}(\mathbf{x}), y)] \leq c$

Challenges

- ✘ *Statistical*: does the solution of the constrained empirical problem generalize?

Constrained learning

$$\begin{array}{l} \hat{P}^* = \min_{\theta} \quad \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) \\ \text{subject to} \quad \frac{1}{N} \sum_{m=1}^M g(f_{\theta}(\mathbf{x}_m), y_m) \leq c \end{array} \quad \xrightarrow{\quad ? \quad} \quad \begin{array}{l} P^* = \min_{\theta} \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \\ \text{subject to} \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{A}} [g(f_{\theta}(\mathbf{x}), y)] \leq c \end{array}$$

Challenges

- ✘ *Statistical*: does the solution of the constrained empirical problem generalize?

Constrained learning

$$\hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n)$$

subject to $\frac{1}{N} \sum_{m=1}^M g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$

?

$$P^* = \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)]$$

subject to $\mathbb{E}_{(x,y) \sim \mathcal{Q}} [g(f_{\theta}(x), y)] \leq c$

Challenges

- ✘ *Statistical*: does the solution of the constrained empirical problem generalize?
- ✘ *Computational*: can we solve the constrained empirical problem?

Constrained learning

$$\hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n)$$

NON-CONVEX ?

subject to $\frac{1}{N} \sum_{m=1}^M g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$

$$P^* = \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)]$$

subject to $\mathbb{E}_{(x,y) \sim \mathcal{Q}} [g(f_{\theta}(x), y)] \leq c$

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \frac{1}{N} \sum_{m=1}^M g(f_{\theta}(\mathbf{x}_m), y_m)$$

Challenges

- ⊗ *Statistical*: does the solution of the constrained empirical problem generalize?
- ⊗ *Computational*: can we solve the constrained empirical problem?

Constrained learning

$$\begin{array}{l} \hat{P}^* = \min_{\theta} \quad \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) \\ \text{subject to} \quad \frac{1}{N} \sum_{m=1}^M g(f_{\theta}(\mathbf{x}_m), y_m) \leq c \end{array} \quad \xrightarrow{\quad ? \quad} \quad \begin{array}{l} P^* = \min_{\theta} \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \\ \text{subject to} \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{A}} [g(f_{\theta}(\mathbf{x}), y)] \leq c \end{array}$$

Challenges

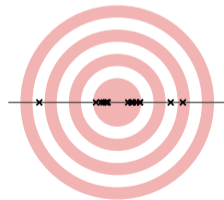
- ✘ *Statistical*: does the solution of the constrained empirical problem generalize?
- ✘ *Computational*: can we solve the constrained empirical problem?

What classical learning theory says?

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) \xrightarrow{\text{"LLN"}} \min_{\theta} \mathbb{E} [\text{Loss}(f_{\theta}(\mathbf{x}), y)]$$

✓ f_{θ} is *probably approximately correct (PAC)* learnable

e.g., linear functions, smooth functions (finite RKHS norm, bandlimited), NNs...
($N \approx 1/\epsilon^2$)

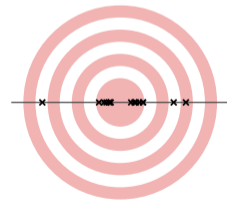


What classical learning theory says?

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) \xrightarrow{\text{"LLN"}} \min_{\theta} \mathbb{E} [\text{Loss}(f_{\theta}(\mathbf{x}), y)]$$

- ✓ f_{θ} is *probably approximately correct (PAC)* learnable
e.g., linear functions, smooth functions (finite RKHS norm, bandlimited), NNs...
($N \approx 1/\epsilon^2$)

- ✗ Constraints?



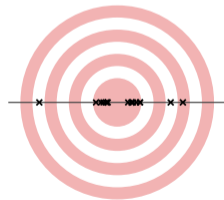
What's in a solution?

Definition (PAC learnability)

f_{θ} is a *probably approximately correct* (PAC) learnable if for every ϵ, δ and every distributions \mathcal{D}, \mathcal{Q} , we can obtain $f_{\theta^{\dagger}}$ from $N_f(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,

- near-optimal

$$P^* - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\ell(f_{\theta^{\dagger}}(\mathbf{x}), y) \right] \leq \epsilon$$



What's in a solution?

Definition (PACC learnability)

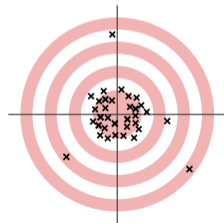
f_{θ} is a *probably approximately correct constrained (PACC)* learnable if for every ϵ, δ and every distributions \mathcal{D}, \mathcal{Y} , we can obtain $f_{\theta^{\dagger}}$ from $N_f(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,

- near-optimal

$$\left| P^* - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\ell(f_{\theta^{\dagger}}(\mathbf{x}), y) \right] \right| \leq \epsilon$$

- approximately feasible

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{Y}} \left[g(f_{\theta^{\dagger}}(\mathbf{x}), y) \right] \leq c + \epsilon$$



When is constrained learning possible?

$$\begin{array}{l} \hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) \\ \text{subject to } \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) \leq c \end{array} \xrightarrow{\text{PACC}} \begin{array}{l} P^* = \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \\ \text{subject to } \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{A}} [g(f_{\theta}(\mathbf{x}), y)] \leq c \end{array}$$

Proposition

f_{θ} is PAC learnable $\not\Rightarrow$ f_{θ} is PAC**C** learnable

Claims

Constrained learning is the right way to learn under requirements

Constrained learning is hard...


...but possible



Constrained learning

$$\hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n)$$

subject to $\frac{1}{N} \sum_{m=1}^M g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$

PAGG 

$$P^* = \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)]$$

subject to $\mathbb{E}_{(x,y) \sim \mathcal{Q}} [g(f_{\theta}(x), y)] \leq c$

Challenges

- ⊗ *Statistical*: does the solution of the constrained empirical problem generalize?
- ⊗ *Computational*: can we solve the constrained empirical problem?

Duality

PRIMAL



DUAL



Duality

$$\hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) \quad \text{subject to} \quad \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$$



DUAL



Duality

$$\hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) \quad \text{subject to} \quad \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$$



$$\hat{D}^* = \max_{\lambda \geq 0} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \left[\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) - c \right]$$

Duality

$$\hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) \quad \text{subject to} \quad \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$$



$$\hat{D}^* = \max_{\lambda \geq 0} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \left[\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) - c \right]$$

- In general, $\hat{D}^* \leq \hat{P}^*$
- But in some cases, $\hat{D}^* = \hat{P}^*$ (strong duality) [e.g., convex optimization]

Duality

$$\hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) \quad \text{subject to} \quad \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$$



$$\hat{D}^* = \max_{\lambda \geq 0} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \left[\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) - c \right]$$

- In general, $\hat{D}^* \leq \hat{P}^*$
- But in some cases, $\hat{D}^* = \hat{P}^*$ (strong duality) [e.g., convex optimization]

Non-convex variational duality

Convex optimization: Primal \longleftrightarrow Dual

Non-convex, finite dimensional optimization: Primal \nleftrightarrow Dual



Non-convex variational duality

Convex optimization: Primal \longleftrightarrow Dual

Non-convex, finite dimensional optimization: Primal \nleftrightarrow Dual

Non-convex, infinite dimensional optimization: Primal \longleftrightarrow Dual



Sparse logistic regression

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} - \sum_{n=1}^N \log \left[1 + \exp \left(y_n \cdot \boldsymbol{\theta}^T \mathbf{x}_n \right) \right]$$

$$\text{s. to } \|\boldsymbol{\theta}\|_0 = \sum_{t=1}^p \mathbb{I}[\theta_t \neq 0] \leq k$$

Discrete, non-convex

[Chen et al., JMLR'19]: NP-hard



Sparse logistic regression

$$\min_{\theta \in \mathbb{R}^p} - \sum_{n=1}^N \log \left[1 + \exp \left(y_n \cdot \theta^T x_n \right) \right]$$

$$\text{s. to } \|\theta\|_0 = \sum_{t=1}^p \mathbb{I}[\theta_t \neq 0] \leq k$$

Discrete, non-convex

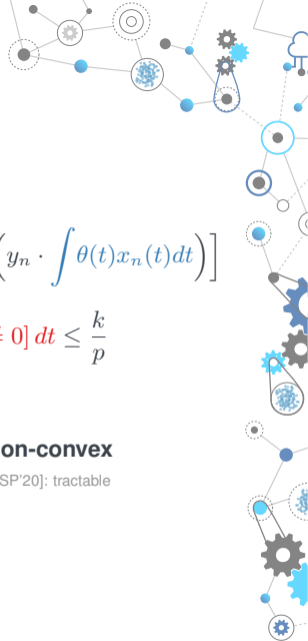
[Chen et al., JMLR'19]: NP-hard

$$\min_{\theta \in L_2} - \sum_{n=1}^N \log \left[1 + \exp \left(y_n \cdot \int \theta(t) x_n(t) dt \right) \right]$$

$$\text{s. to } \|\theta\|_{L_0} = \int \mathbb{I}[\theta(t) \neq 0] dt \leq \frac{k}{p}$$

Continuous, non-convex

[Chamon et al., IEEE TSP'20]: tractable



Sparse logistic regression

$$\min_{\theta \in \mathbb{R}^p} - \sum_{n=1}^N \log \left[1 + \exp \left(y_n \cdot \theta^T x_n \right) \right]$$

$$\text{s. to } \|\theta\|_0 = \sum_{t=1}^p \mathbb{I}[\theta_t \neq 0] \leq k$$

Discrete, non-convex

[Chen et al., JMLR'19]: NP-hard

$$\min_{\theta \in L_2} - \sum_{n=1}^N \log \left[1 + \exp \left(y_n \cdot \int \theta(t) x_n(t) dt \right) \right]$$

$$\text{s. to } \|\theta\|_{L_0} = \int \mathbb{I}[\theta(t) \neq 0] dt \leq \frac{k}{p}$$

Continuous, non-convex

[Chamon et al., IEEE TSP'20]: tractable

How to learn under constraints?



$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n)$$

$$\text{subject to } \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) \leq c$$

PAGG \rightarrow

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)]$$

$$\text{subject to } \mathbb{E}_{(x,y) \sim \mathcal{D}} [g(f_{\theta}(x), y)] \leq c$$

\swarrow

$$\max_{\lambda \geq 0} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \left[\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) - c \right]$$

How to learn under constraints?



$$\begin{array}{ccc} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) & \xrightarrow{\text{PAGG}} & \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \\ \text{subject to } \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) \leq c & & \text{subject to } \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{A}} [g(f_{\theta}(\mathbf{x}), y)] \leq c \\ & & \nearrow \\ \max_{\lambda \geq 0} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \left[\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) - c \right] & & \end{array}$$

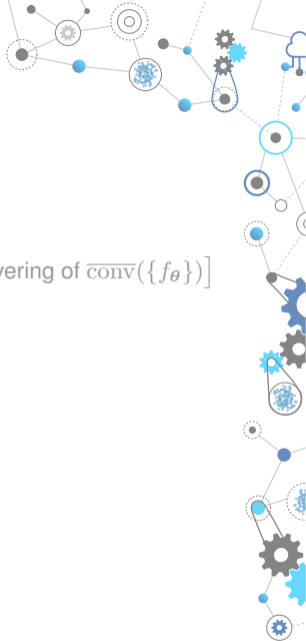
Dual (near-)PACC learning

Theorem

Let f be ν -universal, i.e., for each θ_1, θ_2 , and $\gamma \in [0, 1]$ there exists θ such that

$$\mathbb{E} \left[\left| \gamma f_{\theta_1}(\mathbf{x}) + (1 - \gamma) f_{\theta_2}(\mathbf{x}) - f_{\theta}(\mathbf{x}) \right| \right] \leq \nu$$

$[\{f_{\theta}\} \text{ is a good covering of } \overline{\text{conv}}(\{f_{\theta}\})]$



Dual (near-)PACC learning

Theorem

Let f be ν -universal, i.e., for each θ_1, θ_2 , and $\gamma \in [0, 1]$ there exists θ such that

$$\mathbb{E} \left[|\gamma f_{\theta_1}(\mathbf{x}) + (1 - \gamma) f_{\theta_2}(\mathbf{x}) - f_{\theta}(\mathbf{x})| \right] \leq \nu$$

Then \hat{D}^* is a (near-)PACC learner, i.e., for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , with probability $1 - \delta$,

Near-optimal:
$$|P^* - \hat{D}^*| \leq \tilde{O} \left(\nu + \frac{1}{\sqrt{N}} \right)$$

Approximately feasible:
$$\mathbb{E} \left[g(f_{\theta^\dagger}(\mathbf{x}), y) \right] \leq c + \tilde{O} \left(\nu + \frac{1}{\sqrt{N}} \right)$$

(ℓ strongly convex and g convex)

(mild additional conditions apply)

Dual (near-)PACC learning

Theorem

Let f be ν -universal with VC dimension $d_{\text{VC}} < \infty$, ℓ strongly convex, and g convex. Then, f_{θ^\dagger} is a (near-)PACC solution of (P-CSL) for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , i.e., with probability at least $1 - \delta$,

$$\left| P^* - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\ell(f_{\theta^\dagger}(\mathbf{x}), y) \right] \right| \leq (1 + \Delta)(\epsilon_0 + \epsilon)$$

$$\mathbb{E} \left[g(f_{\theta^\dagger}(\mathbf{x}), y) \right] \leq c + (1 + \Delta)^{3/2} (M\sqrt{\epsilon_0} + \epsilon)$$

$$\epsilon_0 = M\nu \quad \epsilon = B \sqrt{\frac{1}{N} \left[1 + \log \left(\frac{4m(2N)^{d_{\text{VC}}}}{\delta} \right) \right]} \quad \Delta = \max \left(\|\lambda^*\|_1, \|\hat{\lambda}^*\|_1, \|\tilde{\lambda}^*\|_1 \right)$$

Sources of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

Dual (near-)PACC learning

Theorem

Let f be ν -universal with VC dimension $d_{\text{VC}} < \infty$, ℓ strongly convex, and g convex. Then, f_{θ^\dagger} is a (near-)PACC solution of (P-CSL) for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , i.e., with probability at least $1 - \delta$,

$$\left| P^* - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\ell(f_{\theta^\dagger}(\mathbf{x}), y) \right] \right| \leq (1 + \Delta)(\epsilon_0 + \epsilon)$$

$$\mathbb{E} \left[g(f_{\theta^\dagger}(\mathbf{x}), y) \right] \leq c + (1 + \Delta)^{3/2} (M\sqrt{\epsilon_0} + \epsilon)$$

$$\epsilon_0 = M\nu \quad \epsilon = B \sqrt{\frac{1}{N} \left[1 + \log \left(\frac{4m(2N)^{d_{\text{VC}}}}{\delta} \right) \right]} \quad \Delta = \max \left(\|\lambda^*\|_1, \|\hat{\lambda}^*\|_1, \|\tilde{\lambda}^*\|_1 \right)$$

Sources of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

Dual (near-)PACC learning

Theorem

Let f be ν -universal with VC dimension $d_{\text{VC}} < \infty$, ℓ strongly convex, and g convex. Then, f_{θ^\dagger} is a (near-)PACC solution of (P-CSL) for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , i.e., with probability at least $1 - \delta$,

$$\left| P^* - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\ell(f_{\theta^\dagger}(\mathbf{x}), y) \right] \right| \leq (1 + \Delta)(\epsilon_0 + \epsilon)$$

$$\mathbb{E} \left[g(f_{\theta^\dagger}(\mathbf{x}), y) \right] \leq c + (1 + \Delta)^{3/2} (M\sqrt{\epsilon_0} + \epsilon)$$

$$\epsilon_0 = M\nu \quad \epsilon = B \sqrt{\frac{1}{N} \left[1 + \log \left(\frac{4m(2N)^{d_{\text{VC}}}}{\delta} \right) \right]} \quad \Delta = \max \left(\|\lambda^*\|_1, \|\hat{\lambda}^*\|_1, \|\tilde{\lambda}^*\|_1 \right)$$

Sources of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

Dual (near-)PACC learning

Theorem

Let f be ν -universal with VC dimension $d_{\text{VC}} < \infty$, ℓ strongly convex, and g convex. Then, f_{θ^\dagger} is a (near-)PACC solution of (P-CSL) for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , i.e., with probability at least $1 - \delta$,

$$\left| P^* - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\ell(f_{\theta^\dagger}(\mathbf{x}), y) \right] \right| \leq (1 + \Delta)(\epsilon_0 + \epsilon)$$

$$\mathbb{E} \left[g(f_{\theta^\dagger}(\mathbf{x}), y) \right] \leq c + (1 + \Delta)^{3/2} (M\sqrt{\epsilon_0} + \epsilon)$$

$$\epsilon_0 = M\nu \quad \epsilon = B \sqrt{\frac{1}{N} \left[1 + \log \left(\frac{4m(2N)^{d_{\text{VC}}}}{\delta} \right) \right]} \quad \Delta = \max \left(\|\lambda^*\|_1, \|\hat{\lambda}^*\|_1, \|\tilde{\lambda}^*\|_1 \right)$$

Sources of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

Dual (near-)PACC learning

Theorem

Let f be ν -universal with VC dimension $d_{VC} < \infty$, ℓ strongly convex, and g convex. Then, f_{θ^\dagger} is a (near-)PACC solution of (P-CSL) for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , i.e., with probability at least $1 - \delta$,

$$\left| P^* - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\ell(f_{\theta^\dagger}(\mathbf{x}), y) \right] \right| \leq (1 + \Delta)(\epsilon_0 + \epsilon)$$

$$\mathbb{E} \left[g(f_{\theta^\dagger}(\mathbf{x}), y) \right] \leq c + (1 + \Delta)^{3/2} (M\sqrt{\epsilon_0} + \epsilon)$$

$$\epsilon_0 = M\nu \quad \epsilon = B \sqrt{\frac{1}{N} \left[1 + \log \left(\frac{4m(2N)^{d_{VC}}}{\delta} \right) \right]} \quad \Delta = \max \left(\|\lambda^*\|_1, \|\hat{\lambda}^*\|_1, \|\tilde{\lambda}^*\|_1 \right)$$

Sources of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

Dual learning trade-offs

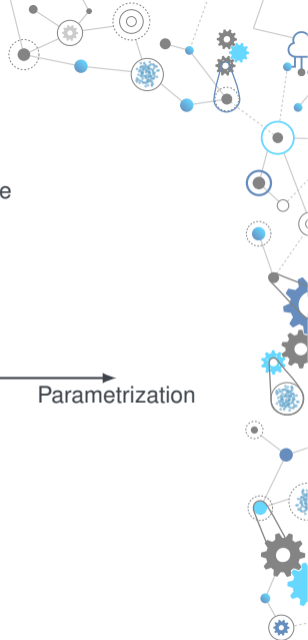
- Unconstrained learning

parametrization \times sample size

Sample size

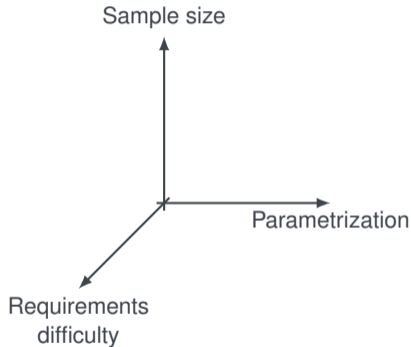


Parametrization



Dual learning trade-offs

- Unconstrained learning
parametrization \times sample size
- Constrained learning
parametrization \times sample size \times requirements



When is constrained learning possible?

Corollary

f_θ is PAC learnable \approx^* f_θ is PAC**C** learnable

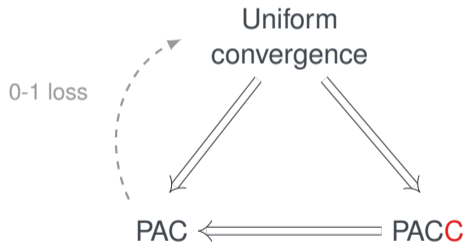
Constrained learning is **essentially as hard as** unconstrained learning

(mild conditions apply)



When is constrained learning possible?

Corollary



(mild conditions apply)

Claims

Constrained learning is the right way to learn under requirements

Constrained learning is hard...

...but possible. How?



Dual learning algorithm

$$\hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \left[\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) - c \right]$$



Dual learning algorithm

- Minimize the primal (\equiv **ERM**)

$$\theta^\dagger \in \operatorname{argmin}_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{n=1}^N \left[\ell(f_\theta(\mathbf{x}_n), y_n) + \lambda g(f_\theta(\mathbf{x}_n), y_n) \right]$$

$$\hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta(\mathbf{x}_n), y_n) + \lambda \left[\frac{1}{N} \sum_{m=1}^N g(f_\theta(\mathbf{x}_m), y_m) - c \right]$$



Dual learning algorithm

- Minimize the primal (\equiv ERM)

$$\theta^+ \approx \theta - \eta \nabla_{\theta} \left[\ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda g(f_{\theta}(\mathbf{x}_n), y_n) \right], \quad n = 1, 2, \dots$$

[Ge et al., ICLR'18; Soltanolkotabi et al., IEEE TIT'18; Mei et al., PNAS'18; Kawaguchi et al., AISTATS'20...]

$$\hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \left[\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) - c \right]$$

Dual learning algorithm

- Minimize the primal (\equiv ERM)

$$\theta^+ \approx \theta - \eta \nabla_{\theta} \left[\ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda g(f_{\theta}(\mathbf{x}_n), y_n) \right], \quad n = 1, 2, \dots$$

- Update the dual

$$\lambda^+ = \left[\lambda + \eta \left(\frac{1}{N} \sum_{m=1}^N g(f_{\theta^+}(\mathbf{x}_m), y_m) - c \right) \right]_+$$

$$\hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), y_n) + \lambda \left[\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(\mathbf{x}_m), y_m) - c \right]$$

A (near-)PACC learner

Theorem

Suppose θ^\dagger is a ρ -approximate solution of the regularized ERM:

$$\theta^\dagger \approx \operatorname{argmin}_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{n=1}^N \left(\ell(f_\theta(\mathbf{x}_n), y_n) + \lambda g(f_\theta(\mathbf{x}_n), y_n) \right).$$

Then, after $T = \left\lceil \frac{\|\lambda^\star\|^2}{2\eta M\nu} \right\rceil + 1$ dual iterations with step size $\eta \leq \frac{2\epsilon}{mB^2}$,

the iterates $(\theta^{(T)}, \lambda^{(T)})$ are such that

$$\left| P^\star - L(\theta^{(T)}, \lambda^{(T)}) \right| \leq (2 + \Delta)(\epsilon_0 + \epsilon) + \rho$$

with probability $1 - \delta$ over sample sets.



In practice...

- 1: Initialize: θ_0, λ_0
- 2: **for** $t = 1, \dots, T$
- 3: $\beta_1 \leftarrow \theta_{t-1}$
- 4: **for** $n = 1, \dots, N$
- 5: $\beta_{n+1} \leftarrow \beta_n - \eta_{\theta} \nabla_{\beta} [\ell(f_{\beta_n}(\mathbf{x}_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(\mathbf{x}_n), y_n)]$
- 6: **end**
- 7: $\theta_t \leftarrow \beta_{N+1}$
- 8: $\lambda_t = \left[\lambda_{t-1} + \eta_{\lambda} \left(\frac{1}{N} \sum_{m=1}^N g(f_{\theta_t}(\mathbf{x}_m), y_m) - c \right) \right]_+$
- 9: **end**
- 10: Output: θ_T, λ_T

SGD

Dual update

 PyTorch

<https://github.com/lfochamon/csl>

Penalty-based vs. dual learning

Penalty-based learning

$$\theta^\dagger \in \operatorname{argmin}_{\theta} \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$$

- Parameter: λ (data-dependent)
- Generalizes with respect to $\text{Loss} + \lambda \text{Penalty}$

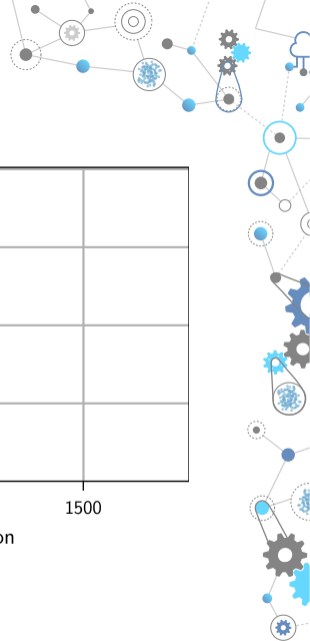
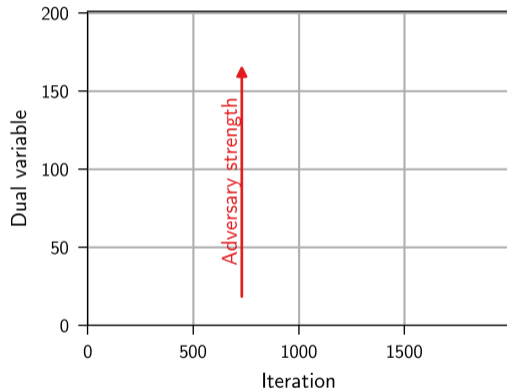
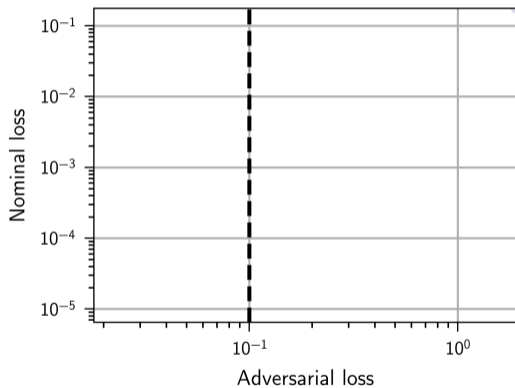
Dual learning

$$\theta^\dagger \in \operatorname{argmin}_{\theta} \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$$

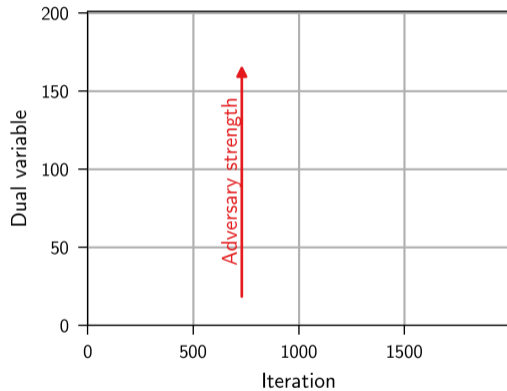
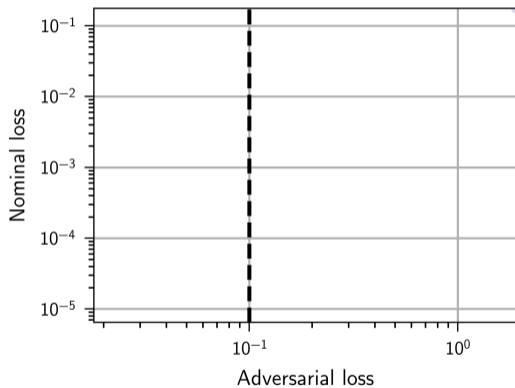
$$\lambda^+ = \left[\lambda + \eta \left(\text{Penalty}(\theta^\dagger) - c \right) \right]_+$$

- Parameter: c (requirement-dependent)
- Generalizes with respect to Loss and $\text{Penalty} \leq c$

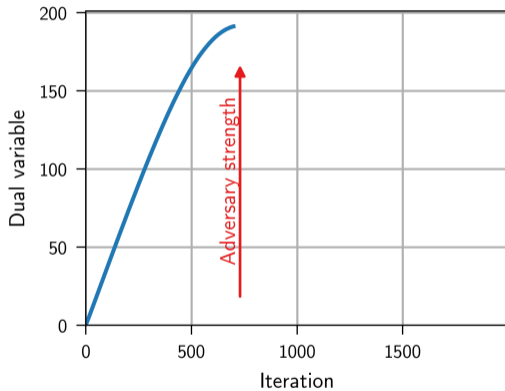
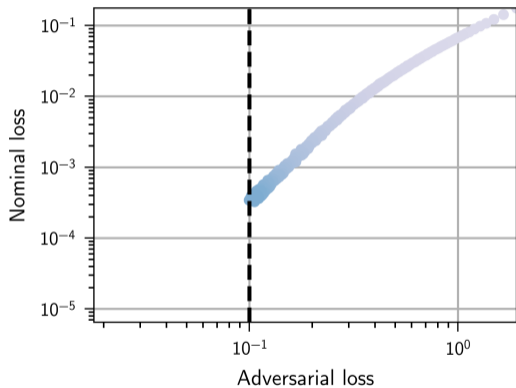
Robust image recognition



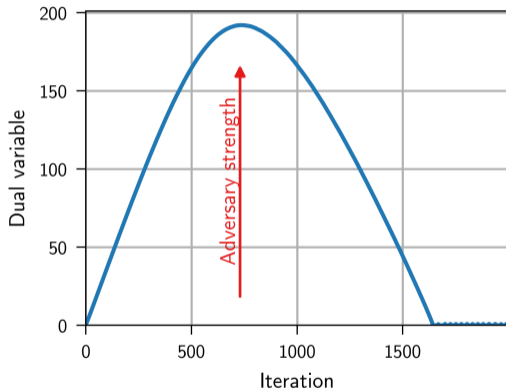
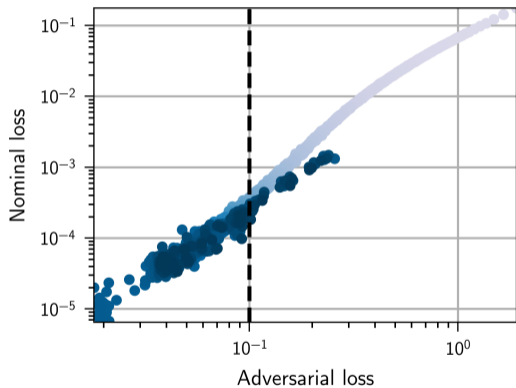
Robust image recognition



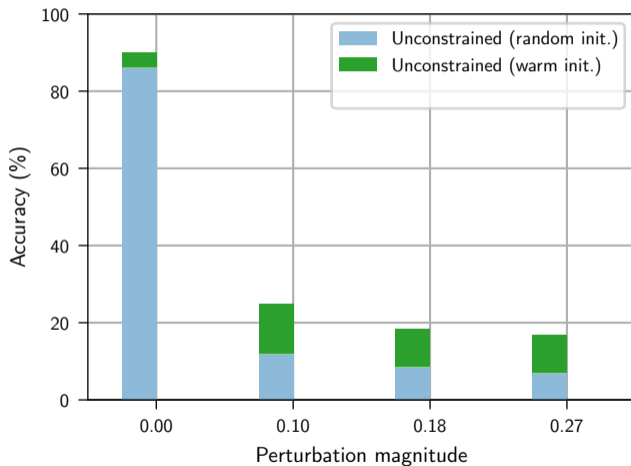
Robust image recognition



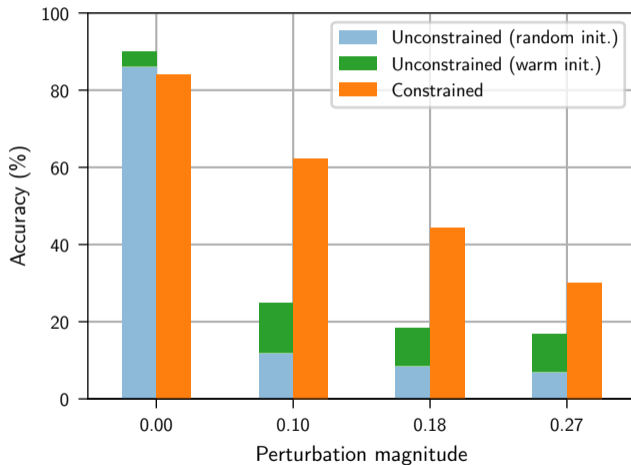
Robust image recognition



Robust image recognition



Robust image recognition



Claims

- **Constrained learning is the right tool to learn under requirements**

- **Constrained learning is hard...**

- **... but possible**



Claims

- **Constrained learning is ~~the right~~ a good tool to learn under requirements**

Constrained learning imposes requirements during training that generalize at test time, e.g.,

- **robustness** [CR, NeurIPS'20; R* \mathcal{C} *PH, NeurIPS'21; RCPH, ICML'22 (spotlight); CPCR, IEEE TIT'23]
- **fairness** [CPCR, ICASSP'20 (best student paper); CR, NeurIPS'20; CPCR, IEEE TIT'23]
- **invariance and data augmentation** [HCR, ICML'23]
- **(manifold) smoothness** [CCHVR, ICML'23]
- **resilience** [HRC, NeurIPS'23]
- **safe RL** [PCCR, NeurIPS'19; PCCR, IEEE TAC'23; CPCR, IEEE TAC'24]
- **learning to solve PDEs** [MC, ICLR'25]
- ...

- **Constrained learning is hard...**

- **... but possible**

Claims

- **Constrained learning is ~~the right~~ a good tool to learn under requirements**

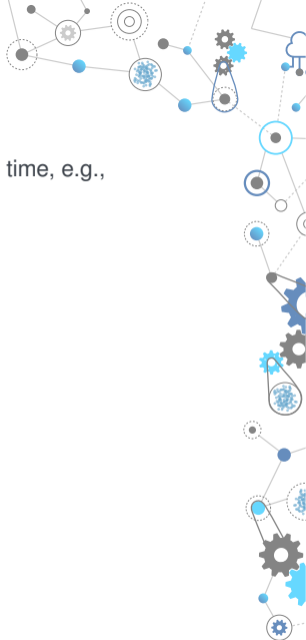
Constrained learning imposes requirements during training that generalize at test time, e.g.,

- **robustness** [CR, NeurIPS'20; R*C*PH, NeurIPS'21; RCPH, ICML'22 (spotlight); CPCR, IEEE TIT'23]
- **fairness** [CPCR, ICASSP'20 (best student paper); CR, NeurIPS'20; CPCR, IEEE TIT'23]
- **invariance and data augmentation** [HCR, ICML'23]
- **(manifold) smoothness** [CCHVR, ICML'23]
- **resilience** [HRC, NeurIPS'23]
- **safe RL** [PCCR, NeurIPS'19; PCCR, IEEE TAC'23; CPCR, IEEE TAC'24]
- **learning to solve PDEs** [MC, ICLR'25]
- ...

- **Constrained learning is hard...**

Constrained, non-convex, statistical optimization problem

- **... but possible**



Claims

- **Constrained learning is ~~the right~~ a good tool to learn under requirements**

Constrained learning imposes requirements during training that generalize at test time, e.g.,

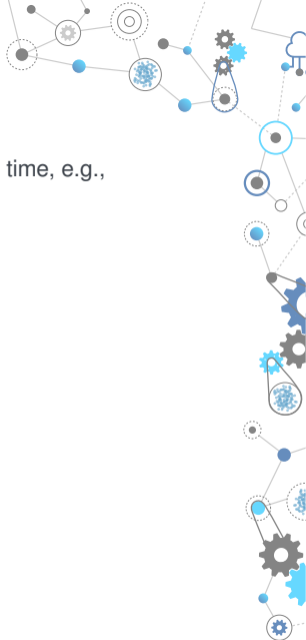
- **robustness** [CR, NeurIPS'20; R***C***PH, NeurIPS'21; RCPH, ICML'22 (spotlight); **CPCR**, IEEE TIT'23]
- **fairness** [CPCR, ICASSP'20 (best student paper); CR, NeurIPS'20; CPCR, IEEE TIT'23]
- **invariance and data augmentation** [HCR, ICML'23]
- **(manifold) smoothness** [CCHVR, ICML'23]
- **resilience** [HRC, NeurIPS'23]
- **safe RL** [PCCR, NeurIPS'19; PCCR, IEEE TAC'23; CPCR, IEEE TAC'24]
- **learning to solve PDEs** [MC, ICLR'25]
- ...

- **Constrained learning is hard...**

Constrained, non-convex, statistical optimization problem

- **... but possible**

We can learn under requirements (essentially) whenever we can learn at all



Claims

- **Constrained learning is ~~the right~~ a good tool to learn under requirements**

Constrained learning imposes requirements during training that generalize at test time, e.g.,

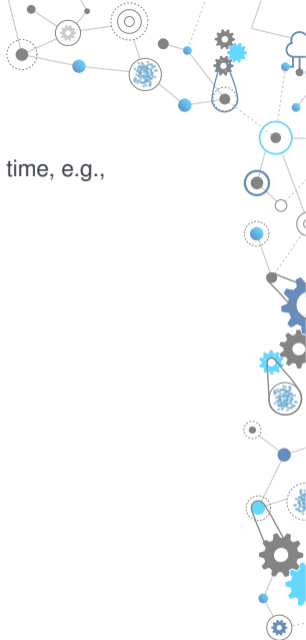
- **robustness** [CR, NeurIPS'20; R*C*PH, NeurIPS'21; RCPH, ICML'22 (spotlight); CPCR, IEEE TIT'23]
- **fairness** [CPCR, ICASSP'20 (best student paper); CR, NeurIPS'20; CPCR, IEEE TIT'23]
- **invariance and data augmentation** [HCR, ICML'23]
- **(manifold) smoothness** [CCHVR, ICML'23]
- **resilience** [HRC, NeurIPS'23]
- **safe RL** [PCCR, NeurIPS'19; PCCR, IEEE TAC'23; CPCR, IEEE TAC'24]
- **learning to solve PDEs** [MC, ICLR'25]
- ...

- **Constrained learning is hard...**

Constrained, non-convex, statistical optimization problem

- **... but possible. How?**

We can learn under requirements (essentially) whenever we can learn at all by solving (*penalized*) *ERM problems*



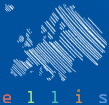


Collaborators:

Aneesh Barthakur, Miguel Calvo-Fullana,
Juan Cerviño, Mark Eisen, Yonina C. Eldar,
Hamed Hassani, Ignacio Hounie,
Dyonisios Kalogerias, Dan D. Lee, Viggo Moro,
George J. Pappas, Santiago Paternain,
Alejandro Ribeiro, Alexander Robey, Luana Ruiz,
Anastasios Tsiamis, Rene Vidal, Clark Zhang

www.luizchamon.com

luiz.chamon@polytechnique.edu



Luiz F. O. Chamon

**learning
under
requirements**