

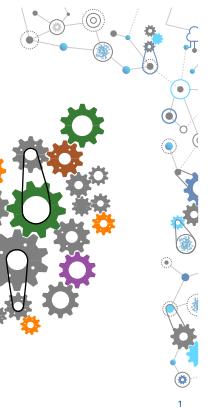
SimTech ML session
June 21, 2023



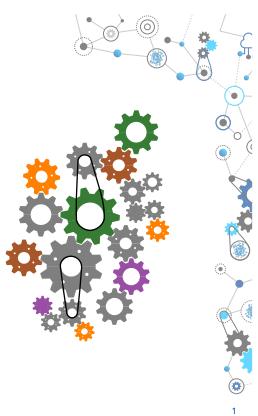
Luiz F. O. Chamon

learning under requirements

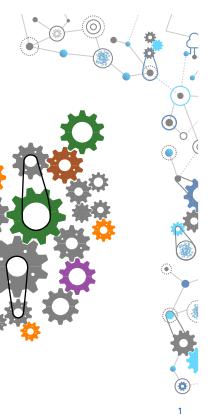
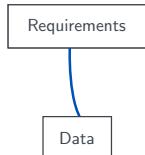
System engineering cycle



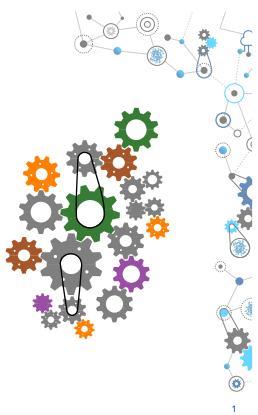
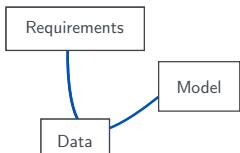
System engineering cycle



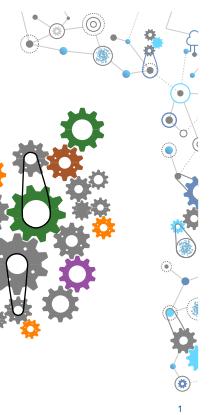
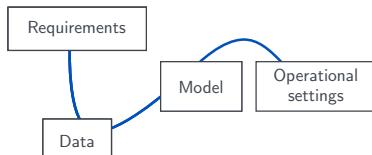
System engineering cycle



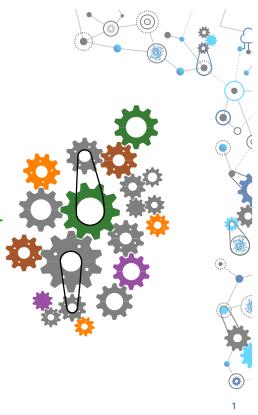
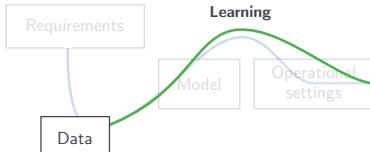
System engineering cycle



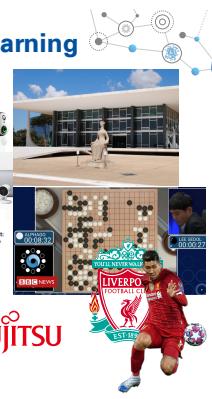
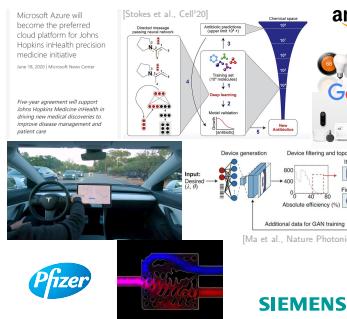
System engineering cycle



The promise of learning



The promise emerging reality of learning



The promise emerging reality of learning

Microsoft Azure will become the preferred cloud platform for Johns Hopkins iHealth precision medicine initiative
June 16, 2016 | Microsoft News Center



[Piggott et al., Nature Photonics'15]

SIEMENS

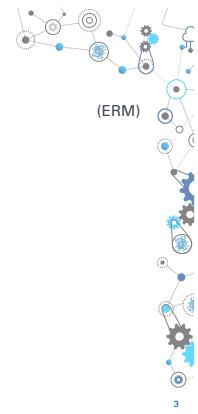
FUJITSU



2

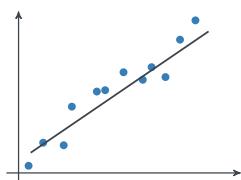
Empirical Risk Minimization

$$\underset{\theta}{\text{minimize}} \quad \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n)$$



Empirical Risk Minimization

$$\underset{\theta}{\text{minimize}} \quad \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n)$$

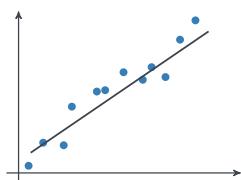


(ERM)

3

Empirical Risk Minimization

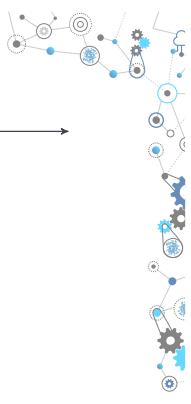
$$\underset{\theta}{\text{minimize}} \quad \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n)$$



3

Learning breakthroughs

→



4

Learning breakthroughs

→

1970s
(theoretical)
Learning theory

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

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

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

4

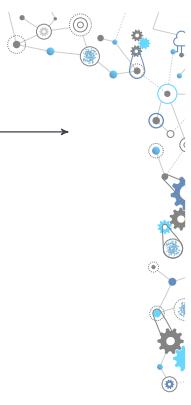
Learning breakthroughs

→

1970s
(theoretical)
Learning theory

2000s
(technological)
Internet

2010s
(computational)
Deep learning



4

The emerging reality of learning

Microsoft Azure will become the preferred cloud platform for Johns Hopkins iHealth precision medicine initiative
June 16, 2016 | Microsoft News Center

[Stokes et al., Cell'20]

amazon

Google

Siemens

Fujitsu

Pfizer

Liverpool

FC

Siemens

Siemens

Siemens

Siemens

Siemens

Device message processing

Input:

Output:

Device generation

Device filtering and topology refinement

Final design

Additional data for GAN training

News

Device

Design

News

Device

Design

News

Device

Design

Chemical vapor

Input:

Output:

Device

Design

News

Device

Design

News

Device

Design

News

Device

Design

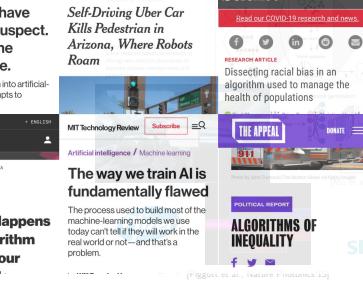
News

Device

5

[Piggott et al., Nature Photonics'15]

The emerging reality limitations of learning



The way we train AI is fundamentally flawed

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.

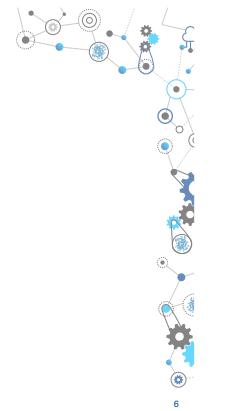
POLITICAL REPORT
ALGORITHMS OF INEQUALITY

TECHNOLOGY
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.
RETAIL
Amazon scraps secret AI recruiting tool that showed bias against women
By Jeffrey Dastin
MARCH 06, 2018 / 10:44 PM / UPDATED 2 YEARS AGO
ARTIFICIAL INTELLIGENCE / MACHINE LEARNING
Facebook's ad-serving algorithm discriminates by gender and race
Even if an advertiser is well-intentioned, the algorithm still prefers certain groups of people over others.

TECHNOLOGY
How Automation Bias Encourages the Use of Flawed Algorithms
By CHLOE HADAVAS
MARCH 06, 2018 / 8:24 PM / UPDATED 2 YEARS AGO
ARTIFICIAL INTELLIGENCE / MACHINE LEARNING
Tesla Autopilot System Found Probably at Fault in 2018 Crash
The National Transportation Safety Board called for improvements in the electric-car company's driver-assist feature and cited failures by other agencies.

5

Improving ERM



Improving ERM

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

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



Improving ERM

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

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

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

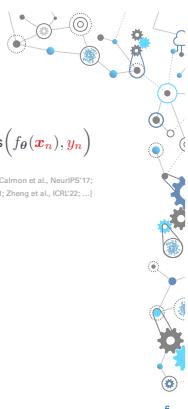
[Kamiran & Calders, KIS'12; Feldman et al., SIGKDD'15; Calmon et al., NeurIPS'17; Chen et al., ICML'20; Müller & Hutter, ICCV'21; Zheng et al., ICLR'22; ...]

Improving ERM

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

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

[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., ICLR'22; ...]



Improving ERM

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

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

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

[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., ICLR'22; ...]

A different paradigm...



Learning is doing exactly what we asked for.



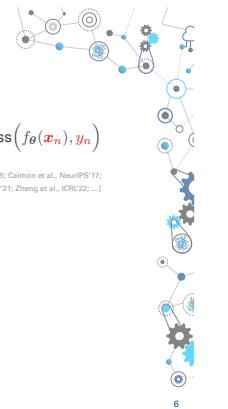
A different paradigm...



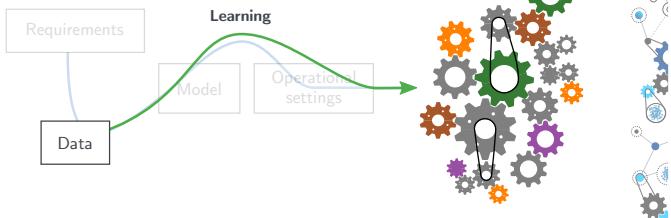
Learning is doing exactly what we asked for.

7

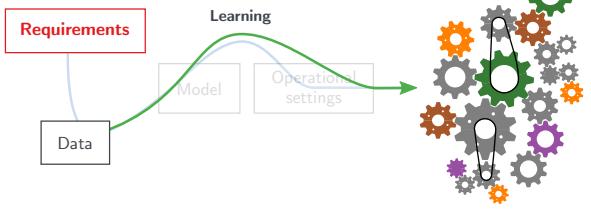
How can we learn models that do what we want?



The promise of learning



The promise of learning



A different paradigm...



Learning is doing exactly what we asked for.

How can we learn models that do what we want?

Constrained learning

9

10

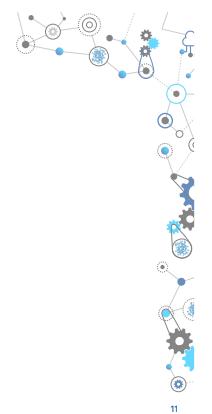
11

Claims

Constrained learning is the right tool to learn under requirements

Constrained learning is hard...

...but possible



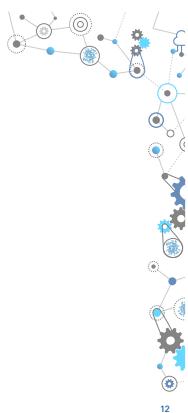
Claims

Constrained learning is the right tool to learn under requirements

Wireless resource allocation
Robust image recognition

Constrained learning is hard...

...but possible



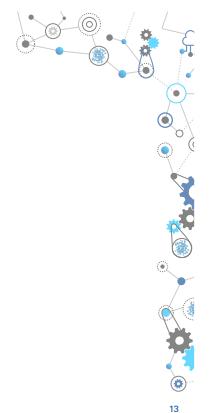
Claims

Constrained learning is the right tool to learn under requirements

Wireless resource allocation
Robust image recognition

Constrained learning is hard...

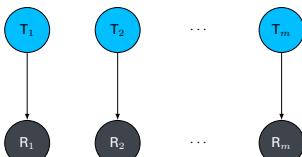
...but possible



Wireless resource allocation

Problem

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



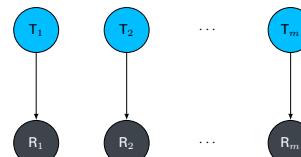
12



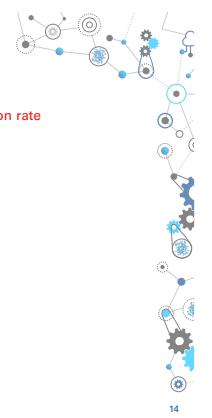
Wireless resource allocation

Problem

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



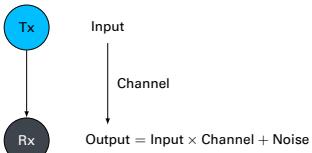
13



Wireless resource allocation

Problem

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

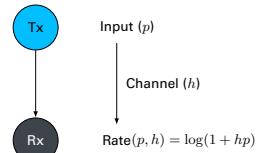


15

Wireless resource allocation

Problem

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

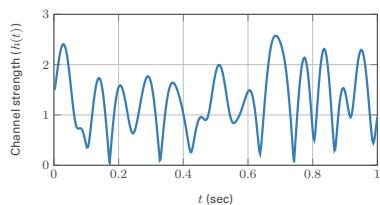


15

Wireless resource allocation

Problem

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

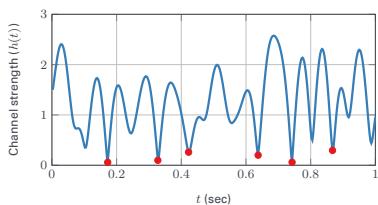


16

Wireless resource allocation

Problem

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

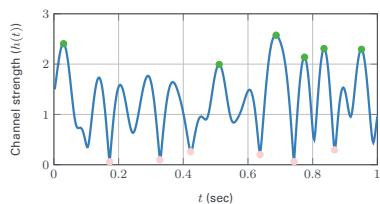


16

Wireless resource allocation

Problem

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

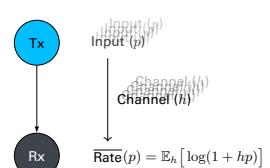


16

Wireless resource allocation

Problem

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



17

Wireless resource allocation

Problem

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



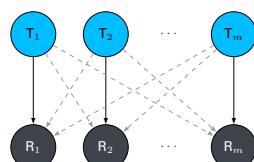
$$\begin{aligned} \min_p \quad & \mathbb{E}_h [p(h)] \\ \text{s. to} \quad & \mathbb{E}_h [\log(1 + hp)] \geq c \end{aligned}$$

17

Wireless resource allocation

Problem

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



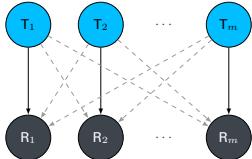
$$\begin{aligned} \min_p \quad & \sum_{i=1}^m \mathbb{E}_h [p_i(h)] \\ \text{s. to} \quad & \mathbb{E}_h \left[\log \left(1 + \frac{h^{ii} p_i(h)}{1 + \sum_{j \in N_i} h^{ij} p_j(h)} \right) \right] \geq c_i \end{aligned}$$

17

Wireless resource allocation

Problem

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



$$\begin{aligned} \min_{\mathbf{p}} \quad & \sum_{i=1}^m \mathbb{E}_{\mathbf{h}}[p_i(\mathbf{h})] \\ \text{s. to} \quad & \mathbb{E}_{\mathbf{h}} \left[\log \left(1 + \frac{\mathbf{h}^{ii} p_i(\mathbf{h})}{1 + \sum_{j \in \mathcal{N}_i} \mathbf{h}^{ij} p_j(\mathbf{h})} \right) \right] \geq c_i \end{aligned}$$

17

Wireless resource allocation

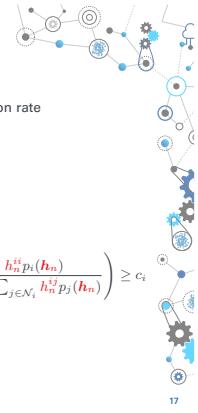
Problem

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

[Eisen, Zhang, C., Lee, and Ribeiro, IEEE TSP'19]

[Eisen, Zhang, C., Lee, and Ribeiro, IEEE TSP'19]

17



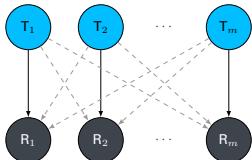
$$\begin{aligned} \min_{\mathbf{p}} \quad & \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^m p_i(\mathbf{h}_n) \\ \text{s. to} \quad & \frac{1}{N} \sum_{n=1}^N \log \left(1 + \frac{\mathbf{h}_n^{ii} p_i(\mathbf{h}_n)}{1 + \sum_{j \in \mathcal{N}_i} \mathbf{h}_n^{ij} p_j(\mathbf{h}_n)} \right) \geq c_i \end{aligned}$$

17

Wireless resource allocation

Problem

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



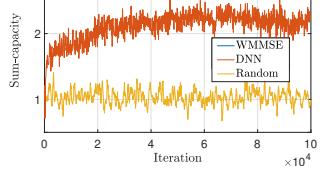
$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^m p_i(\mathbf{h}_n; \theta) \\ \text{s. to} \quad & \frac{1}{N} \sum_{n=1}^N \log \left(1 + \frac{\mathbf{h}_n^{ii} p_i(\mathbf{h}_n; \theta)}{1 + \sum_{j \in \mathcal{N}_i} \mathbf{h}_n^{ij} p_j(\mathbf{h}_n; \theta)} \right) \geq c_i \end{aligned}$$

17

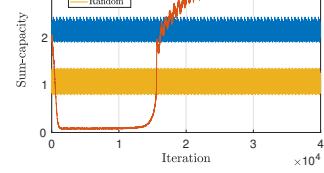
Wireless resource allocation

Problem

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



[Eisen, Zhang, C., Lee, and Ribeiro, IEEE TSP'19; Eisen and Ribeiro, IEEE TSP'20]



18

Claims

Constrained learning is the right tool to learn under requirements

Wireless resource allocation
Robust image recognition

Constrained learning is hard...

...but possible

Robust image recognition

Problem

Learn an image classifier



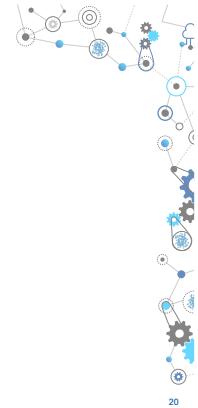
Cello

19

Robust image recognition

Problem

Learn an image classifier



20



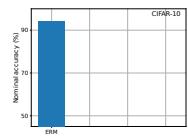
Robust image recognition

Problem

Learn an image classifier



Cello



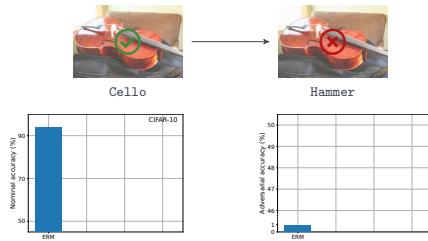
20

20

Robust image recognition

Problem

Learn an image classifier

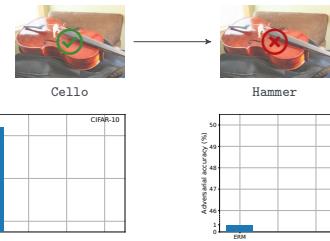


20

Robust image recognition

Problem

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



20

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

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



21

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

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

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



21

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

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

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



21

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

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

22

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

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

- Computing the worst-case perturbations

22



Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

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

22

- Computing the worst-case perturbations

▪ ≈ gradient ascent

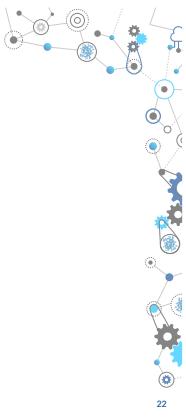
[Szegedy et al., ICLR'14; Goodfellow et al., ICLR'15; Madry et al., ICLR'18; ...]

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

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

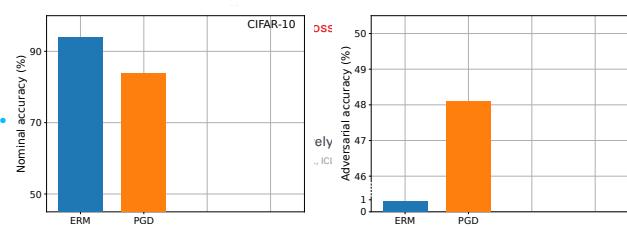


22

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations



22

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

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



22

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

- Computing the worst-case perturbations

- ✖ ≈ gradient ascent: non-convex and (severely) underparametrized
[Szegedy et al., ICLR'14; Goodfellow et al., ICLR'15; Madry et al., ICLR'18; ...]
- Balancing nominal accuracy and robustness
 - Penalty-based method (e.g., [Zhang et al., ICML'19])

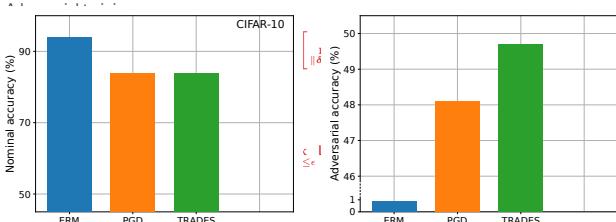
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \rightarrow \min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_\infty \leq \epsilon} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

23

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations



[Zhang et al., ICML'19]

23

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

- Adversarial training

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

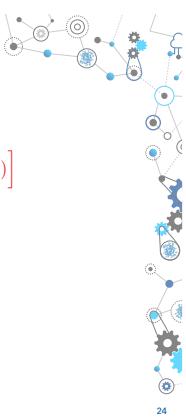
23

Penalty-based methods

Problem

Learn an image classifier that is robust to input perturbations

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



24

- No straightforward relation between λ and **adversarial loss**

Penalty-based methods

Problem

Learn an image classifier that is robust to input perturbations

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

24

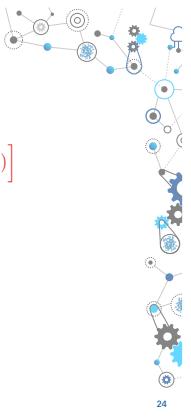
- No straightforward relation between λ and **adversarial loss**
- λ depends on the values of the losses (dataset, model, performance measure)

Penalty-based methods

Problem

Learn an image classifier that is robust to input perturbations

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



24

Adversarial training

Problem

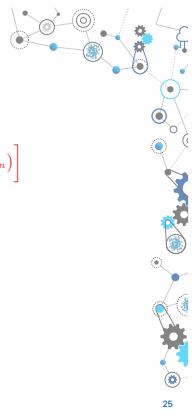
Learn an image classifier that is robust to input perturbations

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

- Computing the worst-case perturbations

✖ gradient ascent: non-convex and (severely) underparametrized
[Szegedy et al., ICLR'14; Goodfellow et al., ICLR'15; Madry et al., ICLR'18; ...]

- Balancing nominal accuracy and robustness
- ✖ Penalty-based method: value of λ and generalization



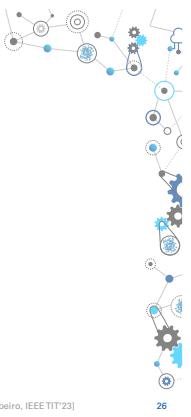
25

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

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



26

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

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



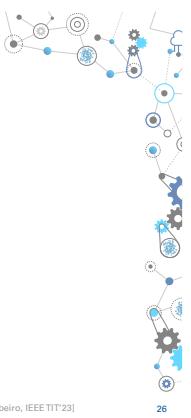
26

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

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



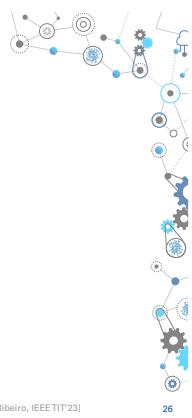
26

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

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

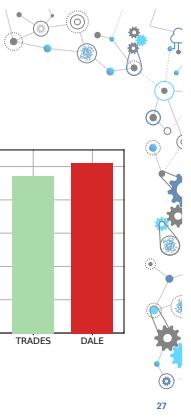
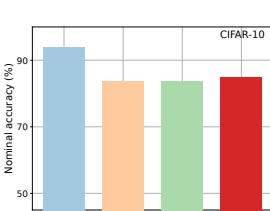


26

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

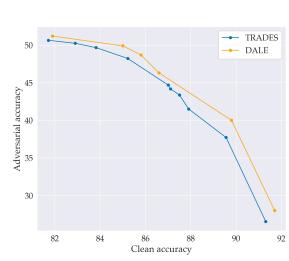


27

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations



[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]



28

Claims

Constrained learning is the right tool to learn under requirements

Constrained learning is hard...

...but possible

29

What is (un)constrained learning?

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

- ℓ, g are bounded, Lipschitz continuous (possibly non-convex) functions
- f_{θ} is a (possibly nonlinear) parametrization [e.g., logistic classifier, (G)(C)NN]
- $\mathcal{D}, \mathfrak{A}, \mathfrak{P}$ unknown

[C, Paternain, Calvo-Fullana, and Ribeiro, ICASSP'20 (best student paper); C, and Ribeiro, NeurIPS'20; C, Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23] 30

What is (un)constrained learning?

$$\begin{aligned} P^* &= \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)] \\ \text{subject to } &\mathbb{E}_{(x,y) \sim \mathfrak{A}} [g(f_{\theta}(x), y)] \leq c \\ &h(f_{\theta}(x), y) \leq u, \quad \mathfrak{P}\text{-a.e.} \end{aligned}$$

29

What is (un)constrained learning?

$$\begin{aligned} P^* &= \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)] \\ \text{subject to } &\mathbb{E}_{(x,y) \sim \mathfrak{A}} [g(f_{\theta}(x), y)] \leq c \\ &h(f_{\theta}(x), y) \leq u, \quad \mathfrak{P}\text{-a.e.} \end{aligned}$$

- ℓ, g are bounded, Lipschitz continuous (possibly non-convex) functions
- f_{θ} is a (possibly nonlinear) parametrization [e.g., logistic classifier, (G)(C)NN]
- $\mathcal{D}, \mathfrak{A}, \mathfrak{P}$ unknown

[C, Paternain, Calvo-Fullana, and Ribeiro, ICASSP'20 (best student paper); C, and Ribeiro, NeurIPS'20; C, Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23] 30

- ℓ, g are bounded, Lipschitz continuous (possibly non-convex) functions
- f_{θ} is a (possibly nonlinear) parametrization [e.g., logistic classifier, (G)(C)NN]
- $\mathcal{D}, \mathfrak{A}, \mathfrak{P}$ unknown

[C, Paternain, Calvo-Fullana, and Ribeiro, ICASSP'20 (best student paper); C, and Ribeiro, NeurIPS'20; C, Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23] 30

What is (un)constrained learning?

$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) \\ \text{subject to } &\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c \\ &h(f_{\theta}(x_r), y_r) \leq u, \quad r = 1, \dots, N \end{aligned}$$

29

Constrained learning challenges

$$\begin{array}{ccc} \hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) & & P^* = \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)] \\ \text{subject to } \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c & \xrightarrow{?} & \text{subject to } \mathbb{E}_{(x,y) \sim \mathfrak{A}} [g(f_{\theta}(x), y)] \leq c \\ & & h(f_{\theta}(x_r), y_r) \leq u \end{array}$$

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

[C, Paternain, Calvo-Fullana, and Ribeiro, ICASSP'20 (best student paper); C, and Ribeiro, NeurIPS'20; C, Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23] 30

Challenges

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

What is (un)constrained learning?

$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) \\ \text{subject to } &\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c \\ &h(f_{\theta}(x_r), y_r) \leq u, \quad r = 1, \dots, N \end{aligned}$$

29

Constrained learning challenges

$$\begin{array}{ccc} \hat{P}^* = \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) & & P^* = \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)] \\ \text{subject to } \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c & \xrightarrow{?} & \text{subject to } \mathbb{E}_{(x,y) \sim \mathfrak{A}} [g(f_{\theta}(x), y)] \leq c \\ & & h(f_{\theta}(x_r), y_r) \leq u \end{array}$$

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

[C, Paternain, Calvo-Fullana, and Ribeiro, ICASSP'20 (best student paper); C, and Ribeiro, NeurIPS'20; C, Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23] 32

Challenges

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

32

Claims

Constrained learning is the right tool to learn under requirements

Constrained learning is hard...

...but possible

34

Constrained learning challenges

$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) & P^* &= \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \\ \text{subject to} & \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c & \xrightarrow{?} & \text{subject to} \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{X}} [g(f_{\theta}(\mathbf{x}), y)] \leq c \\ & h(f_{\theta}(x_r), y_r) \leq u & & h(f_{\theta}(\mathbf{x}), y) \leq u \text{ a.e.} \end{aligned}$$

Challenges

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

35

Constrained learning challenges

$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) & P^* &= \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \\ \text{subject to} & \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c & \xrightarrow{?} & \text{subject to} \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{X}} [g(f_{\theta}(\mathbf{x}), y)] \leq c \\ & h(f_{\theta}(x_r), y_r) \leq u & & h(f_{\theta}(\mathbf{x}), y) \leq u \text{ a.e.} \end{aligned}$$

35

Constrained learning challenges

$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) & P^* &= \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \\ \text{subject to} & \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c & \xrightarrow{?} & \text{subject to} \quad \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{X}} [g(f_{\theta}(\mathbf{x}), y)] \leq c \\ & h(f_{\theta}(x_r), y_r) \leq u & & h(f_{\theta}(\mathbf{x}), y) \leq u \text{ a.e.} \end{aligned}$$

Challenges

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

35

Challenges

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

A PAC learning primer

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

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

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

36

A PAC learning primer

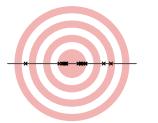
- Classical learning theory [Vapnik & Chervonenkis, TP'71; Valiant, CACM'84]:

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

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

36

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



36

A PAC learning primer

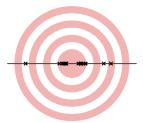
- Classical learning theory [Vapnik & Chervonenkis, TP'71; Valiant, CACM'84]:

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

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

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

- Tools: Rademacher complexity, VC dimensionality...



36

A PAC learning primer

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

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

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

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

• Tools: Rademacher complexity, VC dimensionality...

- Requirements?

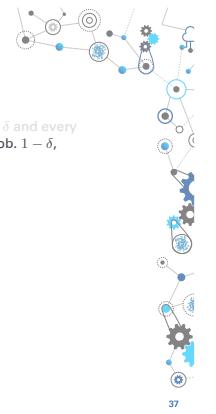


36

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{X} , we can obtain f_{θ^*} from $N_f(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,



37

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{X} , we can obtain f_{θ^*} from $N_f(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,

- 1) near-optimal

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

[C. and Ribeiro, NeurIPS'20; C. Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

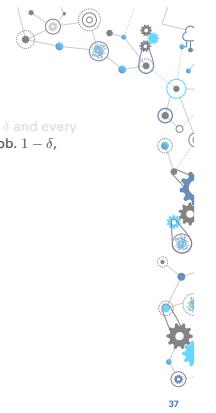


37

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{X} , we can obtain f_{θ^*} from $N_f(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,



37

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{X} , we can obtain f_{θ^*} from $N_f(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,

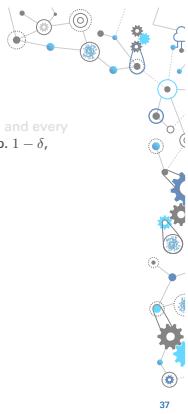
- 1) near-optimal

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

- 2) approximately feasible

$$\mathbb{E}_{(x,y) \sim \mathcal{X}} [g(f_{\theta^*}(x), y)] \leq c + \epsilon$$

[C. and Ribeiro, NeurIPS'20; C. Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

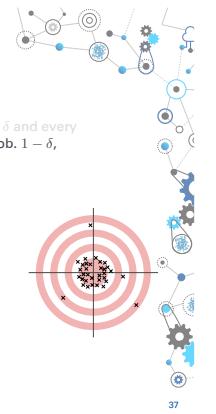


37

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{X} , we can obtain f_{θ^*} from $N_f(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,



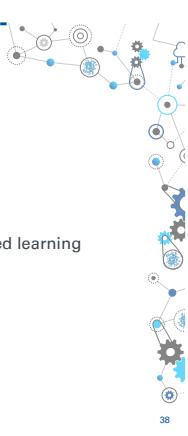
37

When is it possible to learn under constraints?

Theorem

f_{θ} is PAC learnable \iff f_{θ} is PACC learnable

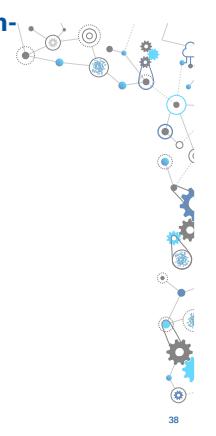
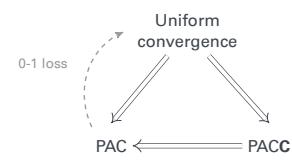
Constrained learning is essentially as hard as unconstrained learning



38

When is it possible to learn under constraints?

Theorem



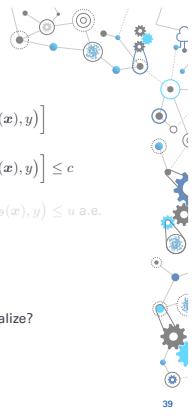
38

[C. and Ribeiro, NeurIPS'20; C. Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

Constrained learning challenges

$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) \\ \text{subject to } & \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c \end{aligned}$$

PAC(C)



Challenges

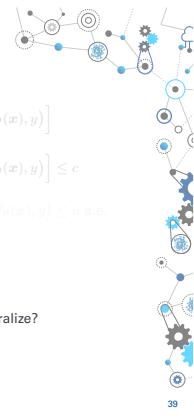
- 1) *Statistical:* does the solution of the constrained empirical problem generalize?
- 2) *Computational:* can we solve the constrained empirical problem?

39

Constrained learning challenges

$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) \\ \text{subject to } & \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c \end{aligned}$$

PAC(C)



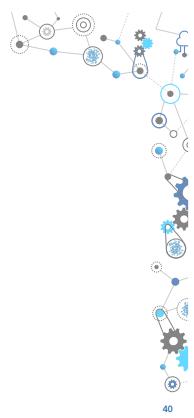
Challenges

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

39

Constrained optimization methods

$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) \\ \text{subject to } & \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c \\ & h(f_{\theta}(x_r), y_r) \leq u \end{aligned}$$



40

Constrained optimization methods

- Feasible update methods
e.g., conditional gradients (Frank-Wolfe)
- Interior point methods
e.g., barriers, projection, polyhedral approx.

$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) \\ \text{subject to } & \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c \\ & h(f_{\theta}(x_r), y_r) \leq u \end{aligned}$$



40

Constrained optimization methods

$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) \\ \text{subject to } & \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c \\ & h(f_{\theta}(x_r), y_r) \leq u \end{aligned}$$

- Feasible update methods
e.g., conditional gradients (Frank-Wolfe)
 - ✖ Tractability [non-convex constraints]
 - ✓ Feasible candidate solution
- Interior point methods
e.g., barriers, projection, polyhedral approx.
 - ✖ Tractability [non-convex constraints]
 - ✓ Feasible candidate solution

40

Constrained optimization methods

- Feasible update methods
e.g., conditional gradients (Frank-Wolfe)
 - ✖ Tractability [non-convex constraints]
 - ✓ Feasible candidate solution
- Interior point methods
e.g., barriers, projection, polyhedral approx.
 - ✖ Tractability [non-convex constraints]
 - ✓ Feasible candidate solution
- Duality
e.g., (augmented) Lagrangian
 - ✓ Tractability
 - ✖ Feasible candidate solution [strong-duality]

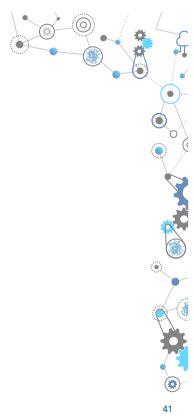
$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) \\ \text{subject to } & \frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c \\ & h(f_{\theta}(x_r), y_r) \leq u \end{aligned}$$



40

Duality

PRIMAL
↔
DUAL



41

Duality

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

↔

DUAL



41

Duality

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

↑

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

Duality

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

↑

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

- Unconstrained optimization is “easier” than constrained optimization

41

41

Duality

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

↑

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

Duality

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

↑

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

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

41

41

- Unconstrained optimization is “easier” than constrained optimization

- In general, $\hat{D}^* \leq \hat{P}^*$

Break Robust learning revisited

Adversarial training: Challenges

Problem

Learn an image classifier that is robust to input perturbations

$$\begin{aligned} \min_{\theta} \quad & \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\|\delta\|_\infty \leq \epsilon} \text{Loss}(f_{\theta}(x + \delta), y) \right] \\ \text{subject to} \quad & \mathbb{E}_{(x,y) \sim \mathcal{D}} [\text{Loss}(f_{\theta}(x), y)] \leq c \end{aligned}$$

- Computing the worst-case perturbations
✖ ≈ gradient ascent: non-convex and (severely) underparametrized
[Szegedy et al., ICLR’14; Goodfellow et al., ICLR’15; Madry et al., ICLR’18; ...]
- Balancing nominal accuracy and robustness
✓ Constrained learning [GR, NeurIPS’20; FG*APH, NeurIPS’21; CPCR, IEEETIT’23]

43

Adversarial training: Challenges

Problem

Learn an image classifier that is robust to input perturbations

$$\begin{aligned} \min_{\theta} \quad & \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\|\delta\|_\infty \leq \epsilon} \text{Loss}(f_{\theta}(x + \delta), y) \right] \\ \text{subject to} \quad & \mathbb{E}_{(x,y) \sim \mathcal{D}} [\text{Loss}(f_{\theta}(x), y)] \leq c \end{aligned}$$

- Computing the worst-case perturbations
✖ ≈ gradient ascent: non-convex and (severely) underparametrized
[Szegedy et al., ICLR’14; Goodfellow et al., ICLR’15; Madry et al., ICLR’18; ...]
- Balancing nominal accuracy and robustness
✓ Constrained learning [GR, NeurIPS’20; FG*APH, NeurIPS’21; CPCR, IEEETIT’23]

43



Semi-infinite constrained learning

Problem

Learn an image classifier that is robust to input perturbations

$$\begin{aligned} \min_{\theta} \quad & \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\|\delta\|_\infty \leq \epsilon} \text{Loss}(f_{\theta}(x + \delta), y) \right] \\ \text{subject to} \quad & \mathbb{E}_{(x,y) \sim \mathcal{D}} [\text{Loss}(f_{\theta}(x), y)] \leq c \end{aligned}$$



44

Semi-infinite constrained learning

Problem

Learn an image classifier that is robust to input perturbations

$$\begin{aligned} \min_{\theta} & \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [l(\mathbf{x}, y)] \\ \text{subject to} & \text{Loss}(f_{\theta}(\mathbf{x} + \delta), y) \leq t(\mathbf{x}, y), \\ & \text{for a.e. } (\mathbf{x}, y) \text{ and } \|\delta\|_{\infty} \leq \epsilon \end{aligned}$$



44

- Epigraph formulation:

$$\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(\mathbf{x} + \delta), y) \leq t \iff \text{Loss}(f_{\theta}(\mathbf{x} + \delta), y) \leq t, \text{ for all } \|\delta\|_{\infty} \leq \epsilon$$

Semi-infinite constrained learning

Problem

Learn an image classifier that is robust to input perturbations

$$\begin{aligned} \min_{\theta} & \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [l(\mathbf{x}, y)] \\ \text{subject to} & \text{Loss}(f_{\theta}(\mathbf{x} + \delta_0), y) \leq t(\mathbf{x}, y) \\ & \text{Loss}(f_{\theta}(\mathbf{x} + \delta_{\sqrt{\epsilon}}), y) \leq t(\mathbf{x}, y) \\ & \text{Loss}(f_{\theta}(\mathbf{x} + \delta_{\epsilon}), y) \leq t(\mathbf{x}, y) \\ & \text{Loss}(f_{\theta}(\mathbf{x} + \delta_{\pi}), y) \leq t(\mathbf{x}, y) \\ & \text{Loss}(f_{\theta}(\mathbf{x} + \delta_4), y) \leq t(\mathbf{x}, y) \\ \max_{\|\delta\|_{\infty} \leq \epsilon} & \text{Loss}(f_{\theta}(\mathbf{x} + \delta), y) \leq t, \text{ for all } \|\delta\|_{\infty} \leq \epsilon \\ \text{Semi-infinite program} & \text{Loss}(f_{\theta}(\mathbf{x} + \delta_{\pi^*}), y) \leq t(\mathbf{x}, y) \\ & \text{Loss}(f_{\theta}(\mathbf{x} + \delta_{\pi^*}), y) \leq t(\mathbf{x}, y) \\ & \text{Loss}(f_{\theta}(\mathbf{x} + \delta_{25}), y) \leq t(\mathbf{x}, y) \end{aligned}$$



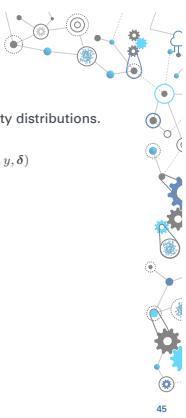
44

Semi-infinite constrained learning

Proposition

Let $(\mathbf{x}, y) \mapsto \ell(f_{\theta}(\mathbf{x}, y)) \in L^2$ and \mathcal{P}^2 be the space of square integrable probability distributions. Then,

$$\begin{aligned} P^* &= \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(\mathbf{x}, y)] \text{ subject to } \ell(f_{\theta}(\mathbf{x} + \delta), y) \leq t(\mathbf{x}, y), \forall (\mathbf{x}, y, \delta) \\ D^* &= \min_{\theta} \max_{\mu \in \mathcal{P}^2} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\int \mu(\mathbf{x}, y, \delta) \ell(f_{\theta}(\mathbf{x} + \delta), y) d\delta \right] \end{aligned}$$



45

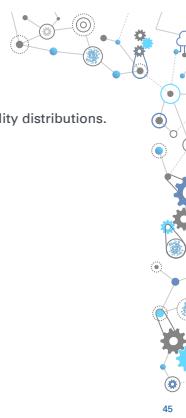
Semi-infinite constrained learning

Proposition

Let $(\mathbf{x}, y) \mapsto \ell(f_{\theta}(\mathbf{x}, y)) \in L^2$ and \mathcal{P}^2 be the space of square integrable probability distributions. Then,

$$\begin{aligned} P^* &= \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \ell(f_{\theta}(\mathbf{x} + \delta), y) \right] \\ D^* &= \min_{\theta} \max_{\mu \in \mathcal{P}^2} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\int \mu(\mathbf{x}, y, \delta) \ell(f_{\theta}(\mathbf{x} + \delta), y) d\delta \right] \end{aligned}$$

[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]



45

[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

Semi-infinite constrained learning

Proposition

Let $(\mathbf{x}, y) \mapsto \ell(f_{\theta}(\mathbf{x}, y)) \in L^2$ and \mathcal{P}^2 be the space of square integrable probability distributions. Then,

$$\begin{aligned} P^* &= \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \ell(f_{\theta}(\mathbf{x} + \delta), y) \right] \\ D^* &= \min_{\theta} \max_{\mu \in \mathcal{P}^2} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\mathbb{E}_{\delta \sim \mu(\delta | \mathbf{x}, y)} [\ell(f_{\theta}(\mathbf{x} + \delta), y)] \right] \end{aligned}$$



45

- $\mu(\mathbf{x}, y)$ is a conditional probability (Radon-Nikodym derivative) of perturbations

[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

Semi-infinite constrained learning

Proposition

Let $(\mathbf{x}, y) \mapsto \ell(f_{\theta}(\mathbf{x}, y)) \in L^2$ and \mathcal{P}^2 be the space of square integrable probability distributions. Then,

$$\begin{aligned} P^* &= \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \ell(f_{\theta}(\mathbf{x} + \delta), y) \right] \\ D^* &= \min_{\theta} \max_{\mu \in \mathcal{P}^2} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\mathbb{E}_{\delta \sim \mu(\delta | \mathbf{x}, y)} [\ell(f_{\theta}(\mathbf{x} + \delta), y)] \right] \end{aligned}$$

Proposition

For all $\xi > 0$, there exists $\gamma(\mathbf{x}, y) < \max_{\|\delta\|_{\infty} \leq \epsilon} \ell(f_{\theta}(\mathbf{x} + \delta), y)$ s.t. $L(\theta, \mu_{\gamma}) \geq \max_{\mu \in \mathcal{P}^2} L(\theta, \mu) - \xi$ for

$$\mu_{\gamma}(\delta | \mathbf{x}, y) \propto [\ell(f_{\theta}(\mathbf{x} + \delta), y) - \gamma(\mathbf{x}, y)]_+$$

[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]



45

Semi-infinite constrained learning

Proposition

Let $(\mathbf{x}, y) \mapsto \ell(f_{\theta}(\mathbf{x}, y)) \in L^2$ and \mathcal{P}^2 be the space of square integrable probability distributions. Then,

$$\begin{aligned} P^* &= \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \ell(f_{\theta}(\mathbf{x} + \delta), y) \right] \\ D^* &= \min_{\theta} \max_{\mu \in \mathcal{P}^2} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\mathbb{E}_{\delta \sim \mu_{\gamma}(\delta | \mathbf{x}, y)} [\ell(f_{\theta}(\mathbf{x} + \delta), y)] \right] \end{aligned}$$



45

Proposition

For all $\xi > 0$, there exists $\gamma(\mathbf{x}, y) < \max_{\|\delta\|_{\infty} \leq \epsilon} \ell(f_{\theta}(\mathbf{x} + \delta), y)$ s.t. $L(\theta, \mu_{\gamma}) \geq \max_{\mu \in \mathcal{P}^2} L(\theta, \mu) - \xi$ for

$$\mu_{\gamma}(\delta | \mathbf{x}, y) \propto [\ell(f_{\theta}(\mathbf{x} + \delta), y) - \gamma(\mathbf{x}, y)]_+$$

[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

Dual Adversarial Learning

$$\begin{aligned} & \delta' \sim \mu_{\gamma}(\cdot | \mathbf{x}, y) \\ & \theta^+ = \theta - \eta \nabla_{\theta} \ell(f_{\theta}(\mathbf{x} + \delta'), y) \end{aligned}$$

$$\mu_{\gamma}(\delta | \mathbf{x}, y) \propto [\ell(f_{\theta}(\mathbf{x} + \delta), y) - \gamma(\mathbf{x}, y)]_+$$

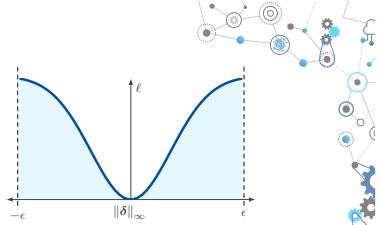


46

Dual Adversarial Learning

$$\begin{cases} \delta' \sim \mu_\gamma(\cdot|x, y) \\ \theta^+ = \theta - \eta \nabla_\theta \ell(f_\theta(x + \delta'), y) \end{cases}$$

$$\mu_\gamma(\delta|x, y) \propto [\ell(f_\theta(x + \delta), y) - \gamma(x, y)]_+$$



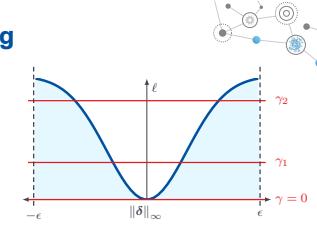
- γ control the smoothness of the approximation

46

Dual Adversarial Learning

$$\begin{cases} \delta' \sim \mu_\gamma(\cdot|x, y) \\ \theta^+ = \theta - \eta \nabla_\theta \ell(f_\theta(x + \delta'), y) \end{cases}$$

$$\mu_\gamma(\delta|x, y) \propto [\ell(f_\theta(x + \delta), y) - \gamma(x, y)]_+$$



- γ control the smoothness of the approximation

$\gamma = 0$
(oversmoothed)

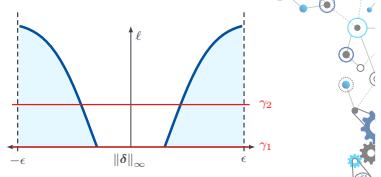
[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

46

Dual Adversarial Learning

$$\begin{cases} \delta' \sim \mu_\gamma(\cdot|x, y) \\ \theta^+ = \theta - \eta \nabla_\theta \ell(f_\theta(x + \delta'), y) \end{cases}$$

$$\mu_\gamma(\delta|x, y) \propto [\ell(f_\theta(x + \delta), y) - \gamma(x, y)]_+$$



- γ control the smoothness of the approximation

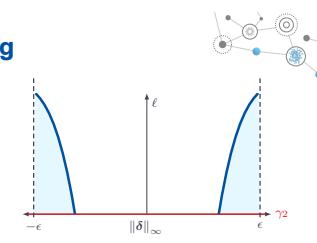
$\gamma = 0$
(oversmoothed)

46

Dual Adversarial Learning

$$\begin{cases} \delta' \sim \mu_\gamma(\cdot|x, y) \\ \theta^+ = \theta - \eta \nabla_\theta \ell(f_\theta(x + \delta'), y) \end{cases}$$

$$\mu_\gamma(\delta|x, y) \propto [\ell(f_\theta(x + \delta), y) - \gamma(x, y)]_+$$



- γ control the smoothness of the approximation

$\gamma = 0$
(oversmoothed)

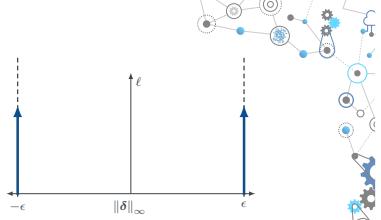
[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

46

Dual Adversarial Learning

$$\begin{cases} \delta' \sim \mu_\gamma(\cdot|x, y) \\ \theta^+ = \theta - \eta \nabla_\theta \ell(f_\theta(x + \delta'), y) \end{cases}$$

$$\mu_\gamma(\delta|x, y) \propto [\ell(f_\theta(x + \delta), y) - \gamma(x, y)]_+$$



- γ control the smoothness of the approximation

$\gamma = 0$
(oversmoothed)

$\gamma \rightarrow \max_\delta \ell(f_\theta(x + \delta), y)$
(undersmoothed)

✖ hard to sample

46

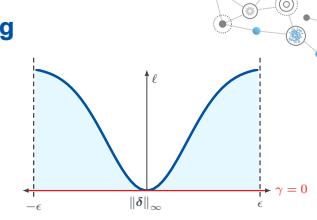
[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

46

Dual Adversarial Learning

$$\begin{cases} \delta' \sim \mu_\gamma(\cdot|x, y) \\ \theta^+ = \theta - \eta \nabla_\theta \ell(f_\theta(x + \delta'), y) \end{cases}$$

$$\mu_\gamma(\delta|x, y) \propto [\ell(f_\theta(x + \delta), y) - \gamma(x, y)]_+$$



- γ control the smoothness of the approximation

$\gamma = 0$
(oversmoothed)

$\gamma \rightarrow \max_\delta \ell(f_\theta(x + \delta), y)$
(undersmoothed)

✖ easy to sample

[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

46

Dual Adversarial Learning

$$\begin{cases} \delta' \sim \ell(f_\theta(x + \delta), y) \\ \theta^+ = \theta - \eta \nabla_\theta \ell(f_\theta(x + \delta'), y) \end{cases}$$

- $\mu_0 \in L^2$ and μ_0 as "differentiable" as $\ell \circ f_\theta \Rightarrow$ Langevin/Hamiltonian Monte Carlo

$$\delta^+ = \operatorname{proj}_{\|\delta\|_\infty \leq \epsilon} \left[\delta + \eta \operatorname{sign} [\nabla_\delta \log (\ell(f_\theta(x + \delta), y))] + \sqrt{2\eta T} \zeta \right], \quad \zeta \sim \text{Laplace}(0, I)$$



47

[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

47

Dual Adversarial Learning

$$\begin{cases} \delta' \sim \ell(f_\theta(x + \delta), y) \\ \theta^+ = \theta - \eta \nabla_\theta \ell(f_\theta(x + \delta'), y) \end{cases}$$

- $\mu_0 \in L^2$ and μ_0 as "differentiable" as $\ell \circ f_\theta \Rightarrow$ Langevin/Hamiltonian Monte Carlo

$$\delta^+ = \operatorname{proj}_{\|\delta\|_\infty \leq \epsilon} \left[\delta + \eta \operatorname{sign} [\nabla_\delta \log (\ell(f_\theta(x + \delta), y))] + \sqrt{2\eta T} \zeta \right], \quad \zeta \sim \text{Laplace}(0, I)$$

- If $\ell \circ f_\theta$ is not differentiable \Rightarrow Metropolis-Hastings

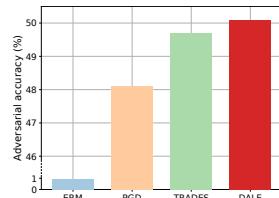
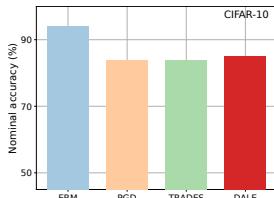
[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

47

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations



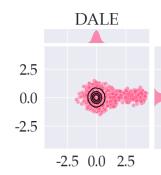
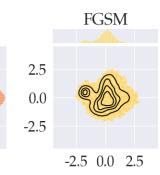
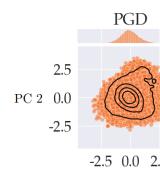
[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

48

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations



[Robey*, C., Pappas, Hassani, and Ribeiro, NeurIPS'21]

49

Constrained learning challenges

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

Challenges

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

50

Duality

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

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

51

Duality

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

51

- Unconstrained optimization is “easier” than constrained optimization
- 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 \longleftrightarrow Dual

52

Non-convex variational duality

Convex optimization: Primal \longleftrightarrow Dual

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

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

52

Sparse logistic regression

$$\begin{aligned} \min_{\theta \in \mathbb{R}^p} & - \sum_{n=1}^N \log [1 + \exp(y_n \cdot \theta^T \mathbf{x}_n)] \\ \text{s.t. } & \|\theta\|_0 = \sum_{t=1}^p \mathbb{I}[\theta_t \neq 0] \leq k \end{aligned}$$

Discrete, non-convex

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

53

Sparse logistic regression

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

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

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

s. to $\|\theta\|_1 = \sum_{t=1}^p |\theta_t| \leq c$

Discrete, convex

53

Sparse logistic regression

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

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

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

s. to $\|\theta\|_1 = \sum_{t=1}^p |\theta_t| \leq c$

Discrete, convex

53

Sparse logistic regression

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

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

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

s. to $\|\theta\|_1 = \sum_{t=1}^p |\theta_t| \leq c$

Discrete, convex
[Foucart & Rauhut, 2013]

53

Sparse logistic regression

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

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

54

Sparse logistic regression

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

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

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

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

Continuous, non-convex
[C., Eldar, and Ribeiro, IEEE TSP'20]: tractable

54

Sparse logistic regression

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

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

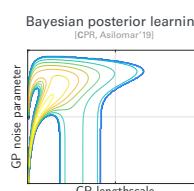
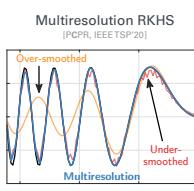
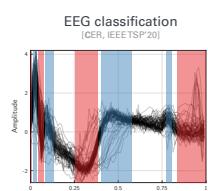
$$\min_{\theta \in L_2} - \sum_{n=1}^N \log \left[1 + \exp \left(y_n \cdot \int \theta(t) x_n(t) dt \right) \right]$$

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

Continuous, non-convex
[C., Eldar, and Ribeiro, IEEE TSP'20]: tractable

54

Non-convex variational duality



55

How can we learn under constraints?

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

↑ ↓

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

56

How can we learn under constraints?

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

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

$$P^* = \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), y)] \text{ subject to } \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{A}} [g(f_{\theta}(\mathbf{x}), y)] \leq c$$

56

How can we learn under constraints?

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

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

$$P^* = \min_{\theta} \int \ell(f_{\theta}(\mathbf{x}), y) d\mathcal{D}(\mathbf{x}, y) \text{ subject to } \int g(f_{\theta}(\mathbf{x}), y) d\mathcal{A}(\mathbf{x}, y) \leq c$$

[C. and Ribeiro, NeurIPS'20; C., Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23]

56

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

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

[\mathcal{H} is a good covering of $\text{conv}(\mathcal{H})$]

56

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

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

Then \hat{D}^* is a (near-)PACC learner, i.e., there exists a solution θ^\dagger that, with probability $1 - \delta$,

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

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

(mild conditions apply)

[C., Paternain, Calvo-Fullana, and Ribeiro, ICASSP'20 (best student paper); C. and Ribeiro, NeurIPS'20; C., Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23] 57

57

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

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

Then \hat{D}^* is a (near-)PACC learner, i.e., there exists a solution θ^\dagger that, with probability $1 - \delta$,

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

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

(if losses are convex) $\mu(f_{\theta^\dagger}(\mathbf{x}), y) \leq r$, with \mathbb{P} -prob. $1 - \tilde{O} \left(\frac{1}{\sqrt{N}} \right)$

(mild conditions apply)

[C., Paternain, Calvo-Fullana, and Ribeiro, ICASSP'20 (best student paper); C. and Ribeiro, NeurIPS'20; C., Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23] 57

57

Dual learning

Theorem

Let f be ν -universal with VC dimension $d_{VC} < \infty$. There exists $(\theta^\dagger, \lambda^\dagger)$ achieving \hat{D}^* such that f_{θ^\dagger} is a (near-)PACC solution of (P-CSL), i.e., with probability at least $1 - \delta$,

$$|P^* - \hat{D}^*| \leq (1 + \Delta)(\epsilon_0 + \epsilon)$$

$$\mathbb{E} [g(f_{\theta^\dagger}(\mathbf{x}), y)] \leq c + \epsilon$$

$\epsilon_0 = M\nu$

$\Delta = \max \left(\|\lambda^*\|_1, \|\hat{\lambda}^*\|_1, \|\tilde{\lambda}^*\|_1 \right)$

$\epsilon = B \sqrt{\frac{1}{N} \left[1 + \log \left(\frac{4m(2N)^{d_{VC}}}{\delta} \right) \right]}$

Source of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

[C., Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23]

58

Dual learning

Theorem

Let f be ν -universal with VC dimension $d_{VC} < \infty$. There exists $(\theta^\dagger, \lambda^\dagger)$ achieving \hat{D}^* such that f_{θ^\dagger} is a (near-)PACC solution of (P-CSL), i.e., with probability at least $1 - \delta$,

$$|P^* - \hat{D}^*| \leq (1 + \Delta)(\epsilon_0 + \epsilon)$$

$$\mathbb{E} [g(f_{\theta^\dagger}(\mathbf{x}), y)] \leq c + \epsilon$$

$\epsilon_0 = M\nu$

$\Delta = \max \left(\|\lambda^*\|_1, \|\hat{\lambda}^*\|_1, \|\tilde{\lambda}^*\|_1 \right)$

$\epsilon = B \sqrt{\frac{1}{N} \left[1 + \log \left(\frac{4m(2N)^{d_{VC}}}{\delta} \right) \right]}$

Source of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

[C., Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23]

58

Dual learning

Theorem

Let f be ν -universal with VC dimension $d_{VC} < \infty$. There exists $(\theta^\dagger, \lambda^\dagger)$ achieving \hat{D}^* such that f_{θ^\dagger} is a (near-)PACC solution of (P-CSL), i.e., with probability at least $1 - \delta$,

$$|P^* - \hat{D}^*| \leq (1 + \Delta)(\epsilon_0 + \epsilon)$$

$$\mathbb{E}[g(f_{\theta^\dagger}(\mathbf{x}), y)] \leq c + \epsilon$$

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

Source of error

parametrization richness (ν) sample size (N)

requirements difficulty (λ^*)

58

Dual learning

Theorem

Let f be ν -universal with VC dimension $d_{VC} < \infty$. There exists $(\theta^\dagger, \lambda^\dagger)$ achieving \hat{D}^* such that f_{θ^\dagger} is a (near-)PACC solution of (P-CSL), i.e., with probability at least $1 - \delta$,

$$|P^* - \hat{D}^*| \leq (1 + \Delta)(\epsilon_0 + \epsilon)$$

$$\mathbb{E}[g(f_{\theta^\dagger}(\mathbf{x}), y)] \leq c + \epsilon$$

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

Source of error

parametrization richness (ν) sample size (N)

requirements difficulty (λ^*)

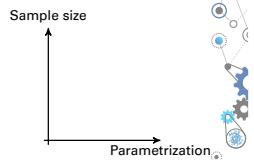
[C, Paternain, Calvo-Fullana, and Ribeiro, IEEEITIT'23]

[C, Paternain, Calvo-Fullana, and Ribeiro, IEEEITIT'23]

58

Dual learning trade-offs

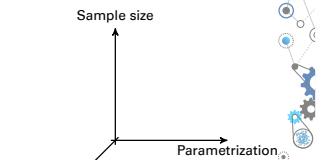
- Unconstrained learning
parametrization \times sample size



59

Dual learning trade-offs

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



[C, Paternain, Calvo-Fullana, and Ribeiro, ICASSP'20 (best student paper); C, and Ribeiro, NeurIPS'20; C, Paternain, Calvo-Fullana, and Ribeiro, IEEEITIT'23] 59

Claims

Constrained learning is the right tool to learn under requirements

Constrained learning is hard...

...but possible. How?

60

**Break
Learning
with invariance**

Learning with invariance

Problem

Learn an image classifier that is invariant to transformation $g \in \mathcal{G}$



62



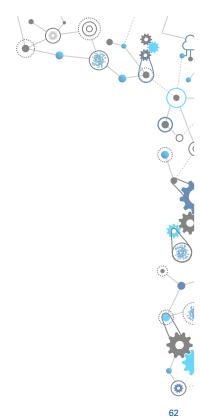
Learning with invariance

Problem

Learn an image classifier that is **invariant to transformation $g \in \mathcal{G}$**



62



Learning with invariance

Problem

Learn an image classifier that is **invariant to transformation $g \in \mathcal{G}$**



- Why?

- Improve accuracy
- Reduce sample complexity
- Robustness



62

Learning with invariance

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

[Musil et al., Chem. Rev'21; Lu et al., SIAM JSC'21; Jumper et al., Nature'21...]

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

[Cohen et al., ICML'16.; Anderson et al., NeurIPS'19; Finzi et al., ICML'20; Li et al., ICLR'21...]

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

[Lageris et al., IEEE TNN'98; Raisi et al., JCP'19...]



63

Learning with invariance

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

[Musil et al., Chem. Rev'21; Lu et al., SIAM JSC'21; Jumper et al., Nature'21...]

✖ Manual design \Rightarrow Limited \mathcal{G}

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

[Cohen et al., ICML'16.; Anderson et al., NeurIPS'19; Finzi et al., ICML'20; Li et al., ICLR'21...]

- ✖ Manual design \Rightarrow Limited \mathcal{G}
- ✖ Computational complexity

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

[Lageris et al., IEEE TNN'98; Raisi et al., JCP'19...]

- ✖ Time-consuming hyperparameter tuning
- ✖ Guarantee?



63

Learning with invariance

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

[Musil et al., Chem. Rev'21; Lu et al., SIAM JSC'21; Jumper et al., Nature'21...]

✖ Manual design \Rightarrow Limited \mathcal{G}

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

[Cohen et al., ICML'16.; Anderson et al., NeurIPS'19; Finzi et al., ICML'20; Li et al., ICLR'21...]

- ✖ Manual design \Rightarrow Limited \mathcal{G}
- ✖ Computational complexity

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

[Lageris et al., IEEE TNN'98; Raisi et al., JCP'19...]

- ✖ Time-consuming hyperparameter tuning
- ✖ Guarantee?



63

Data augmentation

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

Data: $\{(\mathbf{x}_n, y_n)\}_{n=1, \dots, N}$



64

Data augmentation

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

Data: $\{(\mathbf{x}_n, y_n)\}_{n=1, \dots, N}$

$$\min_{\theta} \hat{R}_{\text{aug}}(\theta) \triangleq \frac{1}{N} \sum_{n=1}^N \left[\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \ell(f_{\theta}(g\mathbf{x}_n), y_n) \right]$$

Data: $\bigcup_{g \in \mathcal{G}} \{(g\mathbf{x}_n, y_n)\}_{n=1, \dots, N}$

Data augmentation

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

Data: $\{(\mathbf{x}_n, y_n)\}_{n=1, \dots, N}$

Data augmentation

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

Data: $\{(\mathbf{x}_n, y_n)\}_{n=1, \dots, N}$

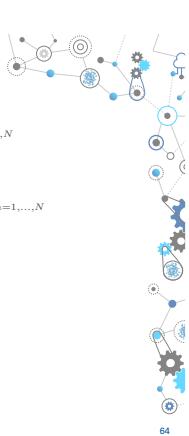
$$\min_{\theta} \hat{R}_{\text{aug}}(\theta) \triangleq \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{G}} [\ell(f_{\theta}(g\mathbf{x}_n), y_n)]$$

Data: $\bigcup_{g \in \mathcal{G}} \{(g\mathbf{x}_n, y_n)\}_{n=1, \dots, N}$

- \mathcal{G} : set of transformations ($\text{id} \in \mathcal{G}$)
- \mathcal{G} : distribution over \mathcal{G} (typically, uniform)

- \mathcal{G} : set of transformations ($\text{id} \in \mathcal{G}$)
- \mathcal{G} : distribution over \mathcal{G} (typically, uniform)

64



64

Data augmentation

$$\min_{\theta} \hat{R}(\theta) \triangleq \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n)$$

$$\min_{\theta} \hat{R}_{\text{aug}}(\theta) \triangleq \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(gx_n), y_n)]$$

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

- \mathcal{G} : set of transformations ($\text{id} \in \mathcal{G}$)
- \mathcal{G} : distribution over \mathcal{G} (typically, uniform)

Data: $\{(x_n, y_n)\}_{n=1, \dots, N}$

Data: $\bigcup_{g \in \mathcal{G}} \{(gx_n, y_n)\}_{n=1, \dots, N}$

- ✓ If...
 - \mathcal{D} is invariant to $g \in \mathcal{G}$
 - \mathcal{G} induces invariance to \mathcal{G}
 - ⇒ reduce sample complexity
- ✗ If \mathcal{G} is chosen poorly ⇒ $\hat{R}_{\text{aug}} \not\approx R$

64

[Chen et al., NeurIPS'20; Bietti et al., NeurIPS'21; Shao et al., NeurIPS'22; Hounie, Chamon, Ribeiro, ICML'23]

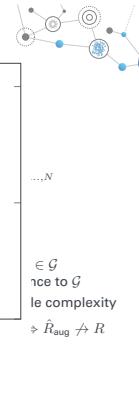
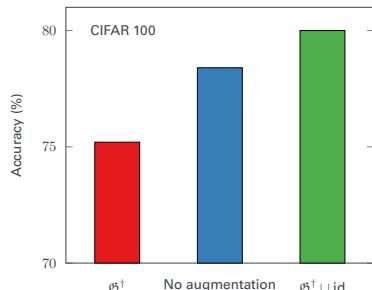
Data augmentation

$$\min_{\theta} \hat{R}(\theta) \triangleq \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n)$$

$$\min_{\theta} \hat{R}_{\text{aug}}(\theta) \triangleq \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(gx_n), y_n)$$

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

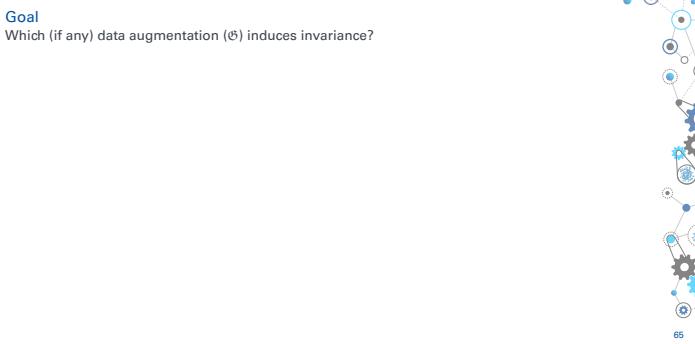
- \mathcal{G} : set of transformations
- \mathcal{G} : distribution over \mathcal{G} (typically, uniform)



From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?



[Hounie, Chamon, Ribeiro, ICML'23]

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance → Robustness → Data Augmentation

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance → Robustness → Data Augmentation

- Which invariance?
 - $f_{\theta}(gx) \stackrel{d}{=} f_{\theta}(x)$, for $x \sim \mathcal{D}$
 - $f_{\theta}(gx) = f_{\theta}(x)$, for all $g \in \mathcal{G}$
 - $\ell(f_{\theta}(gx), y) = \ell(f_{\theta}(x), y)$, for all $g \in \mathcal{G}$

[Hounie, Chamon, Ribeiro, ICML'23]

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance → Robustness → Data Augmentation

- Which invariance?
 - $f_{\theta}(gx) \stackrel{d}{=} f_{\theta}(x)$, for $x \sim \mathcal{D}$
 - $f_{\theta}(gx) = f_{\theta}(x)$, for all $g \in \mathcal{G}$
 - $\ell(f_{\theta}(gx), y) = \ell(f_{\theta}(x), y)$, for all $g \in \mathcal{G}$

[Hounie, Chamon, Ribeiro, ICML'23]

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance → Robustness → Data Augmentation

- Which invariance?
 - $f_{\theta}(gx) \stackrel{d}{=} f_{\theta}(x)$, for $x \sim \mathcal{D}$
 - $f_{\theta}(gx) = f_{\theta}(x)$, for all $g \in \mathcal{G}$
 - $\ell(f_{\theta}(gx), y) = \ell(f_{\theta}(x), y)$, for all $g \in \mathcal{G}$

[Hounie, Chamon, Ribeiro, ICML'23]

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance → Robustness → Data Augmentation

- Which invariance?
 - $\ell(f_{\theta}(gx), y) = \ell(f_{\theta}(x), y)$, for all $g \in \mathcal{G}$ ($\Rightarrow \hat{R}_{\text{aug}} \approx R$)
- What if f_{θ} cannot express invariance to \mathcal{G} ?

[Hounie, Chamon, Ribeiro, ICML'23]



From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

- Which invariance?
✓ $\ell(f_\theta(gx), y) = \ell(f_\theta(x), y)$, for all $g \in \mathcal{G}$ ($\Rightarrow \hat{R}_{\text{aug}} \approx R$)
- What if f_θ cannot express invariance to \mathcal{G} ?
 \Rightarrow Fit the data while maintaining a certain level of invariance

[Hounie, Chamon, Ribeiro, ICML'23]

65

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

- Which invariance?
✓ $\ell(f_\theta(gx), y) = \ell(f_\theta(x), y)$, for all $g \in \mathcal{G}$ ($\Rightarrow \hat{R}_{\text{aug}} \approx R$)
- What if f_θ cannot express invariance to \mathcal{G} ?
 \Rightarrow Fit the data while maintaining a certain level of invariance

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell_{\text{inv}}(f_\theta(x), y)] = \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{g \in \mathcal{G}} \ell(f_\theta(gx), y) \right] - \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_\theta(x), y)] \leq c$$

[Hounie, Chamon, Ribeiro, ICML'23]

65

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

- Which invariance?
✓ $\ell(f_\theta(gx), y) = \ell(f_\theta(x), y)$, for all $g \in \mathcal{G}$ ($\Rightarrow \hat{R}_{\text{aug}} \approx R$)
- What if f_θ cannot express invariance to \mathcal{G} ?
 \Rightarrow Fit the data while maintaining a certain level of invariance

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell_{\text{inv}}(f_\theta(x), y)] = \underbrace{\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{g \in \mathcal{G}} \ell(f_\theta(gx), y) \right]}_{\text{"Adversarial robustness"}} - \underbrace{\mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_\theta(x), y)]}_{R(\theta)} \leq c$$

[Hounie, Chamon, Ribeiro, ICML'23]

65

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

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

subject to Invariance $\leq c$

[Robey*, Chamon*, Pappas, Hassani, and Ribeiro, NeurIPS'21; Hounie, Chamon, Ribeiro, ICML'23]



From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

- Which invariance?
✓ $\ell(f_\theta(gx), y) = \ell(f_\theta(x), y)$, for all $g \in \mathcal{G}$ ($\Rightarrow \hat{R}_{\text{aug}} \approx R$)
- What if f_θ cannot express invariance to \mathcal{G} ?
 \Rightarrow Fit the data while maintaining a certain level of invariance

$$\ell_{\text{inv}}(f_\theta(x), y) = \max_{g \in \mathcal{G}} \ell(f_\theta(gx), y) - \ell(f_\theta(x), y) \geq 0 \quad (\text{id} \in \mathcal{G})$$

[Hounie, Chamon, Ribeiro, ICML'23]

65



From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

- Which invariance?
✓ $\ell(f_\theta(gx), y) = \ell(f_\theta(x), y)$, for all $g \in \mathcal{G}$ ($\Rightarrow \hat{R}_{\text{aug}} \approx R$)
- What if f_θ cannot express invariance to \mathcal{G} ?
 \Rightarrow Fit the data while maintaining a certain level of invariance

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell_{\text{inv}}(f_\theta(x), y)] = \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{g \in \mathcal{G}} \ell(f_\theta(gx), y) \right] - \underbrace{\mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_\theta(x), y)]}_{R(\theta)} \leq c$$

[Hounie, Chamon, Ribeiro, ICML'23]

65



From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

$$\min_{\theta} \text{Accuracy}$$

subject to Invariance $\leq c$

[Robey*, Chamon*, Pappas, Hassani, and Ribeiro, NeurIPS'21; Hounie, Chamon, Ribeiro, ICML'23]

66



From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

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

subject to $\frac{1}{N} \sum_{n=1}^N \left[\max_{g \in \mathcal{G}} \ell(f_\theta(gx_n), y_n) \right] \leq c$

[Robey*, Chamon*, Pappas, Hassani, and Ribeiro, NeurIPS'21; Hounie, Chamon, Ribeiro, ICML'23]

66

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) \\ \text{subject to} \quad & \frac{1}{N} \sum_{n=1}^N \left[\max_{g \in \mathcal{G}} \ell(f_{\theta}(gx_n), y_n) \right] \leq c \\ & \equiv (\text{DALE}) \\ \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{g \sim \mu_n} [\ell(f_{\theta}(gx_n), y_n)], \quad \mu_n = \begin{cases} \text{id} & \text{w.p. } \propto \frac{1}{\lambda} + \ell(f_{\theta}(x_n), y_n) \\ g & \text{w.p. } \propto \ell(f_{\theta}(gx_n), y_n) \end{cases} \end{aligned}$$

[Robey*, Chamon*, Pappas, Hassani, and Ribeiro, NeurIPS'21; Hounie, Chamon, Ribeiro, ICML'23]

66

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

- μ_{γ} is a data augmentation distribution (\mathcal{G}) that yields the best fit for a given invariance level (c)

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{g \sim \mu_n} [\ell(f_{\theta}(gx_n), y_n)], \quad \mu_n = \begin{cases} \text{id} & \text{w.p. } \propto \frac{1}{\lambda} + \ell(f_{\theta}(x_n), y_n) \\ g & \text{w.p. } \propto \ell(f_{\theta}(gx_n), y_n) \end{cases}$$

[Robey*, Chamon*, Pappas, Hassani, and Ribeiro, NeurIPS'21; Hounie, Chamon, Ribeiro, ICML'23]

66

From invariance to data augmentation

Goal

Which (if any) data augmentation (\mathcal{G}) induces invariance?

Invariance \rightarrow Robustness \rightarrow Data Augmentation

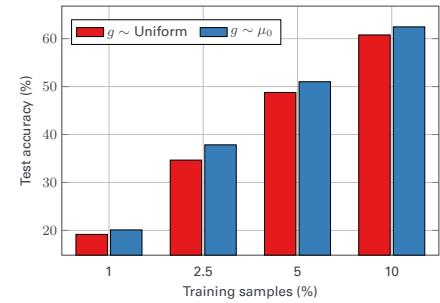
- μ_{γ} is a data augmentation distribution (\mathcal{G}) that yields the best fit for a given invariance level (c)
- No differentiability required (MH sampling)

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{g \sim \mu_n} [\ell(f_{\theta}(gx_n), y_n)], \quad \mu_n = \begin{cases} \text{id} & \text{w.p. } \propto \frac{1}{\lambda} + \ell(f_{\theta}(x_n), y_n) \\ g & \text{w.p. } \propto \ell(f_{\theta}(gx_n), y_n) \end{cases}$$

[Robey*, Chamon*, Pappas, Hassani, and Ribeiro, NeurIPS'21; Hounie, Chamon, Ribeiro, ICML'23]

66

Training on a subset of ImageNet-100



[Hounie, Chamon, Ribeiro, ICML'23]

67

“Identifying” invariances

Dataset	Dual variable (λ)	Synthetic Invariance		
		Rotation	Translation	Scale
MNIST	Rotation	0.000	2.724	0.012
	Translation	1.218	0.439	0.006
	Scale	2.026	4.029	0.003
F-MNIST	Rotation	0.000	3.301	1.352
	Translation	3.572	0.515	0.441
	Scale	4.144	2.725	0.904

[Hounie, Chamon, Ribeiro, ICML'23]

68

Claims

Constrained learning is the right tool to learn under requirements

Constrained learning is hard...

...but possible. How?

Primal-dual algorithms

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

70

Primal-dual algorithms

- Minimize the primal (\equiv ERM)

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

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

70

Primal-dual algorithms

- Minimize the primal (\equiv ERM)

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

[Haeffele et al., CVPR'17; Ge et al., ICLR'18; Mei et al., PNAS'18; Kawaguchi et al., AISTATS'20 ...]



70

Primal-dual algorithms

- Minimize the primal (\equiv ERM)

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

- Update the dual

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



70

[Paternain, C., Calvo-Fullana, and Ribeiro, NeurIPS'19; C. and Ribeiro, NeurIPS'20; C., Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23]

70

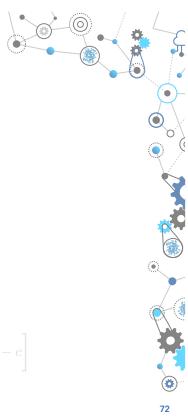
In practice...

- Minimize the primal (\equiv ERM)

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

- Update the dual

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



72

In practice...

```

1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_t \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, N$ 
5:      $\beta_{n+1} \leftarrow \beta_n - \eta \nabla_{\beta} [\ell(f_{\beta_n}(x_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(x_n), y_n)]$ 
6:   end
7:    $\theta_t \leftarrow \beta_{N+1}$ 
8:    $\lambda_t = \left[ \lambda_{t-1} + \eta \left( \frac{1}{N} \sum_{m=1}^N \ell(f_{\theta_t}(x_m), y_m) - c \right) \right]_+$ 
9: end
10: Output:  $\theta_T, \lambda_T$ 

```

PyTorch
<https://github.com/lfochamond/csl>

72

Primal-dual algorithms

- Minimize the primal (\equiv ERM)

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

- Update the dual

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

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

70



A (near-)PACC learner

Theorem

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

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

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

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

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

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

[Paternain, C., Calvo-Fullana, and Ribeiro, NeurIPS'19; C. and Ribeiro, NeurIPS'20; C., Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23]

71

In practice...

- Minimize the primal (\equiv ERM)

$$\theta^+ = \theta - \eta \nabla_{\theta} [\ell(f_{\theta}(x_n), y_n) + \lambda g(f_{\theta}(x_n), y_n)], \quad n = 1, 2, \dots, N$$

- Update the dual

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

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

72



Penalty-based vs. dual learning

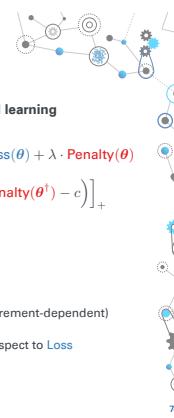
Penalty-based learning

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

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

- Parameter: λ (data-dependent)

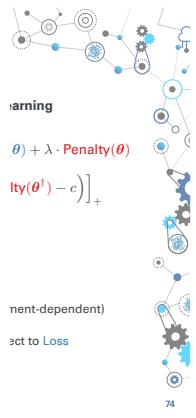
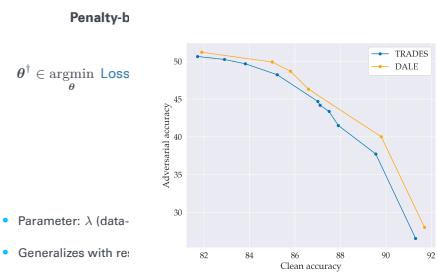
- Generalizes with respect to $\text{Loss} + \lambda \cdot \text{Penalty}$



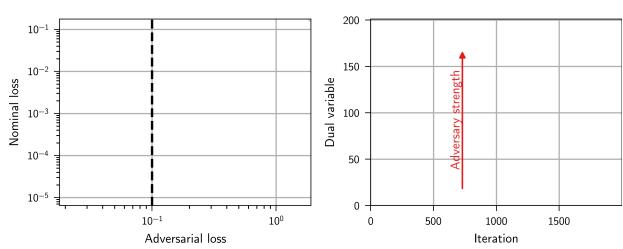
74

73

Penalty-based vs. dual learning

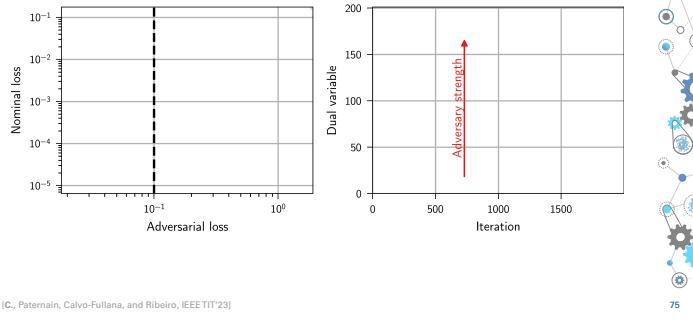


Robust image recognition

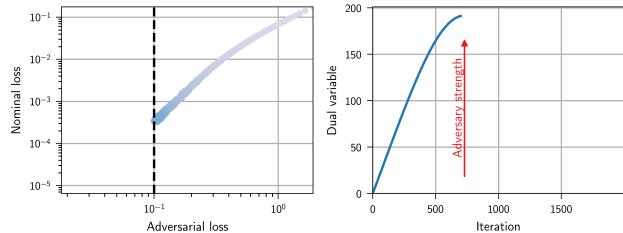


Robust image recognition

Robust image recognition

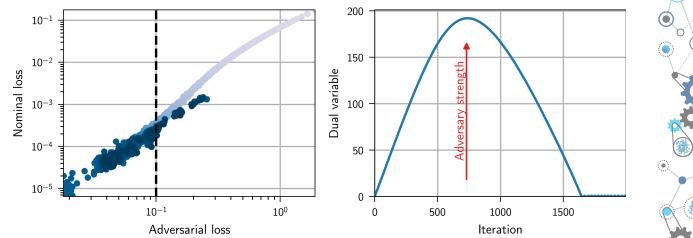


[C., Paternain, Calvo-Fullana, and Ribeiro, IEEEITIT'23]



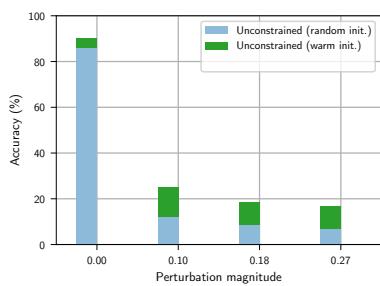
[C., Paternain, Calvo-Fullana, and Ribeiro, IEEEITIT'23]

Robust image recognition



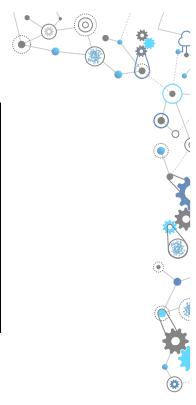
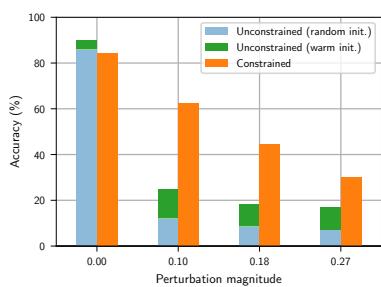
Empirical observations: [Zhang et al., ICML'20; Sitawarin, ArXiv'20]

Robust image recognition



[C., Paternain, Calvo-Fullana, and Ribeiro, IEEEITIT'23]

Robust image recognition

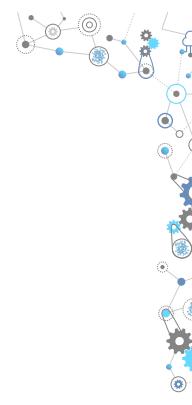
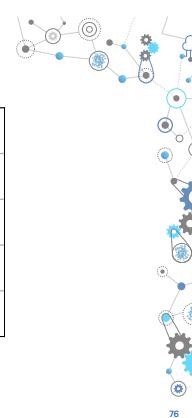
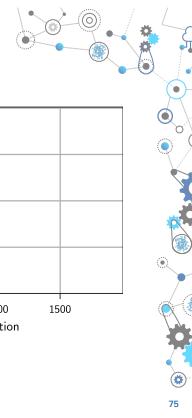


Claims

(1) Constrained learning is the right tool to learn under requirements

(2) Constrained learning is hard...

(3) ...but possible



Claims

(1) Constrained learning is the right A tool to learn under requirements

Constrained learning imposes generalizable requirements organically during training (resource budget, robustness...)

(2) Constrained learning is hard...

(3) ...but possible

77

Claims

(1) Constrained learning is the right A tool to learn under requirements

Constrained learning imposes generalizable requirements organically during training (resource budget, robustness...)

(2) Constrained learning is hard...

Constrained, non-convex, statistical optimization problem

(3) ...but possible

77

Claims

(1) Constrained learning is the right A tool to learn under requirements

Constrained learning imposes generalizable requirements organically during training (resource budget, robustness...)

(2) Constrained learning is hard...

Constrained, non-convex, statistical optimization problem

(3) ...but possible

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

77

Claims

(1) Constrained learning is the right A tool to learn under requirements

Constrained learning imposes generalizable requirements organically during training (resource budget, robustness...)

(2) Constrained learning is hard...

Constrained, non-convex, statistical optimization problem

(3) ...but possible. How?

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

77

The current landscape

- Constrained learning beyond resource allocation and adversarial robustness
 - Fairness [ICPR, ICASSP'20 (best student paper); CRR, NeurIPS'20; CPCR, IEEE TIT'23)], Invariance [ICRR, ICML'23], (manifold) smoothness [ICCV, ICML'23]...
 - Different forms of robustness [ICR, NeurIPS'20; CPCR, IEEE TIT'23; RCPH, ICML'22]
- Constrained (reinforcement) learning [PCCR, NeurIPS'19; PCCR, IEEE TAC'23; CPCR, arXiv'21 (submitted to IEEE TAC)]
- Gradient descent-ascent dynamics and primal recovery [PCCR, NeurIPS'19; CR, NeurIPS'20; CPCR, IEEE TSP'23; CPCR, arXiv'21]
- Constrained inference and control
 - Functional sparsity [ICR, Asilomar'19; CPCR, IEEE TSP'20; CER, IEEE TSP'20]
 - (non-submodular) sensor/actuator selection/scheduling [ICR, NeurIPS'17; CRR, IEEE TAC'21; CAR, IEEE TAC'22]
 - Risk-aware estimation and control [ICPR, ICASSP'20 (best paper award); TKRR, CDC'20]
 - Graph(on) signal processing/neural networks [RCR, NeurIPS'20; RCR, IEEE TSP'21; RCR, arXiv'21]

78

Index

Constrained learning is the right tool to learn under requirements

Wireless resource allocation
Robust image recognition

Constrained learning is hard...

...but possible

Break? (Fair learning)

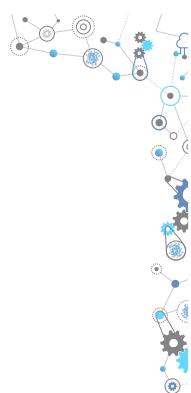
Break Robust learning revisited

...but possible. How?

Break Learning with invariance Fair learning

Approximate primal-dual for rate constraints

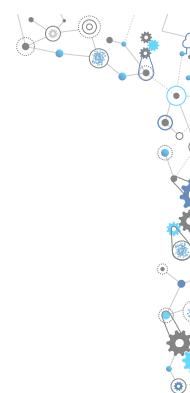
80



Fair learning

Problem

Predict whether an individual will recidivate

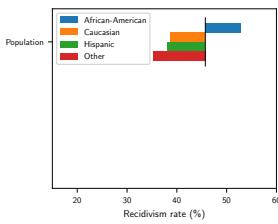


81

Fair learning

Problem

Predict whether an individual will recidivate



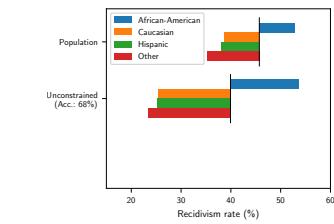
[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

81

Fair learning

Problem

Predict whether an individual will recidivate



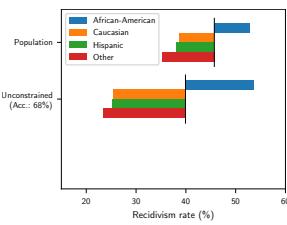
[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

81

Fair learning

Problem

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



[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

81

Fair learning

Problem

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

$$\begin{aligned} \min_{\theta} & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} & \text{Recidivism rate disparity (Race)} \leq c_r \\ & \text{for Race} \in \{\text{African-American, Caucasian, Hispanic, Other}\} \end{aligned}$$

[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

82

Fair learning

Problem

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

$$\begin{aligned} \min_{\theta} & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} & \text{Recidivism rate disparity (Race)} \leq c_r \\ & \text{for Race} \in \{\text{African-American, Caucasian, Hispanic, Other}\} \end{aligned}$$

[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

82

Fair learning

Problem

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

$$\begin{aligned} \min_{\theta} & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} & \Pr[f_{\theta}(x_n) = 1 \mid \text{Race}] \leq \Pr[f_{\theta}(x_n) = 1] + c_r \\ & \text{for Race} \in \{\text{African-American, Caucasian, Hispanic, Other}\} \end{aligned}$$

[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

82

Fair learning

Problem

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

$$\begin{aligned} \min_{\theta} & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} & \Pr[f_{\theta}(x_n) = 1 \mid \text{Race}] \leq \Pr[f_{\theta}(x_n) = 1] + c_r \\ & \text{for Race} \in \{\text{African-American, Caucasian, Hispanic, Other}\} \end{aligned}$$

[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

82

Fair learning

Problem

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

$$\begin{aligned} \min_{\theta} & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} & \Pr[f_{\theta}(x_n) = 1 \mid \text{Race}] \leq \Pr[f_{\theta}(x_n) = 1] + c_r \\ & \text{for Race} \in \{\text{African-American, Caucasian, Hispanic, Other}\} \end{aligned}$$

- Generalization guarantees (randomized solutions)

[Goh et al., NeurIPS'16; Kearns et al., ICML'18; Cotter et al., JMLR'19; ...]

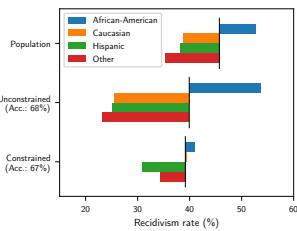
[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

82

Fair learning

Problem

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



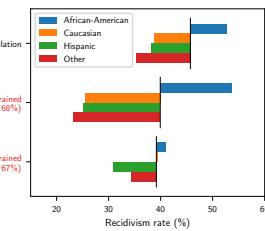
[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

83

Fair learning

Problem

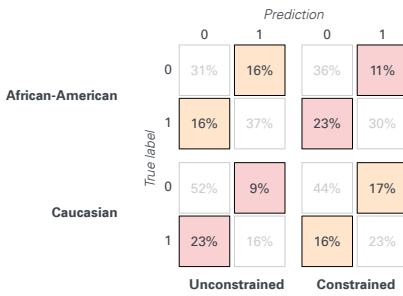
Predict whether an individual will recidivate at the same rate across races



[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

83

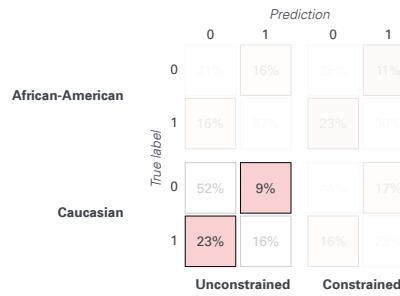
Fair prediction



[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

84

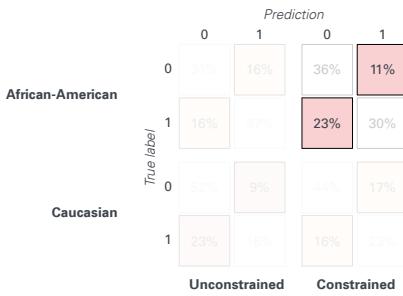
Fair prediction



[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

84

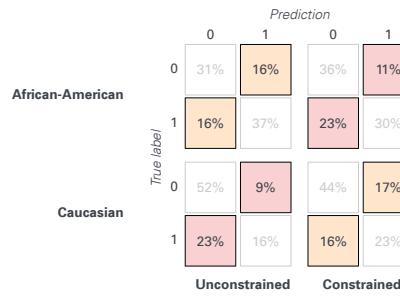
Fair prediction



[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

84

Fair prediction



[Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEETIT'23]

84

Index

Constrained learning is the right tool to learn under requirements

Wireless resource allocation
Robust image recognition

Constrained learning is hard...

...but possible

Break?
(Fair learning)

Break

Robust learning
revisited

...but possible. How?

Break

Learning
with invariance

Fair learning

Approximate primal-dual for rate constraints

85

Approximate primal-dual for rate constraints

- Minimize the primal

$$\theta^+ = \theta - \eta \left[\text{Loss}(f_\theta(x_n), y_n) + \lambda \nabla_\theta \left[\mathbb{I}(f_\theta(x_n) > 0.5) \mathbb{I}(\text{Race}) - \mathbb{I}(f_\theta(x_n) > 0.5) \right] \right]$$

- Update the dual

$$\lambda^+ = \left[\lambda + \eta \left(\frac{1}{N} \sum_{n=1}^N \mathbb{I}(f_{\theta^+}(x_n) > 0.5) \mathbb{I}(\text{Race}) - \mathbb{I}(f_{\theta^+}(x_n) > 0.5) - c \right) \right]_+$$

86

Approximate primal-dual for rate constraints

- Minimize the primal

$$\theta^+ = \theta - \eta \left[\text{Loss}(f_\theta(x_n), y_n) + \lambda \nabla_\theta \left[\mathbb{I}(f_\theta(x_n) > 0.5) \mathbb{I}(\text{Race}) - \mathbb{I}(f_\theta(x_n) > 0.5) \right] \right]$$

- Update the dual

$$\lambda^+ = \left[\lambda + \eta \left(\frac{1}{N} \sum_{n=1}^N \mathbb{I}(f_{\theta^+}(x_n) > 0.5) \mathbb{I}(\text{Race}) - \mathbb{I}(f_{\theta^+}(x_n) > 0.5) - c \right) \right]_+$$

86

Approximate primal-dual for rate constraints

- Minimize (approximately) the primal

$$\theta^+ = \theta - \eta \left[\text{Loss}(f_\theta(x_n), y_n) + \lambda \nabla_\theta \left[\sigma(f_\theta(x_n) - 0.5) \mathbb{I}(\text{Race}) - \sigma(f_\theta(x_n) - 0.5) \right] \right]$$

- Update the dual

$$\lambda^+ = \left[\lambda + \eta \left(\frac{1}{N} \sum_{n=1}^N \mathbb{I}(f_{\theta^+}(x_n) > 0.5) \mathbb{I}(\text{Race}) - \mathbb{I}(f_{\theta^+}(x_n) > 0.5) - c \right) \right]_+$$

86

Approximate primal-dual for rate constraints

- Minimize (approximately) the primal

$$\theta^+ = \theta - \eta \left[\text{Loss}(f_\theta(x_n), y_n) + \lambda \nabla_\theta \left[\sigma(f_\theta(x_n) - 0.5) \mathbb{I}(\text{Race}) - \sigma(f_\theta(x_n) - 0.5) \right] \right]$$

- Update the dual

$$\lambda^+ = \left[\lambda + \eta \left(\frac{1}{N} \sum_{n=1}^N \mathbb{I}(f_{\theta^+}(x_n) > 0.5) \mathbb{I}(\text{Race}) - \mathbb{I}(f_{\theta^+}(x_n) > 0.5) - c \right) \right]_+$$

86

Approximate primal-dual for rate constraints

- Minimize (approximately) the primal

$$\theta^+ = \theta - \eta \left[\text{Loss}(f_\theta(x_n), y_n) + \lambda \nabla_\theta \left[\sigma(f_\theta(x_n) - 0.5) \mathbb{I}(\text{Race}) - \sigma(f_\theta(x_n) - 0.5) \right] \right]$$

- Update the dual

$$\lambda^+ = \left[\lambda + \eta \left(\frac{1}{N} \sum_{n=1}^N \mathbb{I}(f_{\theta^+}(x_n) > 0.5) \mathbb{I}(\text{Race}) - \mathbb{I}(f_{\theta^+}(x_n) > 0.5) - c \right) \right]_+$$

86