

Luiz F. O. Chamon

**supervised and
reinforcement
learning under
requirements**

SIMPAS group meeting
Apr. 7th, 2025



Who am I?



Luiz

- 2025–: École Polytechnique de Paris (Professor)
- 2022–2024: ELLIS-SimTech (Research group leader)
IMPRS-IS faculty
- 2021–2022: Simons Institute, UC Berkeley (Postdoc)
- 2020: University of Pennsylvania (PhD)
- < 2015: University of São Paulo (BSc. & MSc.)
 - I speak 4.5 languages

Agenda

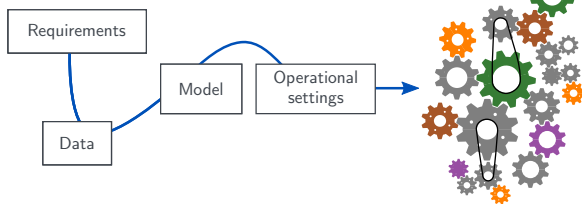
- I. Constrained supervised learning
 - Constrained learning theory
 - Constrained learning algorithms
 - Resilient constrained learning
- Break (10 min)
- II. Constrained reinforcement learning
 - Constrained RL duality
 - Constrained RL algorithms
- Q&A and discussions



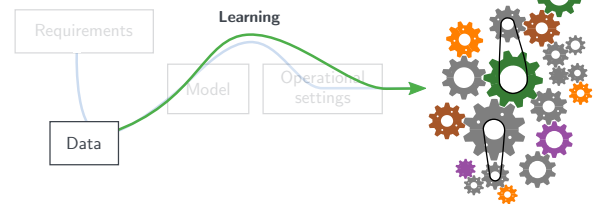
<https://luizchamon.com/sgm>

Why learning under requirements?

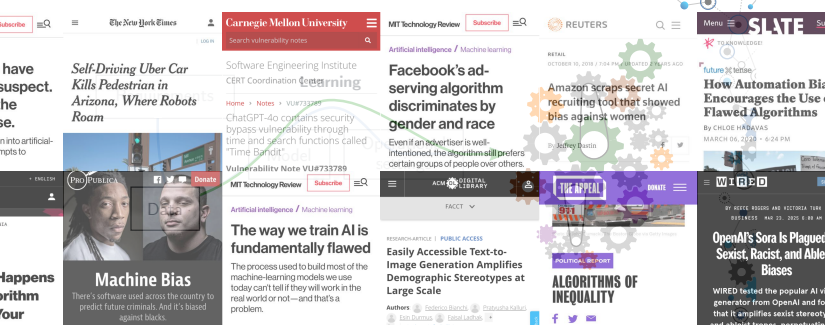
Why learning under requirements?



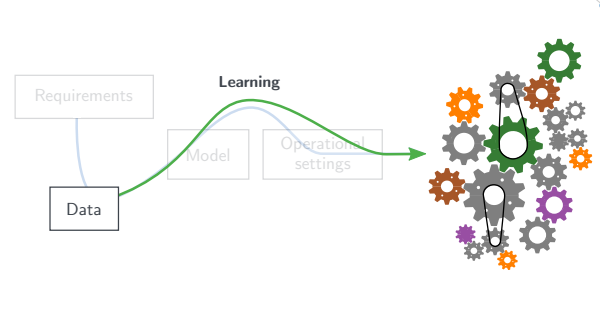
Why learning under requirements?



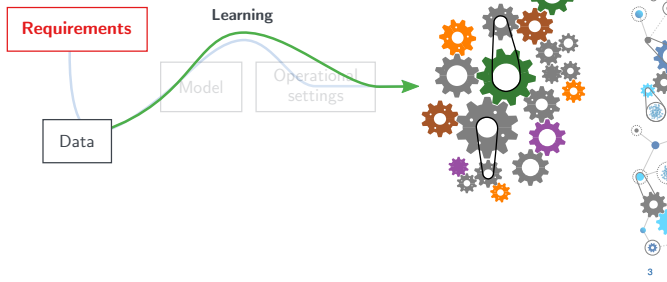
Why learning under requirements?



Why learning under requirements?

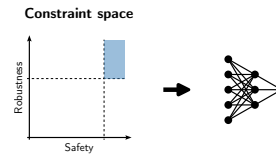


Why learning under requirements?



What is a requirements?

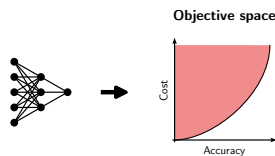
- Requirements are "shall" statements: describe necessary features subject to verification
 - Constraint space: things we decide



[NASA, "Systems engineering handbook," 2019]

What is a requirements?

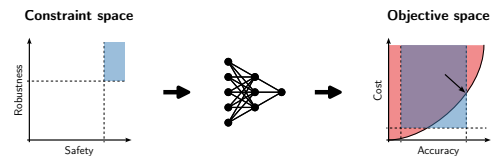
- Requirements are "shall" statements: describe necessary features subject to verification
 - Constraint space: things we decide
- Goals are "should" statements: express recommendations (once "shall" statements are satisfied)
 - Objective space: things the system achieves



[NASA, "Systems engineering handbook," 2019]

What is a requirements?

- Requirements are "shall" statements: describe necessary features subject to verification
 - Constraint space: things we decide
- Goals are "should" statements: express recommendations (once "shall" statements are satisfied)
 - Objective space: things the system achieves



[NASA, "Systems engineering handbook," 2019]

What is (un)constrained learning?

$$P_0^* = \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)]$$

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

[Chamon et al., IEEE ICASSP20 (best student paper); Chamon and Ribeiro, NeurIPS20; Chamon et al., IEEE TIT23]

What is (un)constrained learning?

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

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

$$h(f_{\theta}(x), y) \leq u, \quad \mathbb{Q}\text{-a.e.}$$

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

[Chamon et al., IEEE ICASSP20 (best student paper); Chamon and Ribeiro, NeurIPS20; Chamon et al., IEEE TIT23]

What about penalties?

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

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

$$h(f_{\theta}(x), y) \leq u, \quad \mathbb{Q}\text{-a.e.}$$

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)] + \lambda \mathbb{E}_{(x,y) \sim \mathcal{A}} [g(f_{\theta}(x), y)] + \mathbb{E}_{(x,y) \sim \mathcal{Q}} [\mu(x, y) h(f_{\theta}(x), y)]$$

What about penalties?

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

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

$$h(f_{\theta}(x), y) \leq u, \quad \mathbb{Q}\text{-a.e.}$$

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)] + \lambda \mathbb{E}_{(x,y) \sim \mathcal{A}} [g(f_{\theta}(x), y)] + \mathbb{E}_{(x,y) \sim \mathcal{Q}} [\mu(x, y) h(f_{\theta}(x), y)]$$

- There may not exist (λ, μ) such that the penalized solution is optimal and feasible
- Even if such (λ, μ) exist, they are not easy to find (hyperparameter search, cross-validation...)
- Constrained learning yields stronger guarantees, better performance, better trade-offs...

Applications

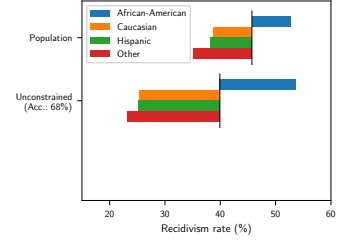
- **Fairness**
(e.g., [Goh et al., NeurIPS'16; Kearns et al., ICML'18; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23])
- **Federated learning**
(e.g., [Shen et al., ICLR'22; Hounie et al., NeurIPS'23])
- **Adversarially robust learning**
(e.g., [Chamon et al., NeurIPS'20; Robey et al., NeurIPS'21; Chamon et al., IEEE TIT'23])
- **Safe learning**
(e.g., [Paternain et al., IEEE TAC'23])
- **Wireless resource allocation**
(e.g., [Eisen et al., IEEE TSP'19; NaderiAlizadeh et al., IEEE TSP'22; Chowdhury et al., Asilomar'23])
- ...

7

Fairness

Problem

Predict whether an individual will recidivate



* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.

8

Fairness: "Equality" of odds

Problem

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

$$\begin{aligned} & \min_{\theta} \quad \text{Prediction error} \\ & \text{subject to} \quad \text{Prediction rate disparity (Race)} \leq c, \\ & \quad \text{for Race} \in \{\text{African-American, Caucasian, Hispanic, Other}\} \end{aligned}$$

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Goh et al., NeurIPS'16; Kearns et al., ICML'18; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23]

9

Fairness: "Equality" of odds

Problem

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

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

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Goh et al., NeurIPS'16; Kearns et al., ICML'18; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23]

9

Fairness: "Equality" of odds

Problem

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

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

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Goh et al., NeurIPS'16; Kearns et al., ICML'18; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23]

9

Applications

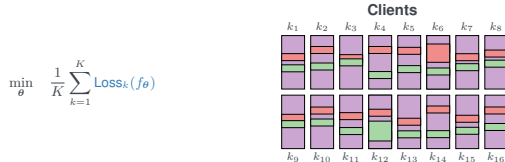
- **Fairness**
(e.g., [Goh et al., NeurIPS'16; Kearns et al., ICML'18; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23])
- **Federated learning**
(e.g., [Shen et al., ICLR'22; Hounie et al., NeurIPS'23])
- **Adversarially robust learning**
(e.g., [Chamon et al., NeurIPS'20; Robey et al., NeurIPS'21; Chamon et al., IEEE TIT'23])
- **Safe learning**
(e.g., [Paternain et al., IEEE TAC'23])
- **Wireless resource allocation**
(e.g., [Eisen et al., IEEE TSP'19; NaderiAlizadeh et al., IEEE TSP'22; Chowdhury et al., Asilomar'23])
- ...

10

Federated learning

Problem

Learn a common model using data from K clients



- k -th client loss: $\text{Loss}_k(f_{\theta}) = \frac{1}{N_k} \sum_{n_k=1}^{N_k} \text{Loss}(f_{\theta}(x_{n_k}), y_{n_k})$

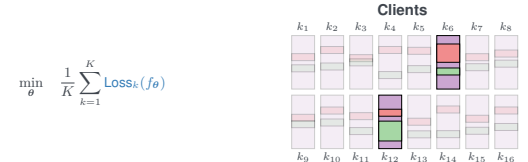
[Shen et al., ICLR'22]

11

Heterogeneous federated learning

Problem

Learn a common model using data from K clients



- k -th client loss: $\text{Loss}_k(f_{\theta}) = \frac{1}{N_k} \sum_{n_k=1}^{N_k} \text{Loss}(f_{\theta}(x_{n_k}), y_{n_k})$

[Shen et al., ICLR'22]

11

Heterogeneous federated learning

Problem

Learn a common model using data from K clients that is good for all clients

$$\min_{\theta} \quad \frac{1}{K} \sum_{k=1}^K \text{Loss}_k(f_{\theta})$$

subject to $\text{Loss}_k(f_{\theta}) \leq \frac{1}{K} \sum_{k=1}^K \text{Loss}_k(f_{\theta}) + c,$
 $k = 1, \dots, K$

- k -th client loss: $\text{Loss}_k(f_{\theta}) = \frac{1}{N_k} \sum_{n_k=1}^{N_k} \text{Loss}(f_{\theta}(x_{n_k}), y_{n_k})$

[Shen et al., ICRL'22]

11

Applications

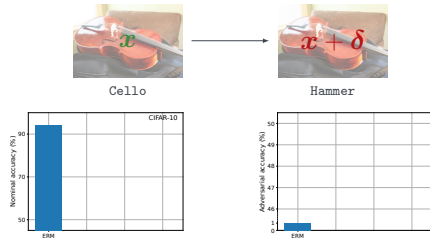
- Fairness
(e.g., [Goh et al., NeurIPS'16; Kearns et al., ICML'16; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23])
- Federated learning
(e.g., [Shen et al., ICLR'22; Hounie et al., NeurIPS'23])
- Adversarially robust learning
(e.g., [Chamon et al., NeurIPS'20; Robey et al., NeurIPS'21; Chamon et al., IEEE TIT'23])
- Safe learning
(e.g., [Paternain et al., IEEE TAC'23])
- Wireless resource allocation
(e.g., [Eisen et al., IEEE TSP'19; NaderiAlizadeh et al., IEEE TSP'22; Chowdhury et al., Asilomar'23])
- ...

12

Robustness

Problem

Learn an accurate classifier that is robust to input perturbations



13

Robustness

Problem

Learn an accurate classifier that is robust to input perturbations



$$\min_{\theta} \quad \text{Nominal loss}$$

subject to $\text{Robustness loss} \leq c$

[Chamon and Ribeiro, NeurIPS'20; Robey*, Chamon*, Pappas, Hassani, and Ribeiro, NeurIPS'21; Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23]

13

Robustness

Problem

Learn an accurate classifier that is robust to input perturbations



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

subject to $\text{Robustness loss} \leq c$

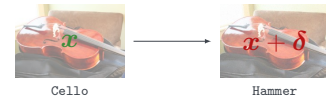
[Chamon and Ribeiro, NeurIPS'20; Robey*, Chamon*, Pappas, Hassani, and Ribeiro, NeurIPS'21; Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23]

13

Robustness

Problem

Learn an accurate classifier that is robust to input perturbations



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

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

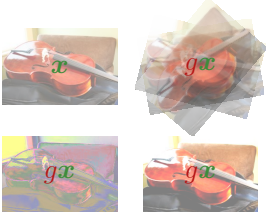
[Chamon and Ribeiro, NeurIPS'20; Robey*, Chamon*, Pappas, Hassani, and Ribeiro, NeurIPS'21; Chamon, Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23]

13

Invariance

Problem

Learn an accurate classifier that is invariant to transformation $g \in \mathcal{G}$, e.g., $\mathcal{G} = \left\{ \begin{array}{l} \text{Rotate, Translate}(Y), \\ \text{Shear}(Y), \text{Crop, Invert,} \\ \text{Solarize, Contrast,} \\ \text{Brightness, Sharpness...} \end{array} \right\}$



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

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

[Hounie, Chamon, Ribeiro, NeurIPS'23]

14

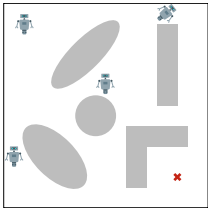
Applications

- Fairness
(e.g., [Goh et al., NeurIPS'16; Kearns et al., ICML'16; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23])
- Federated learning
(e.g., [Shen et al., ICLR'22; Hounie et al., NeurIPS'23])
- Adversarially robust learning
(e.g., [Chamon et al., NeurIPS'20; Robey et al., NeurIPS'21; Chamon et al., IEEE TIT'23])
- Safe learning
(e.g., [Paternain et al., IEEE TAC'23])
- Wireless resource allocation
(e.g., [Eisen et al., IEEE TSP'19; NaderiAlizadeh et al., IEEE TSP'22; Chowdhury et al., Asilomar'23])
- ...

15

Safety

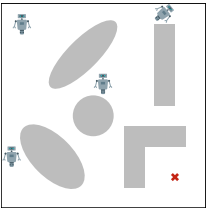
Problem
Find a control policy that navigates the environment effectively and safely



16

Safety

Problem
Find a control policy that navigates the environment effectively and safely



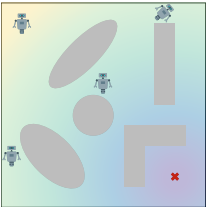
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC 23]

maximize $\pi \in \mathcal{P}(S)$ Task reward
subject to $\mathbb{P}[\text{Colliding with } O_i] \leq \delta,$
for $i = 1, 2, \dots$

17

Safety

Problem
Find a control policy that navigates the environment effectively and safely



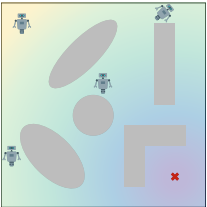
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC 23]

maximize $\pi \in \mathcal{P}(S)$ $\mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right]$
subject to $\mathbb{P}[\text{Colliding with } O_i] \leq \delta,$
for $i = 1, 2, \dots$

17

Safety

Problem
Find a control policy that navigates the environment effectively and safely



[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC 23]

maximize $\pi \in \mathcal{P}(S)$ $\mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right]$
subject to $\mathbb{P} \left(\bigcap_{t=0}^{T-1} \{s_t \notin O_i\} \mid \pi \right) \geq 1 - \delta,$
for $i = 1, 2, \dots$

17

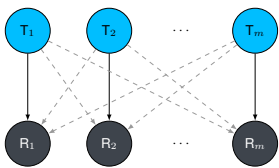
Applications

- Fairness
(e.g., [Goh et al., NeurIPS'16; Kearns et al., ICML'16; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23])
- Federated learning
(e.g., [Shen et al., ICLR'22; Hounie et al., NeurIPS'23])
- Adversarially robust learning
(e.g., [Chamon et al., NeurIPS'20; Robey et al., NeurIPS'21; Chamon et al., IEEE TIT'23])
- Safe learning
(e.g., [Paternain et al., IEEE TAC'23])
- Wireless resource allocation
(e.g., [Eisen et al., IEEE TSP'19; NaderiAlizadeh et al., IEEE TSP'22; Chowdhury et al., Asilomar'23])
- ...

18

Wireless resource allocation

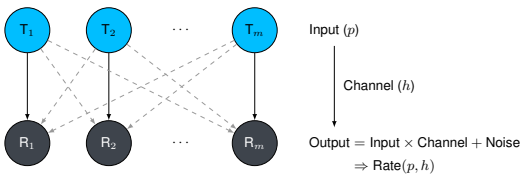
Problem
Allocate the least transmit power to m devices to achieve a communication rate



19

Wireless resource allocation

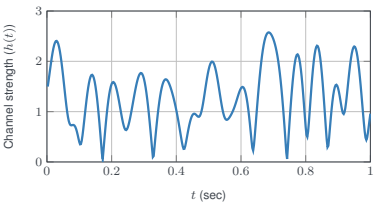
Problem
Allocate the least transmit power to m devices to achieve a communication rate



19

Wireless resource allocation

Problem
Allocate the least transmit power to m devices to achieve a communication rate

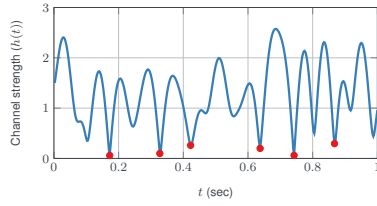


20

Wireless resource allocation

Problem

Allocate the least transmit power to m devices to achieve a communication rate

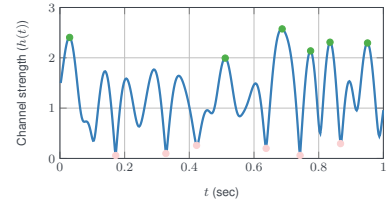


20

Wireless resource allocation

Problem

Allocate the least transmit power to m devices to achieve a communication rate

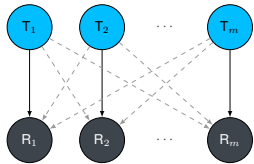


20

Wireless resource allocation

Problem

Allocate the least transmit power to m devices to achieve a communication rate



$$\begin{aligned} \min_{\pi \in \mathcal{P}(S)} \quad & \text{Total transmit power} \\ \text{s. to} \quad & \text{Rate } T_i \rightarrow R_i \geq c_i \end{aligned}$$

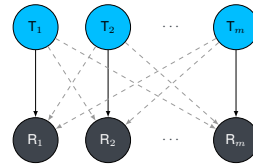
[Eisen, Zhang, Chamon, Lee, and Ribeiro, IEEE TSP'19]

21

Wireless resource allocation

Problem

Allocate the least transmit power to m devices to achieve a communication rate



$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \quad & - \sum_{i=1}^m \mathbb{E}_{h, p \sim \pi(h)} \left[\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^m p_{i,t} \right] \\ \text{s. to} \quad & \text{Rate } T_i \rightarrow R_i \geq c_i \end{aligned}$$

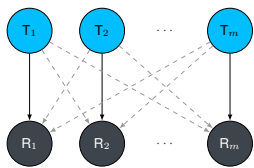
[Eisen, Zhang, Chamon, Lee, and Ribeiro, IEEE TSP'19]

21

Wireless resource allocation

Problem

Allocate the least transmit power to m devices to achieve a communication rate



$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \quad & - \sum_{i=1}^m \mathbb{E}_{h, p \sim \pi(h)} \left[\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^m p_{i,t} \right] \\ \text{s. to} \quad & \mathbb{E}_{h, p \sim \pi(h)} \left[\frac{1}{T} \sum_{t=0}^{T-1} \text{Rate}_i(p_t, h_t) \right] \geq c_i \end{aligned}$$

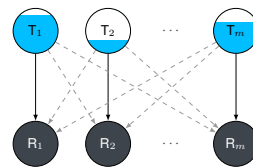
[Eisen, Zhang, Chamon, Lee, and Ribeiro, IEEE TSP'19]

21

Wireless resource allocation

Problem

Allocate power without depleting the battery of m devices to achieve a communication rate



$$\begin{aligned} \min_{\pi \in \mathcal{P}(S)} \quad & \text{Total probability of depleting battery} \\ \text{s. to} \quad & \mathbb{E}_{h, p \sim \pi(h, b)} \left[\frac{1}{T} \sum_{t=0}^{T-1} \text{Rate}_i(p_t, h_t) \right] \geq c_i \end{aligned}$$

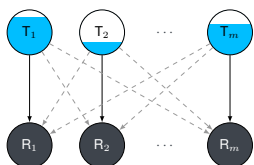
[Chowdhury, Paternain, Verma, Swami, Segarra, Asiloma'23]

21

Wireless resource allocation

Problem

Allocate power without depleting the battery of m devices to achieve a communication rate



$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \quad & - \sum_{i=1}^m \mathbb{E}_{h, p \sim \pi(h, b)} \left[\bigcap_{t=0}^{T-1} \{b_{i,t} = 0\} \right] \\ \text{s. to} \quad & \mathbb{E}_{h, p \sim \pi(h, b)} \left[\frac{1}{T} \sum_{t=0}^{T-1} \text{Rate}_i(p_t, h_t) \right] \geq c_i \end{aligned}$$

[Chowdhury, Paternain, Verma, Swami, Segarra, Asiloma'23]

21

And many more...

- Precision, recall, churn (e.g., [Cotter et al., JMLR'19])
- Scientific priors (e.g., [Lu et al., SIAM J. Sci. Comp.'21; Moro and Chamon, ICLR'25])
- Continual learning (e.g., [Peng et al., ICML'23])
- Active learning (e.g., [Elentier et al., NeurIPS'22])
- Semi-supervised learning (e.g., [Cervino et al., ICML'23])
- Minimum norm interpolation, SVM...

22

Constrained supervised learning

What is (un)constrained learning?

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

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

[Chamon et al., IEEE ICASSP'20 (best student paper); Chamon and Ribeiro, NeurIPS'20; Chamon et al., IEEE TIT'23]

Constrained learning challenges

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

Challenges

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

Constrained learning challenges

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

Challenges

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

Constrained learning challenges

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

Challenges

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

Constrained learning challenges

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

Challenges

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

Agenda

Constrained learning theory

Constrained learning algorithms

What classical learning theory says?

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

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



[Rostamizadeh, Talwalkar, Mohri. Foundations of machine learning, 2012]; [Ben-David, Shalev-Shwartz. Understanding machine learning..., 2014]

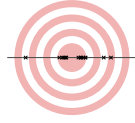
What's in a solution?

Definition (PAC learnability)

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

- near-optimal

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



[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

28

What's in a solution?

Definition (PACC learnability)

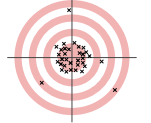
f_θ is a *probably approximately correct constrained (PACC)* learnable if for every ϵ, δ and every distributions \mathcal{D}, \mathcal{A} , we can obtain f_{θ^\dagger} from $N_f(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,

- near-optimal

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

- approximately feasible

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



[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

28

When is constrained learning possible?

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

Proposition

f_θ is PAC learnable \nRightarrow f_θ is PACC learnable

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

29

ECRM is not a PACC learner

Counter-example

$$\begin{aligned} \hat{P}^* &= \min_{\theta \in \Theta} J(\theta) \\ \text{subject to } \theta_2 \mathbb{E}_\tau[r] &\leq \theta_1 - 1 \\ &\quad - \theta_1 \mathbb{E}_\tau[r] \leq \theta_2 - 1 \end{aligned} \quad J(\theta) = \begin{cases} 1/16, & \theta = [1/2, 1/2] \\ 1/8, & \theta = [1, 1] \\ 1/4, & \theta = [1, 0] \end{cases}$$

- $\tau \sim \text{Uniform}(-1/2, 1/2)$

ECRM is not a PACC learner

Counter-example

$$\begin{aligned} \hat{P}^* &= \min_{\theta \in \Theta} J(\theta) = \frac{1}{8} \\ \text{subject to } \theta_2 \mathbb{E}_\tau[r] &\leq \theta_1 - 1 \Rightarrow \theta_1 \geq 1 \\ &\quad - \theta_1 \mathbb{E}_\tau[r] \leq \theta_2 - 1 \Rightarrow \theta_2 \leq 1 \end{aligned} \quad J(\theta) = \begin{cases} 1/16, & \theta = [1/2, 1/2] \\ 1/8, & \theta = [1, 1] \\ 1/4, & \theta = [1, 0] \end{cases}$$

- $\tau \sim \text{Uniform}(-1/2, 1/2)$

ECRM is not a PACC learner

Counter-example

$$\begin{aligned} \hat{P}^* &= \min_{\theta \in \Theta} J(\theta) = \frac{1}{8} \\ \text{subject to } \theta_2 \bar{\tau}_N &\leq \theta_1 - 1 \Rightarrow \theta_1 \geq 1 \\ &\quad - \theta_1 \bar{\tau}_N \leq \theta_2 - 1 \Rightarrow \theta_2 \leq 1 \end{aligned} \quad J(\theta) = \begin{cases} 1/16, & \theta = [1/2, 1/2] \\ 1/8, & \theta = [1, 1] \\ 1/4, & \theta = [1, 0] \end{cases}$$

$$\begin{aligned} \hat{P}_r^* &= \min_{\theta \in \Theta} J(\theta) \\ \text{subject to } \theta_2 \bar{\tau}_N &\leq \theta_1 - 1 + r_1 \\ &\quad - \theta_1 \bar{\tau}_N \leq \theta_2 - 1 + r_2 \end{aligned} \quad \mathbb{P}[|\hat{P}_r^* - P^*| \leq 1/32] \leq 4e^{-0.001N},$$

$$\text{unless } \bar{\tau}_N \leq r_1 < \frac{\bar{\tau}_N + 1}{2} \text{ and } r_2 \geq \bar{\tau}_N$$

- $\tau \sim \text{Uniform}(-1/2, 1/2) \rightarrow \bar{\tau}_N = \frac{1}{N} \sum_{n=1}^N \tau_n$

30

ECRM is not a PACC learner

Counter-example

$$\begin{aligned} \hat{P}^* &= \min_{\theta \in \Theta} J(\theta) = \frac{1}{8} \\ \text{subject to } \theta_2 \mathbb{E}_\tau[r] &\leq \theta_1 - 1 \Rightarrow \theta_1 \geq 1 \\ &\quad - \theta_1 \mathbb{E}_\tau[r] \leq \theta_2 - 1 \Rightarrow \theta_2 \leq 1 \end{aligned} \quad J(\theta) = \begin{cases} 1/16, & \theta = [1/2, 1/2] \\ 1/8, & \theta = [1, 1] \\ 1/4, & \theta = [1, 0] \end{cases}$$

$$\begin{aligned} \hat{P}^* &= \min_{\theta \in \Theta} J(\theta) \\ \text{subject to } \theta_2 \bar{\tau}_N &\leq \theta_1 - 1 \\ &\quad - \theta_1 \bar{\tau}_N \leq \theta_2 - 1 \end{aligned} \quad \mathbb{P}[|\hat{P}^* - P^*| \leq 1/32] = \mathbb{P}[\bar{\tau}_N = 0] = 0$$

- $\tau \sim \text{Uniform}(-1/2, 1/2) \rightarrow \bar{\tau}_N = \frac{1}{N} \sum_{n=1}^N \tau_n$

Constrained learning challenges

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

$$h(f_\theta(\mathbf{x}_r), y_r) \leq u \quad h(f_\theta(\mathbf{x}), y) \leq u \text{ a.e.}$$

Challenges

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

31

Constrained learning challenges

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

Challenges

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

Duality

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

Duality

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

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

An alternative path

$$\begin{aligned}
 \hat{P}^* &= \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) \\
 \text{s. to } &\frac{1}{N} \sum_{n=1}^N g(f_\theta, z_n) \leq c \\
 &\downarrow \text{PAC} \\
 P^* &= \min_{\theta \in \Theta} \mathbb{E}_2[\ell(f_\theta, z)] \\
 \text{s. to } &\mathbb{E}_2[g(f_\theta, z)] \leq c
 \end{aligned}
 \quad \longleftrightarrow \quad
 \begin{aligned}
 \hat{D}^* &= \max_{\lambda \geq 0} \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) + \lambda \left(\frac{1}{N} \sum_{n=1}^N g(f_\theta, z_n) - c \right)
 \end{aligned}$$

Duality

Duality

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

Duality

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

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

An alternative path

$$\begin{aligned}
 \hat{P}^* &= \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) & \xrightarrow{\text{PRIMAL}} & \hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) + \lambda \left(\frac{1}{N} \sum_{n=1}^N g(f_\theta, z_n) - c \right) \\
 \text{s. to } & \frac{1}{N} \sum_{n=1}^N g(f_\theta, z_n) \leq c & & \\
 & \downarrow \text{PAC} & & \\
 P^* &= \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_\theta, z)] & & \\
 \text{s. to } & \mathbb{E}_z [g(f_\theta, z)] \leq c & & \\
 & \downarrow \mathcal{H}_\theta \subset \mathcal{H} & & \\
 \tilde{P}^* &= \min_{\phi \in \mathcal{H}} \mathbb{E} [\ell(\phi, z)] & \xrightarrow{\text{?}} & \bar{D}^* = \max_{\lambda \geq 0} \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] + \lambda (\mathbb{E}_z [g(\phi, z)] - c) \\
 \text{s. to } & \mathbb{E}_z [g(\phi, z)] \leq c & &
 \end{aligned}$$

Non-convex variational duality

Convex optimization: Primal \longleftrightarrow Dual

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

34

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

[Chamon, Eldar, Ribeiro, IEEE TSP'20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

34

Sparse logistic regression

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

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

Discrete, non-convex

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

35

Sparse logistic regression

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

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

Discrete, non-convex

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

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

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

Continuous, non-convex

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

35

Sparse logistic regression

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

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

Discrete, non-convex

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

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

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

Continuous, non-convex

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

35

An alternative path

$$\hat{P}^* = \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, z_n) \quad \xrightarrow{\text{PAC}} \quad \hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, z_n) + \lambda \left(\frac{1}{N} \sum_{n=1}^N g(f_{\theta}, z_n) - c \right)$$

$$\text{s. to } \frac{1}{N} \sum_{n=1}^N g(f_{\theta}, z_n) \leq c$$

$$\downarrow \text{PAC}$$

$$P^* = \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_{\theta}, z)]$$

$$\text{s. to } \mathbb{E}_z [g(f_{\theta}, z)] \leq c$$

$$\hat{P}^* = \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] \quad \xrightarrow{=} \quad \hat{D}^* = \max_{\lambda \geq 0} \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] + \lambda (\mathbb{E}_z [g(\phi, z)] - c)$$

$$\text{s. to } \mathbb{E}_z [g(\phi, z)] \leq c$$

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

36

An alternative path

$$\hat{P}^* = \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, z_n) \quad \xrightarrow{\text{PAC}} \quad \hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, z_n) + \lambda \left(\frac{1}{N} \sum_{n=1}^N g(f_{\theta}, z_n) - c \right)$$

$$\text{s. to } \frac{1}{N} \sum_{n=1}^N g(f_{\theta}, z_n) \leq c$$

$$\downarrow \text{PAC}$$

$$P^* = \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_{\theta}, z)] \quad \xrightarrow{\epsilon_0} \quad D^* = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_{\theta}, z)] + \lambda (\mathbb{E}_z [g(f_{\theta}, z)] - c)$$

$$\text{s. to } \mathbb{E}_z [g(f_{\theta}, z)] \leq c$$

$$\uparrow \epsilon_0$$

$$\hat{P}^* = \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] \quad \xrightarrow{=} \quad \hat{D}^* = \max_{\lambda \geq 0} \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] + \lambda (\mathbb{E}_z [g(\phi, z)] - c)$$

$$\text{s. to } \mathbb{E}_z [g(\phi, z)] \leq c$$

$$\downarrow \epsilon_0$$

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

36

An alternative path

$$\hat{P}^* = \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, z_n) \quad \xrightarrow{\text{PAC}} \quad \hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, z_n) + \lambda \left(\frac{1}{N} \sum_{n=1}^N g(f_{\theta}, z_n) - c \right)$$

$$\text{s. to } \frac{1}{N} \sum_{n=1}^N g(f_{\theta}, z_n) \leq c$$

$$\downarrow \text{PAC}$$

$$P^* = \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_{\theta}, z)] \quad \xrightarrow{\epsilon_0} \quad D^* = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_{\theta}, z)] + \lambda (\mathbb{E}_z [g(f_{\theta}, z)] - c)$$

$$\text{s. to } \mathbb{E}_z [g(f_{\theta}, z)] \leq c$$

$$\uparrow \epsilon_0$$

$$\hat{P}^* = \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] \quad \xrightarrow{=} \quad \hat{D}^* = \max_{\lambda \geq 0} \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] + \lambda (\mathbb{E}_z [g(\phi, z)] - c)$$

$$\text{s. to } \mathbb{E}_z [g(\phi, z)] \leq c$$

$$\downarrow \epsilon_0$$

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

36

Dual (near-)PACC learning

Theorem

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

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

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

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentier, Chamon, Ribeiro, ICLR24]

37

Dual (near-)PACC learning

Theorem

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

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

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

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

(mild additional conditions apply)

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentier, Chamon, Ribeiro, ICLR24]

37

Dual (near-)PACC learning

Theorem

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

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

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

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

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

$$(\ell_0 \text{ strongly convex and } g, h \text{ convex}) \quad h(f_{\theta^\dagger}(x), y) \leq r, \text{ with } \mathfrak{P}\text{-prob. } 1 - \tilde{O} \left(\nu + \frac{1}{\sqrt{N}} \right)$$

(mild additional conditions apply)

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentier, Chamon, Ribeiro, ICLR24]

37

Dual (near-)PACC learning

Theorem

Let f be ν -universal with VC dimension $d_{VC} < \infty$, ℓ_0 strongly convex, and g convex. Then, f_{θ^\dagger} is a (near-)PACC solution of (P-CSL) for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , i.e., with probability at least $1 - \delta$,

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

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

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

Sources of error

parametrization richness (ν) sample size (N) requirements difficulty (λ^*)

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentier, Chamon, Ribeiro, ICLR24]

38

Dual (near-)PACC learning

Theorem

Let f be ν -universal with VC dimension $d_{VC} < \infty$, ℓ_0 strongly convex, and g convex. Then, f_{θ^\dagger} is a (near-)PACC solution of (P-CSL) for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , i.e., with probability at least $1 - \delta$,

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

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

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

Sources of error

parametrization richness (ν) sample size (N) requirements difficulty (λ^*)

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentier, Chamon, Ribeiro, ICLR24]

38

Dual (near-)PACC learning

Theorem

Let f be ν -universal with VC dimension $d_{VC} < \infty$, ℓ_0 strongly convex, and g convex. Then, f_{θ^\dagger} is a (near-)PACC solution of (P-CSL) for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , i.e., with probability at least $1 - \delta$,

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

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

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

Sources of error

parametrization richness (ν) sample size (N) requirements difficulty (λ^*)

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentier, Chamon, Ribeiro, ICLR24]

38

Dual (near-)PACC learning

Theorem

Let f be ν -universal with VC dimension $d_{VC} < \infty$, ℓ_0 strongly convex, and g convex. Then, f_{θ^\dagger} is a (near-)PACC solution of (P-CSL) for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , i.e., with probability at least $1 - \delta$,

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

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

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

Sources of error

parametrization richness (ν) sample size (N) requirements difficulty (λ^*)

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentier, Chamon, Ribeiro, ICLR24]

38

Dual (near-)PACC learning

Theorem

Let f be ν -universal with VC dimension $d_{VC} < \infty$, ℓ_0 strongly convex, and g convex. Then, f_{θ^\dagger} is a (near-)PACC solution of (P-CSL) for all $(\theta^\dagger, \lambda^\dagger)$ that achieve \hat{D}^* , i.e., with probability at least $1 - \delta$,

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

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

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

Sources of error

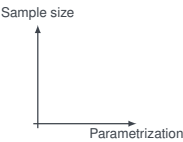
parametrization richness (ν) sample size (N) requirements difficulty (λ^*)

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentier, Chamon, Ribeiro, ICLR24]

38

Dual learning trade-offs

- Unconstrained learning
parametrization × sample size

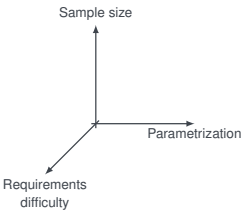


[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

39

Dual learning trade-offs

- Unconstrained learning
parametrization × sample size
- Constrained learning
parametrization × sample size × requirements



[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

39

When is constrained learning possible?

Corollary

f_{θ} is PAC learnable \approx^* f_{θ} is PACC learnable

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

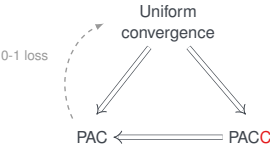
[mild conditions apply]

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

40

When is constrained learning possible?

Corollary



[mild conditions apply]

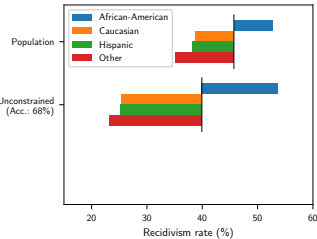
[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

40

Fairness

Problem

Predict whether an individual will recidivate



* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.

41

Fairness: "Equality" of odds

Problem

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

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

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Cotter et al., JMLR'19; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

42

Fairness: "Equality" of odds

Problem

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

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

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Cotter et al., JMLR'19; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

42

Fairness: "Equality" of odds

Problem

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

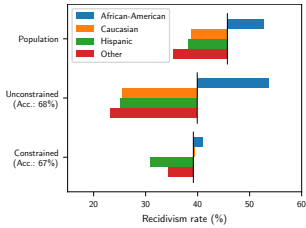
$$\begin{aligned} & \min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ & \text{subject to } \frac{1}{N} \sum_{n=1}^N \sigma(f_{\theta}(x_n) - 0.5) \mathbb{I}[x_n \in \text{Race}] \leq \frac{1}{N} \sum_{n=1}^N \sigma(f_{\theta}(x_n) - 0.5) + c, \\ & \text{for Race} \in \{\text{African-American, Caucasian, Hispanic, Other}\} \end{aligned}$$

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Cotter et al., JMLR'19; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

42

Fairness: “Equality” of odds

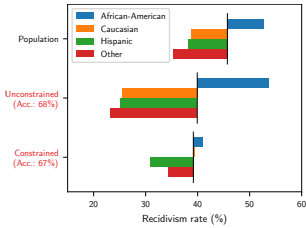
Problem
Predict whether an individual will recidivate **at the same rate across races**



* We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT 23]

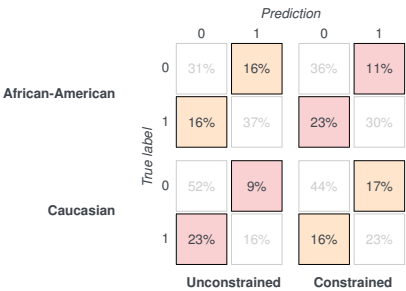
Fairness: “Equality” of odds

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



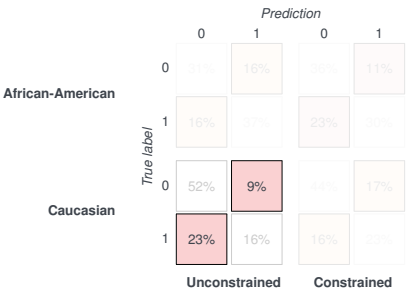
* We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT 23]

Fairness: “Equality” of odds



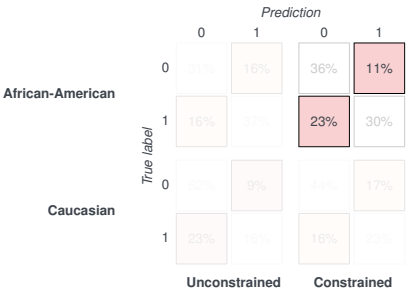
* We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT 23]

Fairness: “Equality” of odds



* We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT 23]

Fairness: “Equality” of odds



* We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT 23]

Agenda

Constrained learning theory

Constrained learning algorithms

Constrained optimization methods

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

subject to

$$\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c$$
$$h(f_{\theta}(x_r), y_r) \leq u$$

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

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

subject to

$$\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c$$
$$h(f_{\theta}(x_r), y_r) \leq u$$

Constrained optimization methods

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

subject to $\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c$
 $h(f_{\theta}(x_r), y_r) \leq u$

- 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
 - ✓ (near-)feasible solution [small duality gap]

46

Dual learning algorithm

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

47

Dual learning algorithm

- Minimize the primal (\equiv ERM)

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

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

47

Dual learning algorithm

- Minimize the primal (\equiv ERM)

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

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

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

47

Dual learning algorithm

- Minimize the primal (\equiv ERM)

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

- Update the dual

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

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

47

A (near-)PACC learner

Theorem

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

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

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

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

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

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

[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

48

In practice...

- Minimize the primal (\equiv ERM)

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

- Update the dual

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

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

49

In practice...

- Minimize the primal (\equiv ERM)

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

- Update the dual

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

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

49

In practice...

```
1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_1 \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, N$ 
5:      $\beta_{n+1} \leftarrow \beta_n - \eta_\theta \nabla_{\beta} [\ell(f_{\beta_n}(\mathbf{x}_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(\mathbf{x}_n), y_n)]$ 
6:   end
7:    $\theta_t \leftarrow \beta_{N+1}$ 
8:    $\lambda_t = \left[ \lambda_{t-1} + \eta_\lambda \left( \frac{1}{N} \sum_{m=1}^N g(f_{\theta_t}(\mathbf{x}_m), y_n) - c \right) \right]_+$ 
9: end
10: Output:  $\theta_T, \lambda_T$ 
```

SGD

Dual update



https://github.com/lfochamon/csl

In practice...

```
1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_1 \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, N$ 
5:      $\beta_{n+1} \leftarrow \beta_n - \eta_\theta \nabla_{\beta} [\ell(f_{\beta_n}(\mathbf{x}_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(\mathbf{x}_n), y_n)]$ 
6:   end
7:    $\theta_t \leftarrow \beta_{N+1}$ 
8:    $\lambda_t = \left[ \lambda_{t-1} + \eta_\lambda \left( \frac{1}{N} \sum_{m=1}^N g(f_{\theta_t}(\mathbf{x}_m), y_n) - c \right) \right]_+$ 
9: end
10: Output:  $\theta_T, \lambda_T$ 
```

Use adaptive method (e.g., ADAM)



https://github.com/lfochamon/csl

In practice...

```
1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_1 \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, N$ 
5:      $\beta_{n+1} \leftarrow \beta_n - \eta_\theta \nabla_{\beta} [\ell(f_{\beta_n}(\mathbf{x}_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(\mathbf{x}_n), y_n)]$ 
6:   end
7:    $\theta_t \leftarrow \beta_{N+1}$ 
8:    $\lambda_t = \left[ \lambda_{t-1} + \eta_\lambda \left( \frac{1}{N} \sum_{m=1}^N g(f_{\theta_t}(\mathbf{x}_m), y_n) - c \right) \right]_+$ 
9: end
10: Output:  $\theta_T, \lambda_T$ 
```

Use adaptive method (e.g., ADAM)
Use different time-scales ($\eta_\lambda = 0.1\eta_\theta$)



https://github.com/lfochamon/csl

In practice...

```
1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_1 \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, N$ 
5:      $\beta_{n+1} \leftarrow \beta_n - \eta_\theta \nabla_{\beta} [\ell(f_{\beta_n}(\mathbf{x}_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(\mathbf{x}_n), y_n)]$ 
6:   end
7:    $\theta_t \leftarrow \beta_{N+1}$ 
8:    $\lambda_t = \left[ \lambda_{t-1} + \eta_\lambda \left( \frac{1}{N} \sum_{m=1}^N g(f_{\theta_t}(\mathbf{x}_m), y_n) - c \right) \right]_+$ 
9: end
10: Output:  $\theta_T, \lambda_T$ 
```

Check slack:
- feasibility: $s_t \leq 0$
- "duality gap": $\lambda_t s_t$
 $s_t = \frac{1}{N} \sum_{n=1}^N g(f_{\theta_t}(\mathbf{x}_n), y_n) - c$

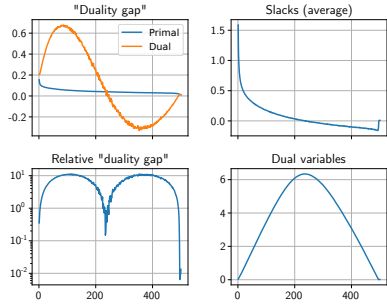
Use adaptive method (e.g., ADAM)
Use different time-scales ($\eta_\lambda = 0.1\eta_\theta$)



https://github.com/lfochamon/csl

In practice...

```
1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_1 \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, l$ 
5:      $\beta_{n+1} \leftarrow \beta_n$ 
6:   end
7:    $\theta_t \leftarrow \beta_{N+1}$ 
8:    $\lambda_t = \left[ \lambda_{t-1} + \dots \right]_+$ 
9: end
10: Output:  $\theta_T, \lambda_T$ 
```



ethod (e.g., ADAM)
ne-scales ($\eta_\lambda = 0.1\eta_\theta$)



https://github.com/lfochamon/csl

Penalty-based vs. dual learning

Penalty-based learning

$\theta^1 \in \operatorname{argmin}_{\theta} \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$

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

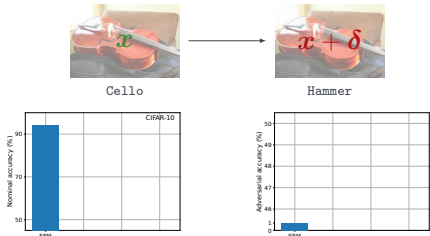
Dual learning

$\theta^1 \in \operatorname{argmin}_{\theta} \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$
 $\lambda^+ = \left[\lambda + \eta (\text{Penalty}(\theta^1) - c) \right]_+$

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

Robust learning

Problem
Learn an accurate classifier that is robust to input perturbations



Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations

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

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



Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations

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

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

\approx gradient ascent

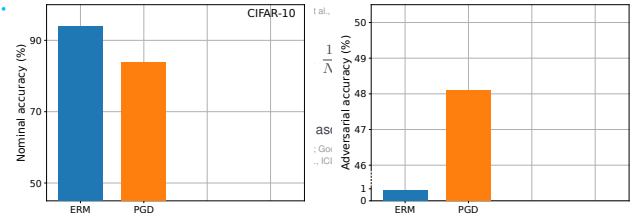
[Szegedy et al., ICLR'14; Goodfellow et al., ICLR'15; Madry et al., ICLR'18; ...]



53

Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations



53

Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations

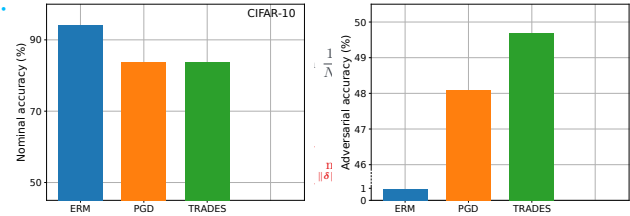
- Adversarial training (e.g., [Zhang et al., ICML'19])

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

54

Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations



54

Constrained learning for robustness

Problem
Learn an accurate classifier that is robust to input perturbations

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

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

[Chamion and Ribeiro, NeurIPS'20; Robey et al., NeurIPS'21; Chamion et al., IEEE TIT'23]

55

Constrained learning for robustness

Problem
Learn an accurate classifier that is robust to input perturbations

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

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

[C. and Ribeiro, NeurIPS'20; Robey, C., Pappas, Hassani, and Ribeiro, NeurIPS'21; C., Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

56

Constrained learning for robustness

Problem
Learn an accurate classifier that is robust to input perturbations

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

subject to $\frac{1}{N} \sum_{n=1}^N u_n \leq c$

$\text{Loss}(f_{\theta}(x_n + \delta_n), y_n) \leq u_n, \text{ for all } \|\delta_n\|_{\infty} \leq \epsilon$

[C. and Ribeiro, NeurIPS'20; Robey, C., Pappas, Hassani, and Ribeiro, NeurIPS'21; C., Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

56

Constrained learning for robustness

Problem
Learn an accurate classifier that is robust to input perturbations

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

subject to $\frac{1}{N} \sum_{n=1}^N u_n \leq c$

Sampling (e.g., LMC)

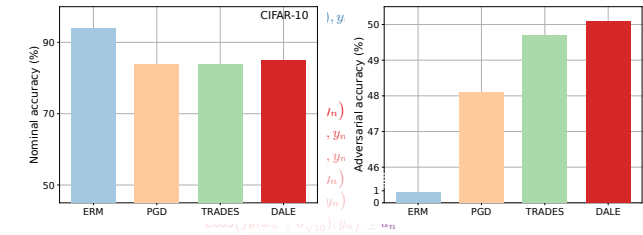
$\begin{cases} \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq u_n \\ \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq u_n \\ \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{3}}), y_n) \leq u_n \\ \text{Loss}(f_{\theta}(x_n + \delta_{\epsilon}), y_n) \leq u_n \\ \text{Loss}(f_{\theta}(x_n + \delta_{\epsilon}), y_n) \leq u_n \\ \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{10}}), y_n) \leq u_n \end{cases}$

[C. and Ribeiro, NeurIPS'20; Robey, C., Pappas, Hassani, and Ribeiro, NeurIPS'21; C., Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

56

Constrained learning for robustness

Problem
Learn an accurate classifier that is robust to input perturbations

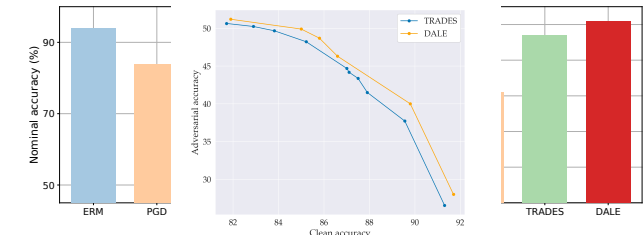


[C. and Ribeiro, NeurIPS'20; Robey*, C.*, Pappas, Hassani, and Ribeiro, NeurIPS'21; C., Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

56

Constrained learning for robustness

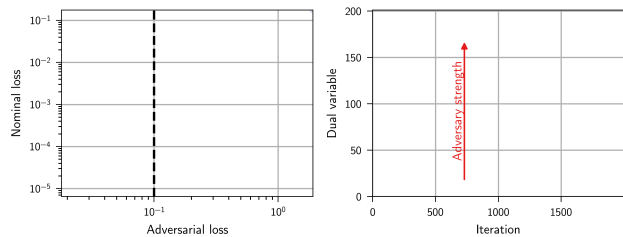
Problem
Learn an accurate classifier that is robust to input perturbations



[C. and Ribeiro, NeurIPS'20; Robey*, C.*, Pappas, Hassani, and Ribeiro, NeurIPS'21; C., Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

56

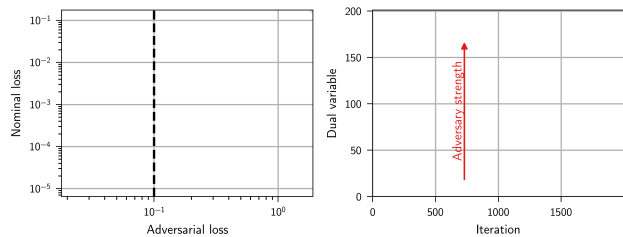
Constrained learning for robustness



[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

57

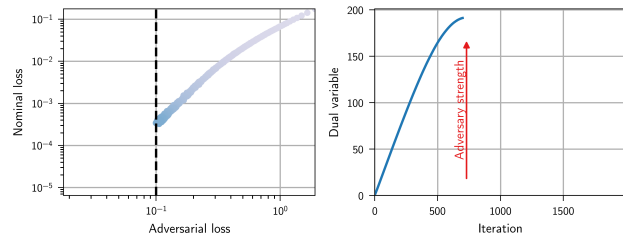
Constrained learning for robustness



[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

57

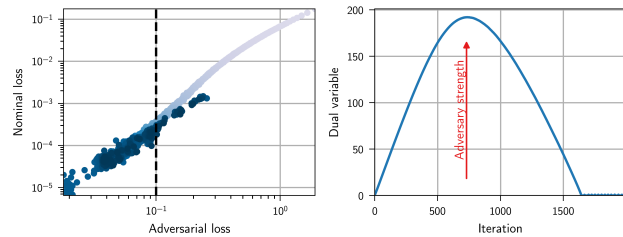
Constrained learning for robustness



[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

57

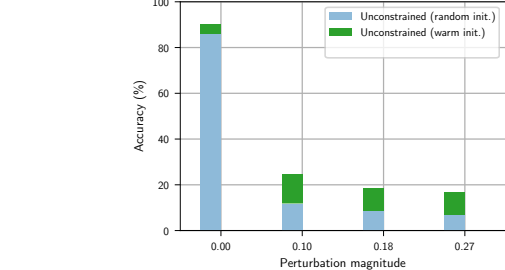
Constrained learning for robustness



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

57

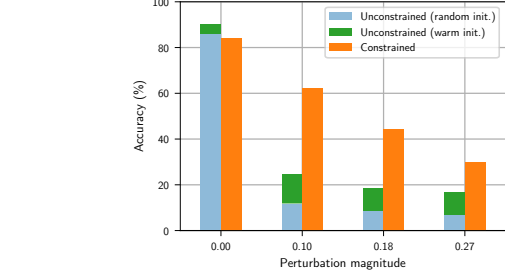
Constrained learning for robustness



[Chamon et al., IEEE TIT'23]

58

Constrained learning for robustness



[Chamon et al., IEEE TIT'23]

58

Penalty-based vs. dual learning

Penalty-based learning

$$\theta^1 \in \operatorname{argmin}_{\theta} \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$$

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

Dual learning

$$\theta^1 \in \operatorname{argmin}_{\theta} \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$$
$$\lambda^+ = \left[\lambda + \eta \left(\text{Penalty}(\theta^1) - c \right) \right]_+$$

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

59

Summary

- Constrained learning is ~~the~~ a tool to learn under requirements
- Constrained learning is hard...
- ...but possible. How?

60

Summary

- Constrained learning is ~~the~~ a tool to learn under requirements
Constrained learning imposes generalizable requirements organically during training, e.g., fairness [Chamon and Ribeiro, NeurIPS 20; Chamon et al., IEEE TIT'23], heterogeneity [Shen et al., ICRL22]. ...
- Constrained learning is hard...
- ...but possible. How?

60

Summary

- Constrained learning is ~~the~~ a tool to learn under requirements
Constrained learning imposes generalizable requirements organically during training, e.g., fairness [Chamon and Ribeiro, NeurIPS 20; Chamon et al., IEEE TIT'23], heterogeneity [Shen et al., ICRL22]. ...
- Constrained learning is hard...
Constrained, non-convex, statistical optimization problem
- ...but possible. How?

60

Summary

- Constrained learning is ~~the~~ a tool to learn under requirements
Constrained learning imposes generalizable requirements organically during training, e.g., fairness [Chamon and Ribeiro, NeurIPS 20; Chamon et al., IEEE TIT'23], heterogeneity [Shen et al., ICRL22]. ...
- Constrained learning is hard...
Constrained, non-convex, statistical optimization problem
- ...but possible. How?
We can learn under requirements (essentially) whenever we can learn at all by solving (penalized) ERM problems.

60

Robust/resilient constraints

Agenda

- Resilient constrained learning
 - Semi-infinite learning
 - Probabilistic robustness
-
- 62

Agenda

- Resilient constrained learning
 - Semi-infinite learning
 - Probabilistic robustness
-
- 63

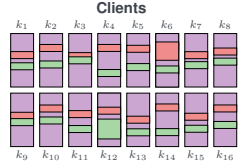
Heterogeneous federated learning

Problem

Learn a common model using data from K clients that is good for all clients

$$\min_{\theta} \frac{1}{K} \sum_{k=1}^K \text{Loss}_k(f_{\theta})$$

subject to $\text{Loss}_k(f_{\theta}) \leq \frac{1}{K} \sum_{k=1}^K \text{Loss}_k(f_{\theta}) + c$
 $k = 1, \dots, K$



- k -th client loss: $\text{Loss}_k(\phi) = \frac{1}{N_k} \sum_{n_k=1}^{N_k} \text{Loss}(f_{\theta}(x_{n_k}), y_{n_k})$

64

Heterogeneous federated learning

Problem

Learn a common model using data from K clients that is good for all clients

$$\min_{\theta} \frac{1}{K} \sum_{k=1}^K \text{Loss}_k(f_{\theta})$$

subject to $\text{Loss}_k(f_{\theta}) \leq \frac{1}{K} \sum_{k=1}^K \text{Loss}_k(f_{\theta}) + c_k$
 $k = 1, \dots, K$



- k -th client loss: $\text{Loss}_k(\phi) = \frac{1}{N_k} \sum_{n_k=1}^{N_k} \text{Loss}(f_{\theta}(x_{n_k}), y_{n_k})$

64

Resilient constrained learning

Definition (Resilience)

(ecology) ability of an ecosystem to adapt its function to accommodate operating conditions

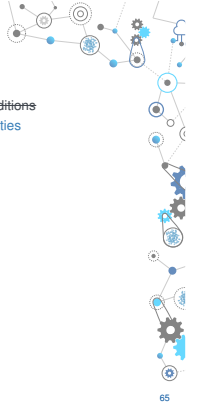


65

Resilient constrained learning

Definition (Resilience)

(ecology) ability of an ecosystem to adapt its function to accommodate operating conditions
 (learning) learning system specification data properties



65

Resilient constrained learning

Definition (Resilience)

(ecology) ability of an ecosystem to adapt its function to accommodate operating conditions
 (learning) learning system specification data properties

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

subject to $\mathbb{E}_{(x,y) \sim \mathcal{Q}_i} [g_i(f_{\theta}(x_m), y_m)] \leq c_i$



65

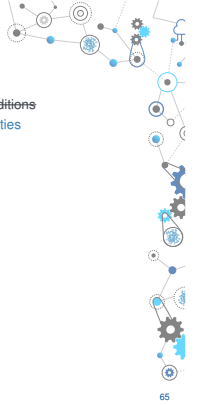
Resilient constrained learning

Definition (Resilience)

(ecology) ability of an ecosystem to adapt its function to accommodate operating conditions
 (learning) learning system specification data properties

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

subject to $\mathbb{E}_{(x,y) \sim \mathcal{Q}_i} [g_i(f_{\theta}(x_m), y_m)] \leq c_i + \mathbf{r}_i$



65

Resilient constrained learning

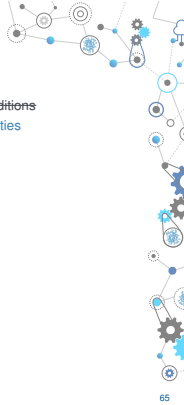
Definition (Resilience)

(ecology) ability of an ecosystem to adapt its function to accommodate operating conditions
 (learning) learning system specification data properties

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

subject to $\mathbb{E}_{(x,y) \sim \mathcal{Q}_i} [g_i(f_{\theta}(x_m), y_m)] \leq c_i + \mathbf{r}_i$

- Larger relaxations \mathbf{r} decrease the objective $P^*(\mathbf{r})$ (benefit), but increase specification violation $c_i + \mathbf{r}_i$ (cost)



65

Resilient constrained learning

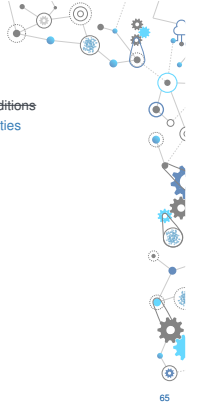
Definition (Resilience)

(ecology) ability of an ecosystem to adapt its function to accommodate operating conditions
 (learning) learning system specification data properties

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

subject to $\mathbb{E}_{(x,y) \sim \mathcal{Q}_i} [g_i(f_{\theta}(x_m), y_m)] \leq c_i + \mathbf{r}_i$

- Larger relaxations \mathbf{r} decrease the objective $P^*(\mathbf{r})$ (benefit), but increase specification violation $c_i + \mathbf{r}_i$ (cost)
- Resilience is a compromise!



65

Resilient constrained learning

Definition (Resilient equilibrium)

For a strictly convex function $h(r)$, we say the relaxation r^* achieves the resilient equilibrium if

$$\nabla h(r^*) \in -\partial P^*(r^*) \quad \leftarrow (\partial: \text{subdifferential})$$

In words: at the resilient equilibrium the **marginal cost of relaxing** equals the **marginal gain of relaxing**

[Hounie, Chamon, Ribeiro, NeurIPS'23]

66

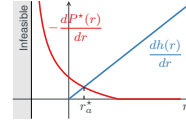
Resilient constrained learning

Definition (Resilient equilibrium)

For a strictly convex function $h(r)$, we say the relaxation r^* achieves the resilient equilibrium if

$$\nabla h(r^*) \in -\partial P^*(r^*) \quad \leftarrow (\partial: \text{subdifferential})$$

In words: at the resilient equilibrium the **marginal cost of relaxing** equals the **marginal gain of relaxing**



[Hounie, Chamon, Ribeiro, NeurIPS'23]

66

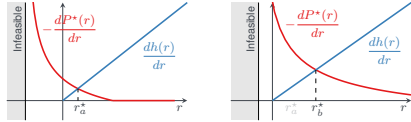
Resilient constrained learning

Definition (Resilient equilibrium)

For a strictly convex function $h(r)$, we say the relaxation r^* achieves the resilient equilibrium if

$$\nabla h(r^*) \in -\partial P^*(r^*) \quad \leftarrow (\partial: \text{subdifferential})$$

In words: at the resilient equilibrium the **marginal cost of relaxing** equals the **marginal gain of relaxing**



[Hounie, Chamon, Ribeiro, NeurIPS'23]

66

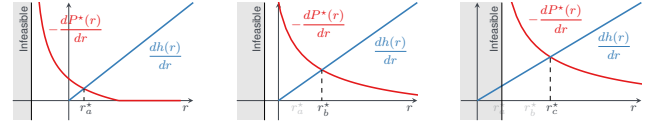
Resilient constrained learning

Definition (Resilient equilibrium)

For a strictly convex function $h(r)$, we say the relaxation r^* achieves the resilient equilibrium if

$$\nabla h(r^*) \in -\partial P^*(r^*) \quad \leftarrow (\partial: \text{subdifferential})$$

In words: at the resilient equilibrium the **marginal cost of relaxing** equals the **marginal gain of relaxing**



[Hounie, Chamon, Ribeiro, NeurIPS'23]

66

Resilient constrained learning

Definition (Resilient equilibrium)

For a strictly convex function $h(r)$, we say the relaxation r^* achieves the resilient equilibrium if

$$\nabla h(r^*) \in -\partial P^*(r^*) = \lambda^*(r^*)$$

In words: at the resilient equilibrium the **marginal cost of relaxing** equals the **marginal gain of relaxing**

- ✓ After relaxing, $\lambda^*(r^*)$ is *smaller* than $\lambda^*(0)$
 \Rightarrow Resilient constrained learning “generalizes better” (lower sample complexity)

[Hounie, Chamon, Ribeiro, NeurIPS'23]

67

Resilient constrained learning

Definition (Resilient equilibrium)

For a strictly convex function $h(r)$, we say the relaxation r^* achieves the resilient equilibrium if

$$\nabla h(r^*) \in -\partial P^*(r^*) = \lambda^*(r^*)$$

In words: at the resilient equilibrium the **marginal cost of relaxing** equals the **marginal gain of relaxing**

- ✓ After relaxing, $\lambda^*(r^*)$ is *smaller* than $\lambda^*(0)$
 \Rightarrow Resilient constrained learning “generalizes better” (lower sample complexity)
- ✓ The resilient equilibrium *exists and is unique* (because h is strictly convex)

[Hounie, Chamon, Ribeiro, NeurIPS'23]

67

Resilient constrained learning

Definition (Resilient equilibrium)

For a strictly convex function $h(r)$, we say the relaxation r^* achieves the resilient equilibrium if

$$\begin{aligned} P^*(r^*) &= \min_{\theta, r} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\text{Loss}(f_\theta(x), y)] + h(r) \\ \text{subject to } &\mathbb{E}_{(x,y) \sim \mathcal{Q}_i} [g_i(f_\theta(x_m), y_m)] \leq c_i + r_i \end{aligned}$$

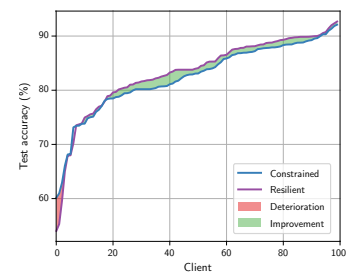
In words: at the resilient equilibrium the **marginal cost of relaxing** equals the **marginal gain of relaxing**

- ✓ After relaxing, $\lambda^*(r^*)$ is *smaller* than $\lambda^*(0)$
 \Rightarrow Resilient constrained learning “generalizes better” (lower sample complexity)
- ✓ The resilient equilibrium *exists and is unique* (because h is strictly convex)

[Hounie, Chamon, Ribeiro, NeurIPS'23]

67

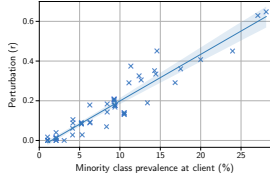
Heterogeneous federated learning



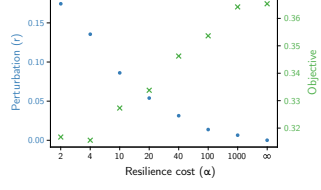
[Hounie, Chamon, Ribeiro, NeurIPS'23]

68

Heterogeneous federated learning

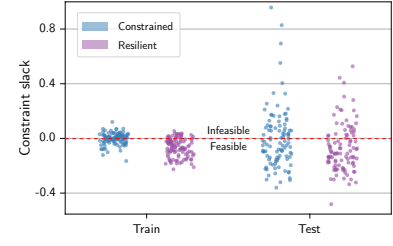


[Hounie, Charmon, Ribeiro, NeurIPS'23]



69

Heterogeneous federated learning



[Hounie, Charmon, Ribeiro, NeurIPS'23]

70

Agenda

Resilient constrained learning

Semi-infinite learning

Probabilistic robustness

71

Semi-infinite constrained learning

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

72

Semi-infinite constrained learning

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \\ \text{subject to} \quad & \text{Loss}(f_{\theta}(x_n + \delta), y_n) \leq t(x_n, y_n), \\ & \text{for all } (x_n, y_n) \text{ and } \delta \in \Delta \end{aligned}$$

• Epigraph formulation:

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

72

Semi-infinite constrained learning

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)]$$

$$\begin{aligned} \text{subject to} \quad & \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_e), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_z), y_n) \leq t(x_n, y_n) \end{aligned}$$

• Epigraph formulation:

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

• Semi-infinite program

$$\begin{aligned} & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^*}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^*}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{2\pi^*}), y_n) \leq t(x_n, y_n) \end{aligned}$$

72

Duality

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right] \\ \iff \\ \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \text{ s.t. } \text{Loss}(f_{\theta}(x_n + \delta), y_n) \leq t(x_n, y_n), \forall (x_n, y_n, \delta) \\ \iff \\ \min_{\theta} \sup_{\mu \in \mathcal{P}} \quad & \frac{1}{N} \sum_{n=1}^N \underbrace{\int_{\Delta} \mu_n(\delta) \text{Loss}(f_{\theta}(x_n + \delta), y_n) d\delta}_{L(\theta; \mu_n)} \end{aligned}$$

73

Duality

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right] \\ \iff \\ \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \text{ s.t. } \text{Loss}(f_{\theta}(x_n + \delta), y_n) \leq t(x_n, y_n), \forall (x_n, y_n, \delta) \\ \iff \\ \min_{\theta} \sup_{\mu \in \mathcal{P}} \quad & \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbb{E}_{\delta \sim \mu} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]}_{L(\theta; \mu_n)} \end{aligned}$$

73

From optimization to sampling

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

$$\approx \min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbb{E}_{\delta \sim \mu_y(\cdot | x_n, y_n)} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]}_{L(\theta, \mu)}$$

Proposition

For all $\epsilon > 0$, there exists $\gamma(x, y) < \max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x + \delta), y)$ s.t. $L(\theta, \mu_{\gamma}) \geq \sup_{\mu \in \mathcal{P}^2} L(\theta, \mu) - \epsilon$ for

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

74

From optimization to sampling

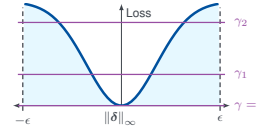
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

$$\approx \min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbb{E}_{\delta \sim \mu_y(\cdot | x_n, y_n)} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]}_{L(\theta, \mu)}$$

Proposition

For any approximation error, $\exists \gamma(x, y)$ such that

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



[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS'21]

75

From optimization to sampling

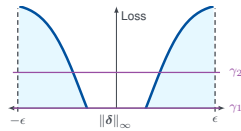
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

$$\approx \min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbb{E}_{\delta \sim \mu_y(\cdot | x_n, y_n)} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]}_{L(\theta, \mu)}$$

Proposition

For any approximation error, $\exists \gamma(x, y)$ such that

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



[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS'21]

75

From optimization to sampling

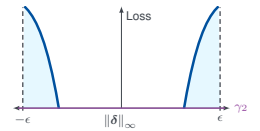
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

$$\approx \min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbb{E}_{\delta \sim \mu_y(\cdot | x_n, y_n)} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]}_{L(\theta, \mu)}$$

Proposition

For any approximation error, $\exists \gamma(x, y)$ such that

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



[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS'21]

75

From optimization to sampling

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

$$\stackrel{!}{=} \min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbb{E}_{\delta \sim \mu_y(\cdot | x_n, y_n)} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]}_{L(\theta, \mu)}$$

Proposition

For any approximation error, $\exists \gamma(x, y)$ such that

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



[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS'21]

75

From optimization to sampling

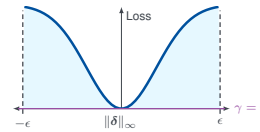
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

$$\approx \min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbb{E}_{\delta \sim \mu_y(\cdot | x_n, y_n)} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]}_{L(\theta, \mu)}$$

Proposition

For any approximation error, $\exists \gamma(x, y)$ such that

$$\mu_0(\delta | x, y) \propto \text{Loss}(f_{\theta}(x + \delta), y)$$



[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS'21]

75

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

$$\max_{\lambda \geq 0} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) + \lambda \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

⚙️ Balancing nominal accuracy and robustness \Rightarrow Dual constrained learning

- ❌ Computing the worst-case perturbations
 - gradient ascent \rightarrow non-convex, underparametrized

76

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

$$\max_{\lambda \geq 0} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) + \lambda \left[\mathbb{E}_{\delta \sim \mu_0(\cdot | x_n, y_n)} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)] \right]$$

⚙️ Balancing nominal accuracy and robustness \Rightarrow Dual constrained learning

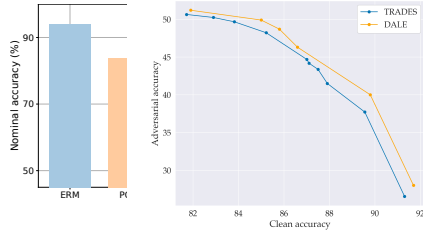
- ✅ Computing the worst-case perturbations
 - gradient ascent \rightarrow non-convex, underparametrized \Rightarrow sampling

76

Dual Adversarial Learning

Problem

Learn an image classifier that is robust to input perturbations



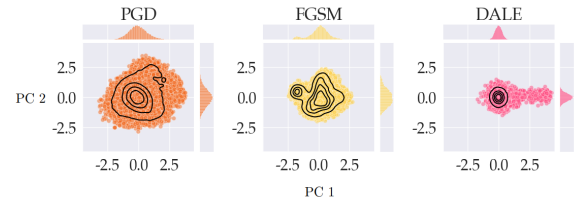
[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS21]

77

Dual Adversarial Learning

Problem

Learn an image classifier that is robust to input perturbations



[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS21]

78

Dual Adversarial Learning

```

1: for  $n = 1, \dots, N$ :
2:    $\delta_n \sim \text{Random}(\Delta)$ 
3:   for  $k = 1, \dots, K$ :
4:      $\zeta \sim \text{Laplace}(0, I)$ 
5:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss}(f_{\theta_t}(x_n + \delta_n), y_n) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss}(f_{\theta}(x_n), y_n) + \lambda \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) - c \right) \right]_+$ 

```

HMC sampling:
 $\delta \sim \mu_0(\cdot | x_n, y_n)$

SGD

GA

[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS21]

79

Dual Adversarial Learning

```

1: for  $n = 1, \dots, N$ :
2:    $\delta_n \sim \text{Random}(\Delta)$ 
3:   for  $k = 1, \dots, K$ :
4:      $\zeta \sim \text{Laplace}(0, I)$ 
5:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss}(f_{\theta_t}(x_n + \delta_n), y_n) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss}(f_{\theta}(x_n), y_n) + \lambda \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) - c \right) \right]_+$ 

```

HMC sampling:
 $\delta \sim \mu_0(\cdot | x_n, y_n)$

SGD

GA

[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS21]

79

Dual Adversarial Learning

```

1: for  $n = 1, \dots, N$ :
2:    $\delta_n \sim \text{Random}(\Delta)$ 
3:   for  $k = 1, \dots, K$ :
4:      $\zeta \sim \text{Laplace}(0, I)$ 
5:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss}(f_{\theta_t}(x_n + \delta_n), y_n) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss}(f_{\theta}(x_n), y_n) + \lambda \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) - c \right) \right]_+$ 

```

HMC sampling:
 $\delta \sim \mu_0(\cdot | x_n, y_n)$

SGD

GA

[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS21]

79

Dual Adversarial Learning

```

1: for  $n = 1, \dots, N$ :
2:    $\delta_n \sim \text{Random}(\Delta)$ 
3:   for  $k = 1, \dots, K$ :
4:      $\zeta \sim \text{Laplace}(0, I)$ 
5:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss}(f_{\theta_t}(x_n + \delta_n), y_n) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss}(f_{\theta}(x_n), y_n) + \lambda \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) - c \right) \right]_+$ 

```

HMC sampling:
 $\delta \sim \mu_0(\cdot | x_n, y_n)$

SGD

GA

[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS21]

79

Dual Adversarial Learning

```

1: for  $n = 1, \dots, N$ :
2:    $\delta_n \sim \text{Random}(\Delta)$ 
3:   for  $k = 1, \dots, K$ :
4:      $\zeta \sim \text{Laplace}(0, I)$ 
5:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss}(f_{\theta_t}(x_n + \delta_n), y_n) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss}(f_{\theta}(x_n), y_n) + \lambda \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) - c \right) \right]_+$ 

```

HMC sampling:
 $\delta \sim \mu(\cdot | x_n, y_n)$

SGD

GA

[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS21]

80

Dual Adversarial Learning

```

1: for  $n = 1, \dots, N$ :
2:    $\delta_n \sim \text{Random}(\Delta)$ 
3:   for  $k = 1, \dots, K$ :
4:      $\zeta \sim \text{Laplace}(0, I)$ 
5:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss}(f_{\theta_t}(x_n + \delta_n), y_n) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss}(f_{\theta}(x_n), y_n) + \lambda \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) - c \right) \right]_+$ 

```

Gaussian
[Lopes et al., arXiv/19]
[Rusak et al., ECCV/20]
Patches
[Zhong et al., AAAI/20]
[Yun et al., ICCV/19]
...

SGD

GA

[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS21]

80

Dual Adversarial Learning

```

1: for  $n = 1, \dots, N$ :
2:    $\delta_n \sim \text{Random}(\Delta)$ 
3:   for  $k = 1, \dots, K$ :
4:      $\zeta \sim \text{Laplace}(0, I)$ 
5:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss}(f_{\theta}(x_n), y_n) + \lambda \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n + \delta_n), y_n) - c \right) \right]_+$ 

```

$T \rightarrow 0$: "PGD"
[Szegedy et al., ICLR'14]
[Goodfellow et al., ICLR'15]
[Madry et al., ICLR'18]

SGD

GA

[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS'21]

80

Invariance

Problem

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

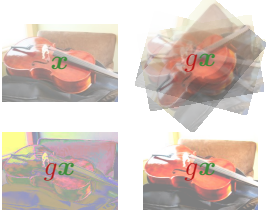


81

Invariance

Problem

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



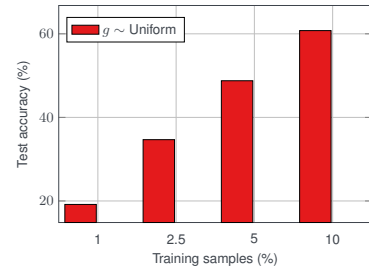
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{g \sim \mathcal{G}} \left[\text{Loss}(f_{\theta}(g(x_n)), y_n) \right]$$

$$\mathcal{G} = \left\{ \begin{array}{l} \bullet \text{ Identity} \\ \bullet \text{ ShearX(Y), Flip, Rotate, TranslateX(Y), Cutout, Crop} \\ \bullet \text{ AutoContrast, Invert, Equalize, Color, Solarize, Posterize, Contrast, Brightness, Sharpness} \end{array} \right\}$$

[Houie, Chamon, Ribeiro, ICML'23]

81

Training on a subset of ImageNet-100



[Houie, Chamon, Ribeiro, ICML'23]

82

Invariance

Problem

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

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{g \sim \mathcal{G}} \left[\text{Loss}(f_{\theta}(g(x_n)), y_n) \right]$$

$$\mathcal{G} = \left\{ \begin{array}{l} \bullet \text{ Identity} \\ \bullet \text{ ShearX(Y), Flip, Rotate, TranslateX(Y), Cutout, Crop} \\ \bullet \text{ AutoContrast, Invert, Equalize, Color, Solarize, Posterize, Contrast, Brightness, Sharpness} \end{array} \right\}$$

[Houie, Chamon, Ribeiro, ICML'23]

83

Invariance

Problem

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

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{g \in \mathcal{G}} \text{Loss}(f_{\theta}(g(x_n)), y_n) \right]$$

$$\mathcal{G} = \left\{ \begin{array}{l} \bullet \text{ Identity} \\ \bullet \text{ ShearX(Y), Flip, Rotate, TranslateX(Y), Cutout, Crop} \\ \bullet \text{ AutoContrast, Invert, Equalize, Color, Solarize, Posterize, Contrast, Brightness, Sharpness} \end{array} \right\}$$

[Houie, Chamon, Ribeiro, ICML'23]

83

Invariance

Problem

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

$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\mathbb{E}_{g \sim \mu_0(\cdot | x_n, y_n)} \left[\text{Loss}(f_{\theta}(g(x_n)), y_n) \right] \right]$$

$$\mathcal{G} = \left\{ \begin{array}{l} \bullet \text{ Identity} \\ \bullet \text{ ShearX(Y), Flip, Rotate, TranslateX(Y), Cutout, Crop} \\ \bullet \text{ AutoContrast, Invert, Equalize, Color, Solarize, Posterize, Contrast, Brightness, Sharpness} \end{array} \right\}$$

[Houie, Chamon, Ribeiro, ICML'23]

83

Invariance

Problem

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

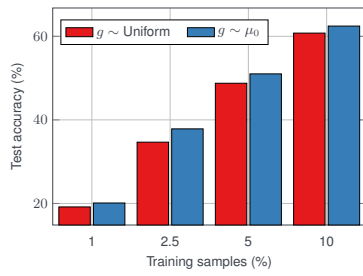
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to } \frac{1}{N} \sum_{n=1}^N \left[\mathbb{E}_{g \sim \mu_0(\cdot | x_n, y_n)} \left[\text{Loss}(f_{\theta}(g(x_n)), y_n) \right] \right] \leq c$$

$$\mathcal{G} = \left\{ \begin{array}{l} \bullet \text{ Identity} \\ \bullet \text{ ShearX(Y), Flip, Rotate, TranslateX(Y), Cutout, Crop} \\ \bullet \text{ AutoContrast, Invert, Equalize, Color, Solarize, Posterize, Contrast, Brightness, Sharpness} \end{array} \right\}$$

[Houie, Chamon, Ribeiro, ICML'23]

83

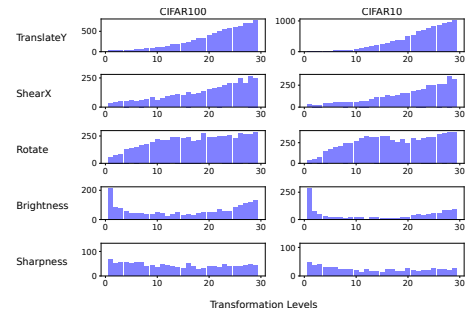
Training on a subset of ImageNet-100



[Hounie, Chamon, Ribeiro, ICML23]

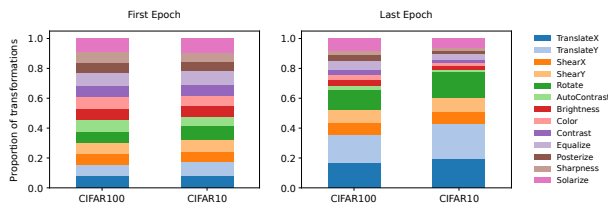
84

Not all transformations are created equal



85

Not all transformations are created equal



86

“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, ICML23]

87

Agenda

Resilient constrained learning

Semi-infinite learning

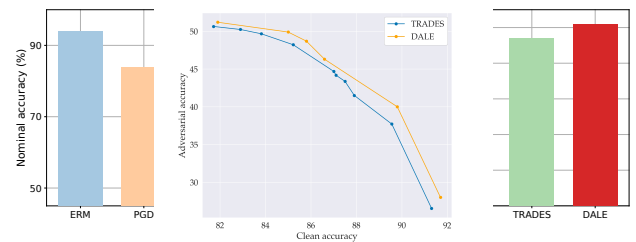
Probabilistic robustness

88

Constrained learning for robustness

Problem

Learn an **accurate** classifier



[Chamon and Ribeiro, NeurIPS20; Robey et al., NeurIPS21; Chamon et al., IEEE TIT23]

89

Constrained learning for robustness

Problem

Learn an **accurate** classifier that is (mostly) robust to input perturbations

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

[Chamon and Ribeiro, NeurIPS20; Robey et al., NeurIPS21; Chamon et al., IEEE TIT23]

90

“Softer” robustness

- Softmax or *log-sum-exp* [Li et al., ICLR21]

$$\min_{\theta} \mathbb{E}_{(x,y)} \left[\frac{1}{\tau} \log \left(\mathbb{E}_{\delta \sim m} \left[e^{\tau \cdot \text{Loss}(f_{\theta}(x+\delta), y)} \right] \right) \right]$$

- $\tau \rightarrow 0$: classical learning (with randomized data augmentation)
- $\tau \rightarrow \infty$: adversarial robustness (ess sup)

- L_p norms [Pice et al., NeurIPS21]

$$\min_{\theta} \mathbb{E}_{(x,y)} \left[\mathbb{E}_{\delta \sim m} \left[\left| \text{Loss}(f_{\theta}(x+\delta), y) \right|^{\tau} \right]^{\frac{1}{\tau}} \right]$$

- $\tau = 1$: classical learning (with randomized data augmentation)
- $\tau \rightarrow \infty$: adversarial robustness (ess sup)

91

“Softer” robustness

- Softmax or *log-sum-exp* [Li et al., ICLR'21]

$$\min_{\theta} \mathbb{E}_{(x,y)} \left[\frac{1}{\tau} \log \left(\mathbb{E}_{\delta \sim m} \left[e^{\tau \cdot \text{Loss}(f_{\theta}(x+\delta), y)} \right] \right) \right]$$

- L_p norms [Rice et al., NeurIPS'21]

$$\min_{\theta} \mathbb{E}_{(x,y)} \left[\mathbb{E}_{\delta \sim m} \left[\left| \text{Loss}(f_{\theta}(x+\delta), y) \right|^{\tau} \right]^{1/\tau} \right]$$

- Computationally challenging (especially as $\tau \rightarrow \infty$, i.e., stronger robustness)
- No guaranteed advantages (lower sample complexity? improved trade-offs?)

91

Towards probabilistic robustness

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \left[\ell(x_n, y_n) \right] \\ \text{subject to} \quad & \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_1), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_z), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_4), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi/2}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^*}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^*}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{2\pi}), y_n) \leq \ell(x_n, y_n) \end{aligned}$$

92

Towards probabilistic robustness

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \left[\ell(x_n, y_n) \right] \\ \text{subject to} \quad & \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_1), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_z), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_4), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi/2}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^*}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^*}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{2\pi}), y_n) \leq \ell(x_n, y_n) \end{aligned}$$

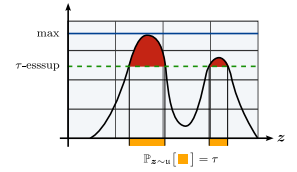
92

Probabilistic robustness

- Probabilistic robustness

$$\min_{\theta} \mathbb{E}_{(x,y)} \left[\tau\text{-esssup}_{\delta \in \Delta} \text{Loss}(f_{\theta}(x+\delta), y) \right]$$

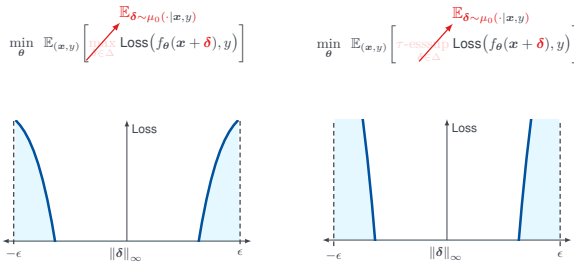
- $\tau = 1/2$: classical learning (for symmetric m)
- $\tau = 0$: adversarial robustness (ess sup)



[Robey, Chamon, Pappas, Hassani, ICML'22 (spotlight)]

93

Probabilistic robustness



[Robey, Chamon, Pappas, Hassani, ICML'22 (spotlight)]

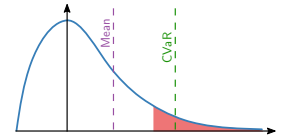
94

Probabilistic robustness and Risk

- Conditional value at risk:

$$\begin{aligned} \text{CVaR}_{\rho}(f) &= \mathbb{E}_z [f(z) \mid f(z) \geq F_z^{-1}(\rho)] \\ &= \inf_{\alpha \in \mathbb{R}} \alpha + \frac{\mathbb{E}_z [f(z) - \alpha]_+}{1 - \rho} \end{aligned}$$

- $\text{CVaR}_0(f) = \mathbb{E}_z [f(z)]$
- $\text{CVaR}_1(f) = \text{ess sup}_z f(z)$



Proposition

CVaR is the tightest convex upper bound of τ -esssup, i.e.,
 $\tau\text{-esssup}_z f(z) \leq \text{CVaR}_{1-\tau}(f)$ with equality when $\rho = 0$ or $\rho = 1$.

[Shapiro et al. Lectures on Stochastic Programming, 2014; Kalogieras et al., IEEE ICASSP'20]

95

Probabilistically robust learning

```

1: for  $n = 1, \dots, N$ :
2:    $\alpha_0 = 0$ 
3:   for  $t = 1, \dots, T$ :
4:      $\delta_t \sim \text{Random}(\Delta)$ 
5:      $\alpha \leftarrow \alpha - \frac{\eta}{\tau} \left( \tau - \mathbb{I} [\text{Loss}(f_{\theta}(x_n + \delta_t), y_n) \geq \alpha] \right)$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss}(f_{\theta}(x_n + \delta_T), y_n) - \alpha \right]_+$ 
8: end
    
```

SGD (CVaR)

SGD (θ)

[Robey, Chamon, Pappas, Hassani, ICML'22 (spotlight)]

96

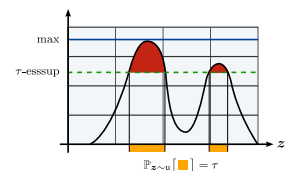
Probabilistic robustness

- Probabilistic robustness

$$\min_{\theta} \mathbb{E}_{(x,y)} \left[\tau\text{-esssup}_{\delta \in \Delta} \text{Loss}(f_{\theta}(x+\delta), y) \right]$$

- $\tau = 1/2$: classical learning (for symmetric m)
- $\tau = 0$: adversarial robustness (ess sup)

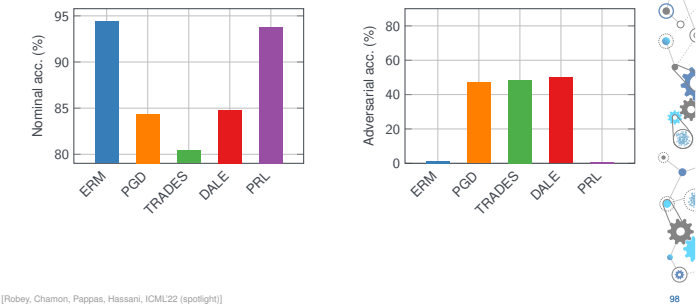
- Potentially better sample complexity
[Robey et al., ICML'22 (spotlight)] ✓
[Raman et al., NeurIPS ML Safety Workshop'22] ✓
- Better performance trade-off
[Robey et al., ICML'22 (spotlight)] ✓



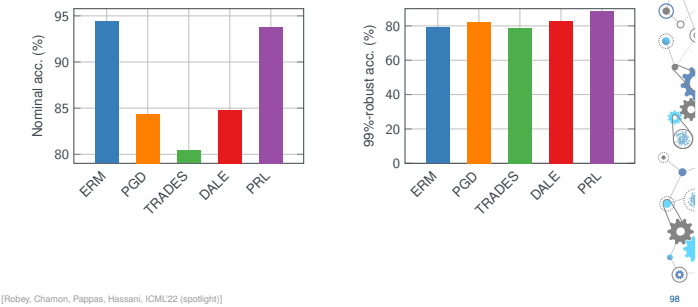
[Robey, Chamon, Pappas, Hassani, ICML'22 (spotlight)]

97

Probabilistically robust learning



Probabilistically robust learning



Summary

- Semi-infinite constrained learning is ~~the~~ a tool to enforce worst-case requirements
- Semi-infinite constrained learning...
- ...but possible. How?

Summary

- Semi-infinite constrained learning is ~~the~~ a tool to enforce worst-case requirements
e.g., robustness [Robey et al., NeurIPS'21], invariance [Hourie et al., ICML'23], smoothness [Cervino et al., ICML'23], ...
- Semi-infinite constrained learning...
- ...but possible. How?

Summary

- Semi-infinite constrained learning is ~~the~~ a tool to enforce worst-case requirements
e.g., robustness [Robey et al., NeurIPS'21], invariance [Hourie et al., ICML'23], smoothness [Cervino et al., ICML'23], ...
- Semi-infinite constrained learning...
Learning problem with an infinite number of constraints
- ...but possible. How?

Summary

- Semi-infinite constrained learning is ~~the~~ a tool to enforce worst-case requirements
e.g., robustness [Robey et al., NeurIPS'21], invariance [Hourie et al., ICML'23], smoothness [Cervino et al., ICML'23], ...
- Semi-infinite constrained learning...
Learning problem with an infinite number of constraints
- ...but possible. How?
Using a hybrid sampling–optimization algorithm or, in the case of probabilistic robustness, a *tight* convex relaxation (CVaR) [Robey et al., ICML'22]

Agenda

- I. Constrained supervised learning
 - Constrained learning theory
 - Constrained learning algorithms
 - Resilient constrained learning
- Break (10 min)
- II. Constrained reinforcement learning
 - Constrained RL duality
 - Constrained RL algorithms
- Q&A and discussions



<https://luizchamon.com/sgm>