

IMPRS tutorial
Sep. 19, 2024

supervised and reinforcement learning under requirements

Luiz F. O. Chamon



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 (German = ε)

Acknowledgment: this material is based on a tutorial prepared in collaboration with Miguel Calvo-Fullana (UPF), Santiago Paternain (RPI), and Alejandro Ribeiro (Penn).

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

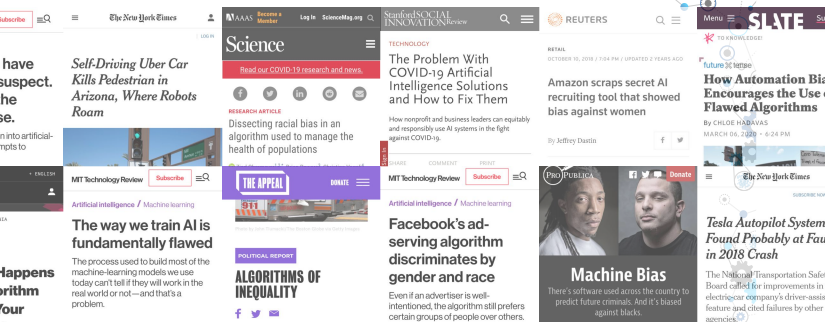
Q&A and discussions



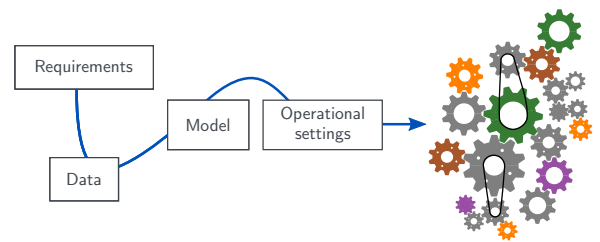
<https://luizchamon.com/imprs2024>

Why learning under requirements?

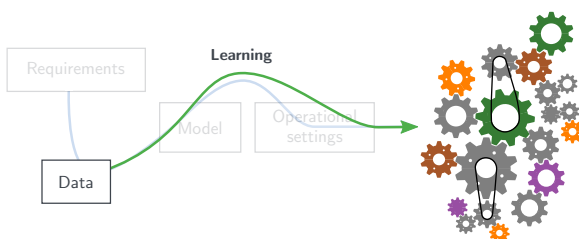
Why learning under requirements?



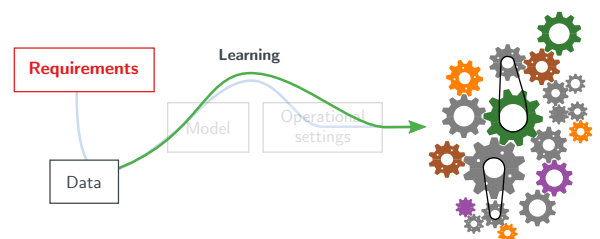
Why learning under requirements?



Why learning under requirements?

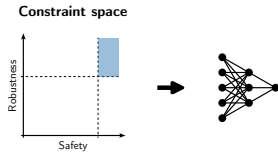


Why learning under requirements?



What is a requirements?

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

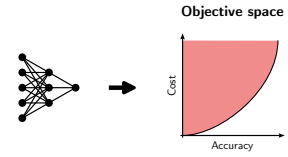


[NASA, "Systems engineering handbook," 2019]

4

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

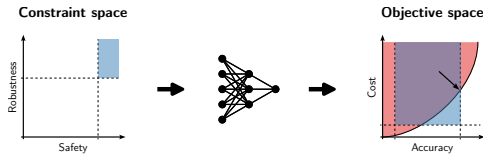


[NASA, "Systems engineering handbook," 2019]

4

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]

4

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]
- $\mathcal{D}, \mathcal{A}, \mathcal{U}$ unknown

[Chamon et al., IEEE ICASSP20 (best student paper); Chamon and Ribeiro, NeurIPS'20; Chamon et al., IEEE TIT'23]

5

What is (un)constrained learning?

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

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

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

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

[Chamon et al., IEEE ICASSP20 (best student paper); Chamon and Ribeiro, NeurIPS'20; Chamon et al., IEEE TIT'23]

5

What about penalties?

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

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

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

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

[Chamon et al., IEEE ICASSP20 (best student paper); Chamon and Ribeiro, NeurIPS'20; Chamon et al., IEEE TIT'23]

6

What about penalties?

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

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

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

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

- 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...

6

Applications

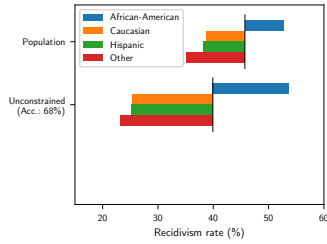
- Fairness (e.g., [Goh et al., NeurIPS'16; Kearns et al., ICML'18; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23])
- Federated learning (e.g., [Shen et al., ICLR'22; Hounie et al., NeurIPS'23])
- Adversarially robust learning (e.g., [Chamon et al., NeurIPS'20; Robey et al., NeurIPS'21; Chamon et 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**

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

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

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

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

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

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

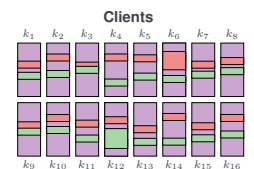
10

Federated learning

Problem

Learn a common model using data from K clients

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



- 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]

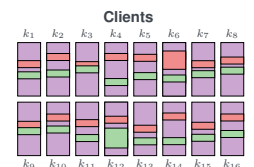
11

Heterogeneous federated learning

Problem

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

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



- 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

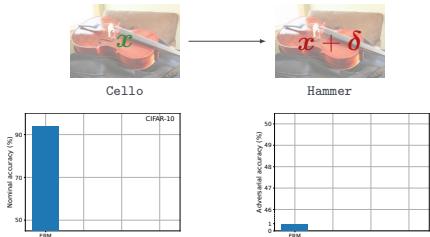
Applications

- 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; Hourie 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'21; 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



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

[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



$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} \quad & \text{Robustness loss} \leq c \end{aligned}$$

[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



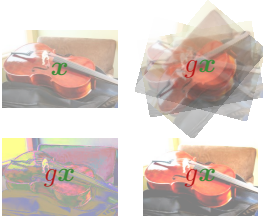
$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} \quad & \frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right] \leq c \end{aligned}$$

[Chamon and Ribeiro, 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., $\mathcal{G} = \left\{ \begin{array}{l} \text{Rotate, Translate(X(Y),} \\ \text{ShearX(Y), Crop, Invert,} \\ \text{Solarize, Contrast,} \\ \text{Brightness, Sharpness...} \end{array} \right\}$



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

[Hourie, Chamon, Ribeiro, NeurIPS'23]

14

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; Hourie 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'21; 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



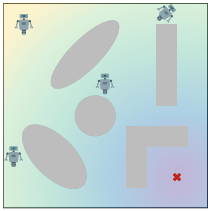
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

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

17

Safety

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



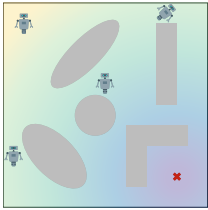
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

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 $\mathbb{P}[\text{Colliding with } \mathcal{O}_i] \leq \delta,$
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 $\mathbb{P} \left(\bigcap_{t=0}^{T-1} \{s_t \notin \mathcal{O}_i\} \mid \pi \right) \geq 1 - \delta_i,$
for $i = 1, 2, \dots$

17

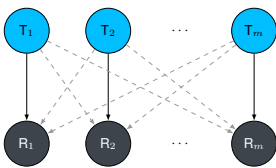
Applications

- 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; Hourie 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

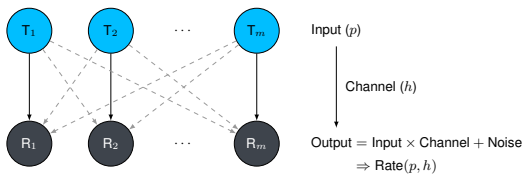
Problem
Allocate the least transmit power to m devices to achieve a communication rate



19

Wireless resource allocation

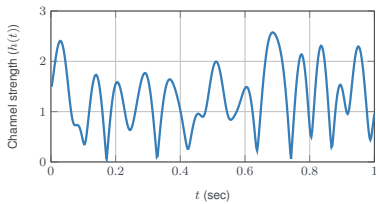
Problem
Allocate the least transmit power to m devices to achieve a communication rate



19

Wireless resource allocation

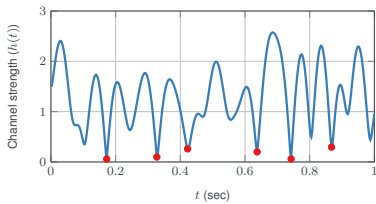
Problem
Allocate the least transmit power to m devices to achieve a communication rate



20

Wireless resource allocation

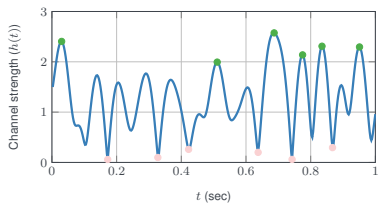
Problem
Allocate the least transmit power to m devices to achieve a communication rate



20

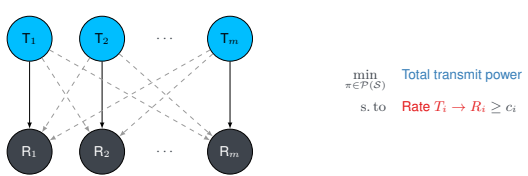
Wireless resource allocation

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



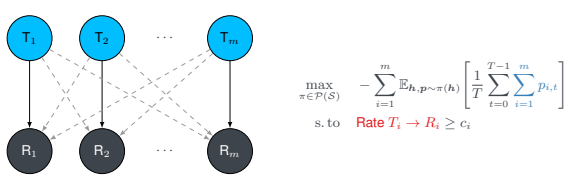
Wireless resource allocation

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



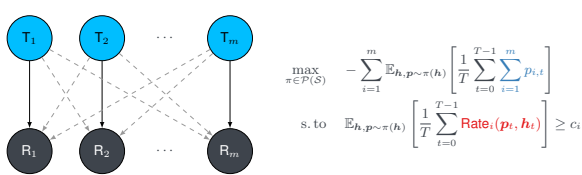
Wireless resource allocation

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



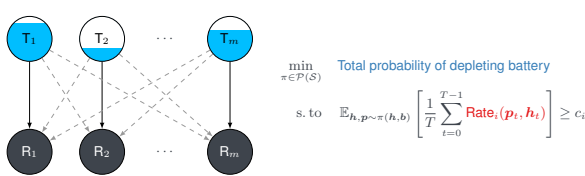
Wireless resource allocation

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



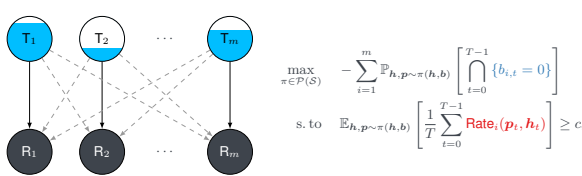
Wireless resource allocation

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



Wireless resource allocation

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



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, arXiv'24])
- 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...

Constrained supervised learning

What is (un)constrained learning?

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

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

$$h(f_{\theta}(x_r), y_r) \leq u, \quad r = 1, \dots, N$$

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

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

24

Constrained learning challenges

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

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

$$h(f_{\theta}(x_r), y_r) \leq u$$

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

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

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

Challenges

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

25

Constrained learning challenges

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

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

$$h(f_{\theta}(x_r), y_r) \leq u$$

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

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

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

Challenges

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

25

Constrained learning challenges

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

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

$$h(f_{\theta}(x_r), y_r) \leq u$$

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

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

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

Challenges

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

25

Constrained learning challenges

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

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

$$h(f_{\theta}(x_r), y_r) \leq u$$

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

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

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

Challenges

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

25

Agenda

Constrained learning theory

Constrained learning algorithms

Resilient constrained learning

26

What classical learning theory says?

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

- ✓ 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]

27

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^*} from $N_f(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,

- near-optimal

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

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT23]

28

What's in a solution?

Definition (PACC learnability)

f_θ is a *probably approximately correct* **(PACC)** learnable if for every ϵ, δ and every distributions \mathcal{D}, \mathcal{Q} , 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

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



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

28

When is constrained learning possible?

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

Proposition

f_θ is PAC learnable $\nRightarrow f_\theta$ is PACC learnable

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

29

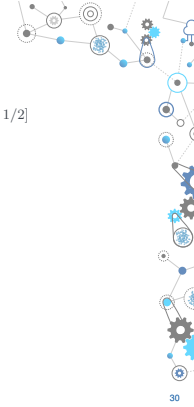
ECRM is not a PACC learner

Counter-example

$$\begin{aligned} P^* &= \min_{\theta \in \Theta} J(\theta) \\ \text{subject to } \theta_2 \mathbb{E}_{\tau}[\tau] &\leq \theta_1 - 1 \\ &\quad - \theta_1 \mathbb{E}_{\tau}[\tau] \leq \theta_2 - 1 \end{aligned}$$

$$J(\theta) = \begin{cases} 1/16, & \theta = [1/2, 1/2] \\ 1/8, & \theta = [1, 1] \\ 1/4, & \theta = [1, 0] \end{cases}$$

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



30

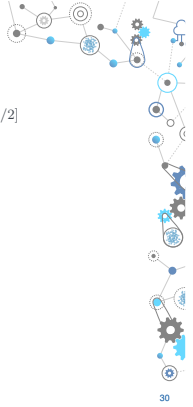
ECRM is not a PACC learner

Counter-example

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

$$J(\theta) = \begin{cases} 1/16, & \theta = [1/2, 1/2] \\ 1/8, & \theta = [1, 1] \\ 1/4, & \theta = [1, 0] \end{cases}$$

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



30

ECRM is not a PACC learner

Counter-example

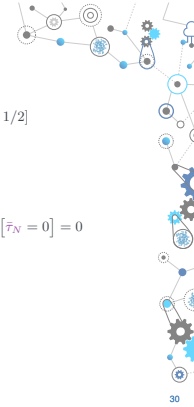
$$\begin{aligned} P^* &= \min_{\theta \in \Theta} J(\theta) = \frac{1}{8} \\ \text{subject to } \theta_2 \mathbb{E}_{\tau}[\tau] &\leq \theta_1 - 1 \Rightarrow \theta_1 \geq 1 \\ &\quad - \theta_1 \mathbb{E}_{\tau}[\tau] \leq \theta_2 - 1 \Rightarrow \theta_2 \leq 1 \end{aligned}$$

$$J(\theta) = \begin{cases} 1/16, & \theta = [1/2, 1/2] \\ 1/8, & \theta = [1, 1] \\ 1/4, & \theta = [1, 0] \end{cases}$$

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

$$\mathbb{P} [|\hat{P}^* - P^*| \leq 1/32] = \mathbb{P} [\bar{\tau}_N = 0] = 0$$

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



30

ECRM is not a PACC learner

Counter-example

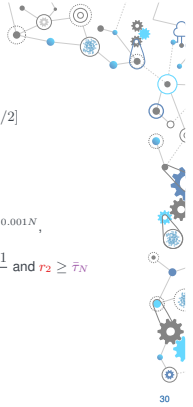
$$\begin{aligned} P^* &= \min_{\theta \in \Theta} J(\theta) = \frac{1}{8} \\ \text{subject to } \theta_2 \mathbb{E}_{\tau}[\tau] &\leq \theta_1 - 1 \Rightarrow \theta_1 \geq 1 \\ &\quad - \theta_1 \mathbb{E}_{\tau}[\tau] \leq \theta_2 - 1 \Rightarrow \theta_2 \leq 1 \end{aligned}$$

$$J(\theta) = \begin{cases} 1/16, & \theta = [1/2, 1/2] \\ 1/8, & \theta = [1, 1] \\ 1/4, & \theta = [1, 0] \end{cases}$$

$$\begin{aligned} \hat{P}^* &= \min_{\theta \in \Theta} J(\theta) \\ \text{subject to } \theta_2 \bar{\tau}_N &\leq \theta_1 - 1 + \mathbf{r}_1 \\ &\quad - \theta_1 \bar{\tau}_N \leq 1 - \theta_2 + \mathbf{r}_2 \end{aligned}$$

$$\mathbb{P} [|\hat{P}^* - P^*| \leq 1/32] \leq 4e^{-0.001N}, \text{ unless } \bar{\tau}_N \leq \mathbf{r}_1 < \frac{\bar{\tau}_N + 1}{2} \text{ and } \mathbf{r}_2 \geq \bar{\tau}_N$$

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



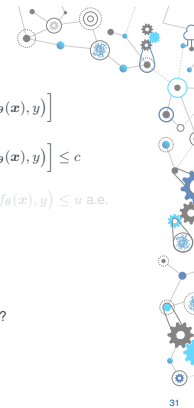
30

Constrained learning challenges

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

Challenges

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



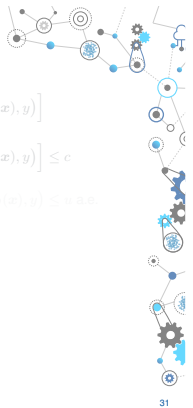
31

Constrained learning challenges

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

Challenges

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



31

Duality

PRIMAL
↕
DUAL

32

Duality

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

32

Duality

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

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

32

Duality

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

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

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

32

Duality

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

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

32

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

An alternative path

$$\hat{P}^* = \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, \mathbf{z}_n) \text{ s.t. } \frac{1}{N} \sum_{n=1}^N g(f_{\theta}, \mathbf{z}_n) \leq c$$

$$\hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, \mathbf{z}_n) + \lambda \left(\frac{1}{N} \sum_{n=1}^N g(f_{\theta}, \mathbf{z}_n) - c \right)$$

PAC

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

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

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

33

An alternative path

$$\hat{P}^* = \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, \mathbf{z}_n) \text{ s.t. } \frac{1}{N} \sum_{n=1}^N g(f_{\theta}, \mathbf{z}_n) \leq c$$

$$\hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, \mathbf{z}_n) + \lambda \left(\frac{1}{N} \sum_{n=1}^N g(f_{\theta}, \mathbf{z}_n) - c \right)$$

PAC

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

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

$$\mathcal{H}_{\theta} \subset \mathcal{H}$$

$$\hat{P}^* = \min_{\phi \in \mathcal{H}} \mathbb{E}_{\mathbf{z}} [\ell(\phi, \mathbf{z})] \text{ s.t. } \mathbb{E}_{\mathbf{z}} [g(\phi, \mathbf{z})] \leq c$$

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

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

33

Non-convex variational duality

Convex optimization: Primal ↔ Dual

Non-convex, finite dimensional optimization: Primal ↔ Dual

34

Non-convex variational duality



[Chamon, Eldar, Ribeiro, IEEE TSP20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

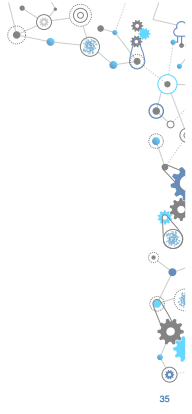
34

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

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

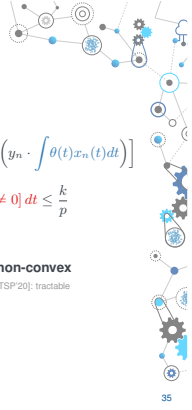
$$\text{s. to } \|\theta\|_0 = \sum_{i=1}^p \mathbb{I}[\theta_i \neq 0] \leq k$$

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 TSP20]: tractable



35

Sparse logistic regression

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

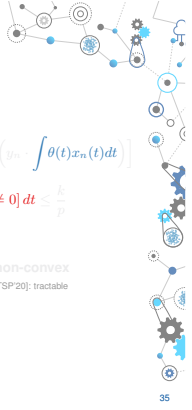
$$\text{s. to } \|\theta\|_0 = \sum_{i=1}^p \mathbb{I}[\theta_i \neq 0] \leq k$$

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 TSP20]: tractable



35

An alternative path

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

PAC

$$P^* = \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_{\theta}, z)] \quad \text{s. to } \mathbb{E}_z [g(f_{\theta}, z)] \leq c$$

=

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

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

36

An alternative path

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

PAC

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

$\uparrow \epsilon_0$ $\downarrow \epsilon_0$

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

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

36

An alternative path

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

PAC

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

$\uparrow \epsilon_0$ $\downarrow \epsilon_0$ $O(\epsilon)$

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

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

36

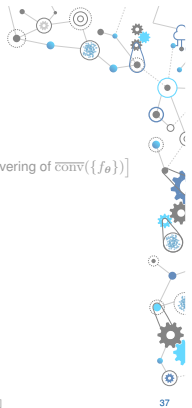
Dual (near-)PACC learning

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 $\overline{\text{conv}}(\{f_{\theta}\})$]



37

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

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

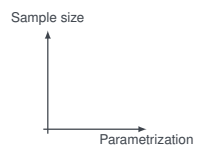
parametrization richness (ν) sample size (N) requirements difficulty (λ^*)

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

38

Dual learning trade-offs

- Unconstrained learning
- parametrization \times sample size

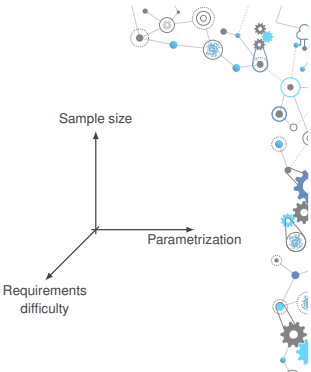


[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

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 TIT23]

39

When is constrained learning possible?

Corollary

f_{θ} is PAC learnable \approx^* f_{θ} is PAC $\color{red}{C}$ learnable

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

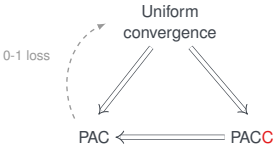
[mild conditions apply]

[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT23]

40

When is constrained learning possible?

Corollary



[mild conditions apply]

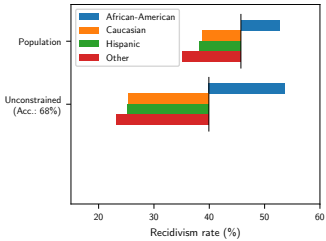
[Chamon and Ribeiro, NeurIPS20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT23]

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**

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

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

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Cotter et al., JMLR19; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT23]

42

Fairness: "Equality" of odds

Problem

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

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

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

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Cotter et al., JMLR19; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT23]

42

Fairness: "Equality" of odds

Problem

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

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

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

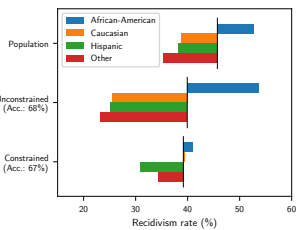
* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Cotter et al., JMLR19; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT23]

42

Fairness: "Equality" of odds

Problem

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

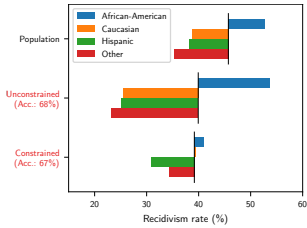


* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT23]

43

Fairness: “Equality” of odds

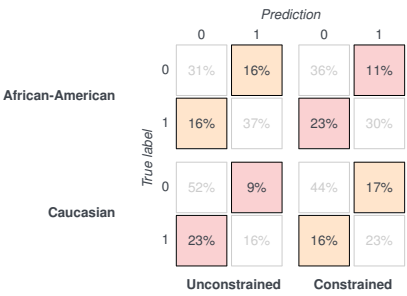
Problem
Predict whether an individual will recidivate **at the same rate across races**



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

43

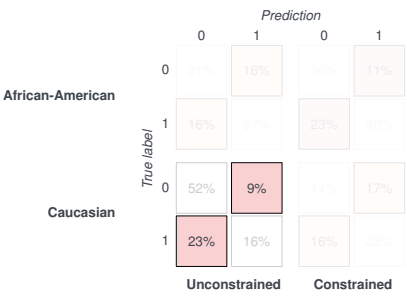
Fairness: “Equality” of odds



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

44

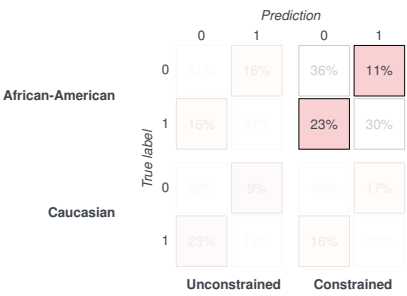
Fairness: “Equality” of odds



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

44

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]

44

Agenda

Constrained learning theory

Constrained learning algorithms

Resilient constrained learning

45

Constrained optimization methods

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

46

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

46

Constrained optimization methods

- Feasible update methods
e.g., conditional gradients (Frank-Wolfe)
 - ✗ Tractability [non-convex constraints]
 - ✓ Feasible candidate solution
 - Interior point methods
e.g., barriers, projection, polyhedral approx.
 - ✗ Tractability [non-convex constraints]
 - ✓ Feasible candidate solution
 - Duality
e.g., (augmented) Lagrangian
 - ✓ Tractability
 - ✓ (near-)feasible solution [small duality gap]
- $$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) \\ \text{subject to } &\frac{1}{N} \sum_{m=1}^N g(f_{\theta}(x_m), y_m) \leq c \\ &h(f_{\theta}(x_r), y_r) \leq u \end{aligned}$$

46

Dual learning algorithm

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

Dual learning algorithm

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

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

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

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

A (near-)PACC learner

Theorem

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

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

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

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

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

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

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

In practice...

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

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

- Update the dual

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

Dual learning algorithm

- Minimize the primal (\equiv ERM)

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

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

Dual learning algorithm

- Update the dual

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

In practice...

- Minimize the primal (\equiv ERM)

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

- Update the dual

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

In practice...

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

In practice...

```
1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_t \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, N$ 
5:      $\beta_{n+1} \leftarrow \beta_n - \eta_\theta \nabla_{\beta} [\ell(f_{\beta_n}(\mathbf{x}_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(\mathbf{x}_n), y_n)]$ 
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>

50

In practice...

```
1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_t \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, N$ 
5:      $\beta_{n+1} \leftarrow \beta_n - \eta_\theta \nabla_{\beta} [\ell(f_{\beta_n}(\mathbf{x}_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(\mathbf{x}_n), y_n)]$ 
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>

50

In practice...

```
1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_t \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, N$ 
5:      $\beta_{n+1} \leftarrow \beta_n - \eta_\theta \nabla_{\beta} [\ell(f_{\beta_n}(\mathbf{x}_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(\mathbf{x}_n), y_n)]$ 
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 ($\eta_\lambda = 0.1\eta_\theta$)

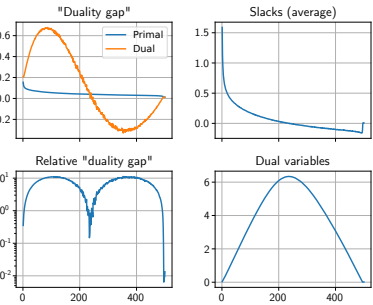


<https://github.com/lfochamon/csl>

50

In practice...

```
1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_t \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, l$ 
5:      $\beta_{n+1} \leftarrow \beta_n$ 
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:  $\theta_T, \lambda_T$ 
```



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



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

50

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

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

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

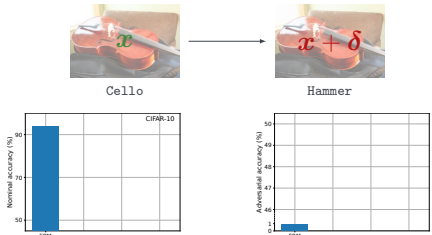
- Parameter: c (requirement-dependent)
- Generalizes with respect to Loss and $\text{Penalty} \leq c$

51

Robust learning

Problem

Learn an accurate classifier that is robust to input perturbations



52

Adversarial training

Problem

Learn an accurate classifier that is robust to input perturbations

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

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

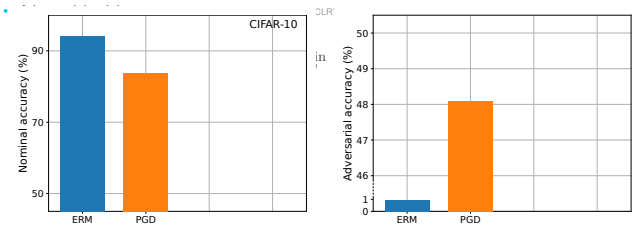


53

Adversarial training

Problem

Learn an accurate classifier that is robust to input perturbations



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

53

Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations

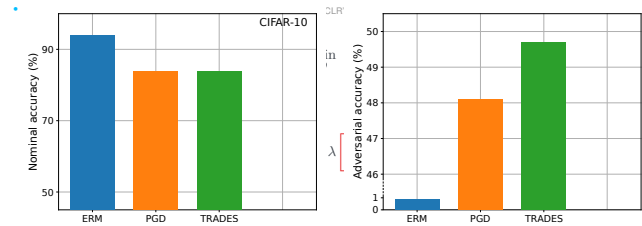
- Adversarial training [Szegedy et al., ICLR'14; Goodfellow et al., ICLR'15; Madry et al., ICLR'18, ...]

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

54

Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations



[Zhang et al., ICML'19]

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

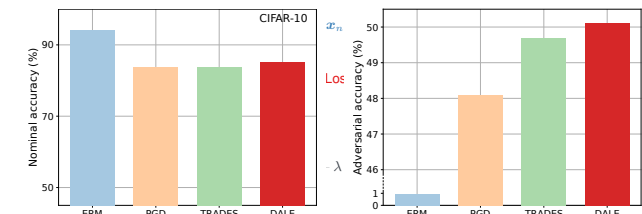
$$\frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right] \leq c$$

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

55

Constrained learning for robustness

Problem
Learn an accurate classifier that is robust to input perturbations

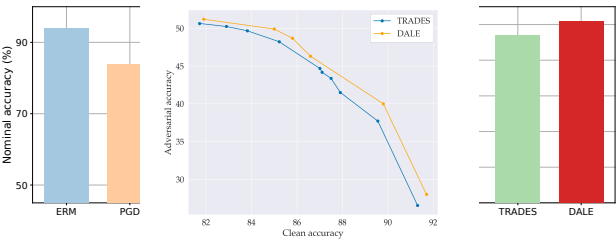


[Chamon and Ribeiro, NeurIPS'20; Robey et al., NeurIPS'21; Chamon et al., IEEE TIT'23]

55

Constrained learning for robustness

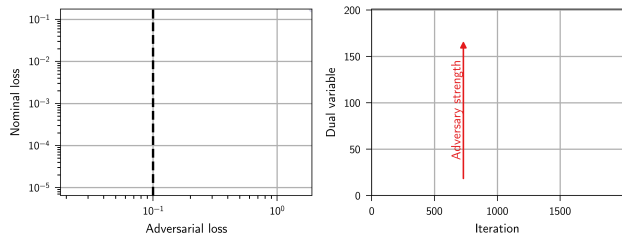
Problem
Learn an accurate classifier



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

55

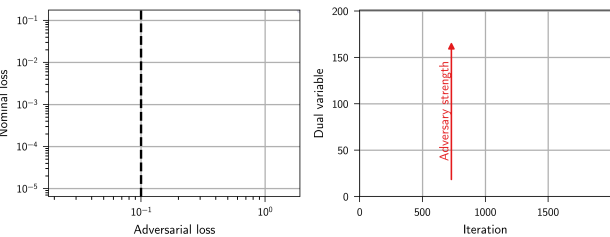
Constrained learning for robustness



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

56

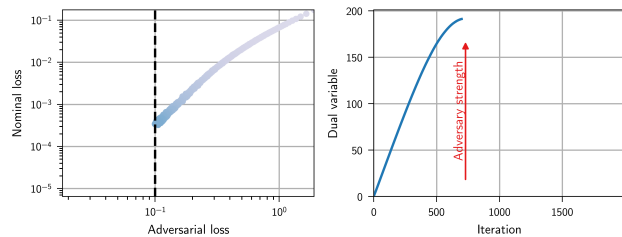
Constrained learning for robustness



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

56

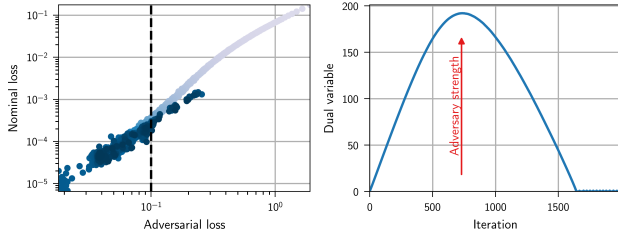
Constrained learning for robustness



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

56

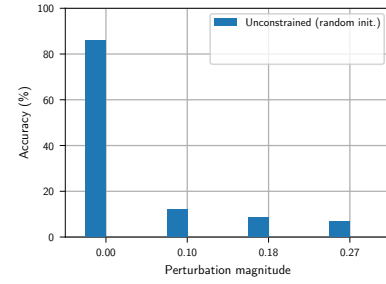
Constrained learning for robustness



Empirical observations: [Zhang et al., ICML20; Sitawarin, arXiv20]

56

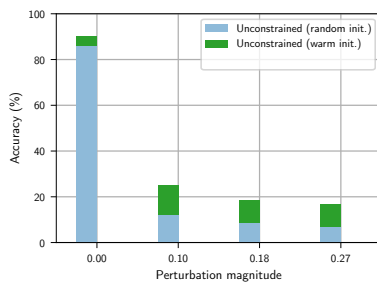
Constrained learning for robustness



[Chamion et al., IEEE TIT'23]

57

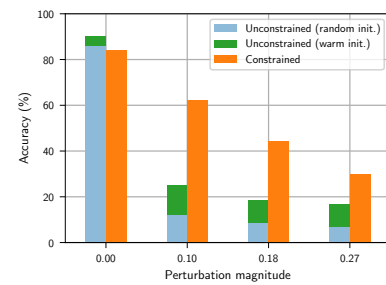
Constrained learning for robustness



[Chamion et al., IEEE TIT'23]

57

Constrained learning for robustness



[Chamion et al., IEEE TIT'23]

57

Penalty-based vs. dual learning

Penalty-based learning

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

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

Dual learning

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

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

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

58

Agenda

Constrained learning theory

Constrained learning algorithms

Resilient constrained learning

59

Heterogeneous federated learning

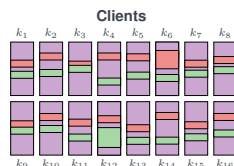
Problem

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

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

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

$$k = 1, \dots, K$$



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

60

Heterogeneous federated learning

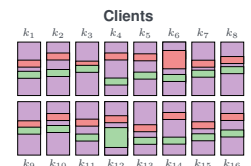
Problem

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

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

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

$$k = 1, \dots, K$$



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

60

Resilient constrained learning

Definition (Resilience)

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

Resilient constrained learning

Definition (Resilience)

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

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

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

Resilient constrained learning

Definition (Resilience)

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

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

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

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

Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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

Resilient constrained learning

Definition (Resilience)

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

Resilient constrained learning

Definition (Resilience)

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

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

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

Resilient constrained learning

Definition (Resilience)

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

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

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

- 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!

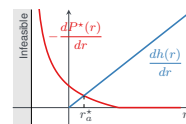
Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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



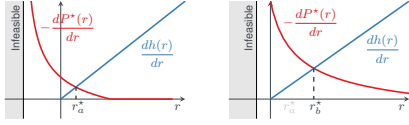
Resilient constrained learning

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, Charmon, Ribeiro, NeurIPS'23]

62

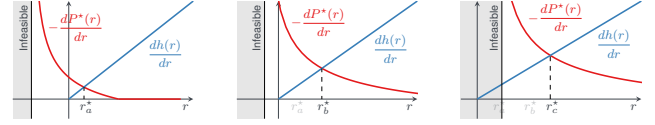
Resilient constrained learning

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, Charmon, Ribeiro, NeurIPS'23]

62

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, Charmon, Ribeiro, NeurIPS'23]

63

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, Charmon, Ribeiro, NeurIPS'23]

63

Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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

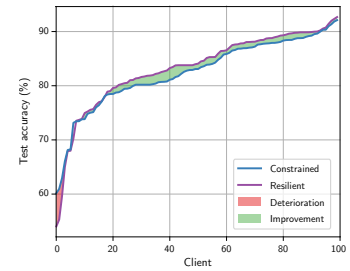
In words: at the resilient equilibrium the **marginal cost of relaxing** equals the **marginal gain of relaxing**

- 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, Charmon, Ribeiro, NeurIPS'23]

63

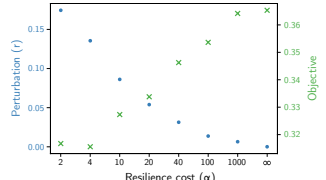
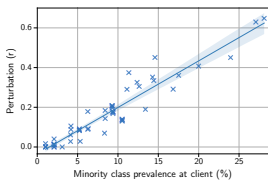
Heterogeneous federated learning



[Hounie, Charmon, Ribeiro, NeurIPS'23]

64

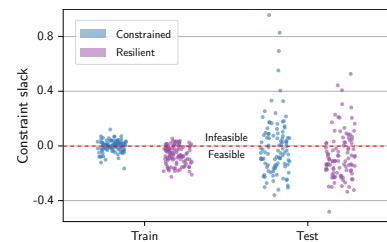
Heterogeneous federated learning



[Hounie, Charmon, Ribeiro, NeurIPS'23]

65

Heterogeneous federated learning



[Hounie, Charmon, Ribeiro, NeurIPS'23]

66

Summary

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

67

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., ICRL'22]. ...
- Constrained learning is hard...
- ...but possible. How?

67

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., ICRL'22]. ...
- Constrained learning is hard...
Constrained, non-convex, statistical optimization problem
- ...but possible. How?

67

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., ICRL'22]. ...
- 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*. Resilient learning can then be used to adapt the requirements to the task difficulty [Hourie et al., NeurIPS'23]

67

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/imprs2024>

68