

Miguel Calvo-Fullana
Universitat Pompeu Fabra, Spain

Luiz F. O. Chamon
Universität Stuttgart, Germany

Santiago Paternain
Rensselaer Polytechnic Institute, USA

Alejandro Ribeiro
University of Pennsylvania, USA

AAAI tutorial
Feb. 20, 2023

supervised and reinforcement learning under requirements

Agenda

I. Constrained supervised learning

II. Robustness-constrained learning

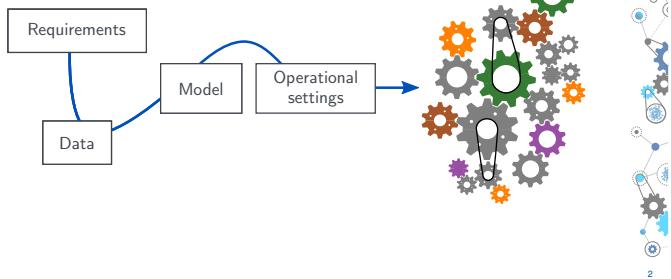
Break (30 min)

III. Constrained reinforcement learning

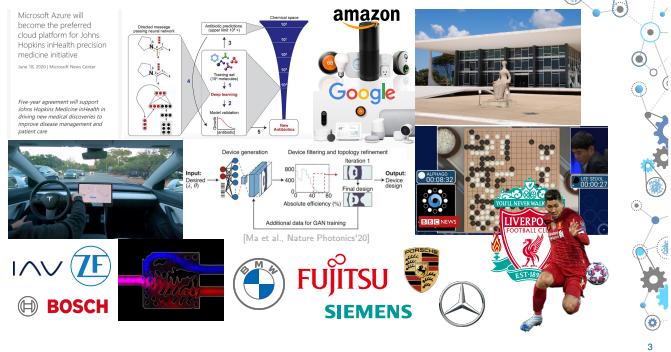


<https://luizchamon.com/aaai>

Why requirements?



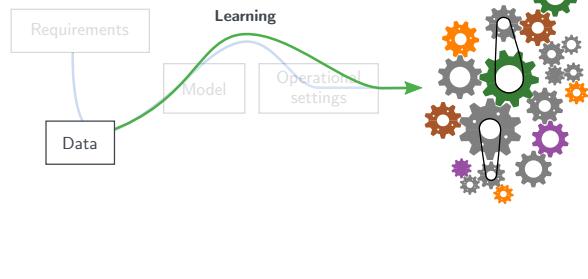
Why requirements?



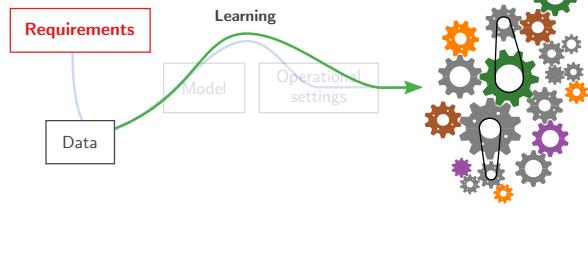
Why requirements?

This collage presents a range of news stories that highlight the challenges of AI. It includes articles from The New York Times, MIT Technology Review, The Appeal, and Slate. Topics covered include the self-driving Uber car accident, AI bias in recruiting, automation bias in AI, and the Tesla Autopilot system's fault. The stories emphasize the need for ethical and safe AI development.

Why requirements?

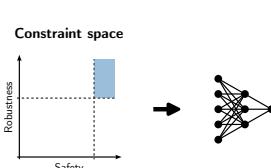


Why requirements?



What is a requirements?

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



[NASA, "Systems engineering handbook," 2019]

1

2

3

4

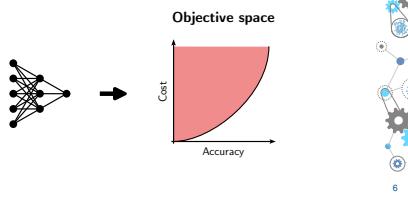
5

6

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]

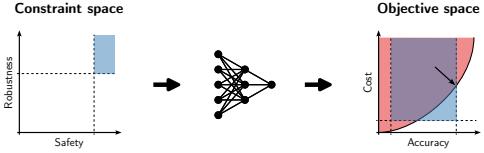


6

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]



6

What is (un)constrained learning?

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

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

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



7

What is (un)constrained learning?

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

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

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



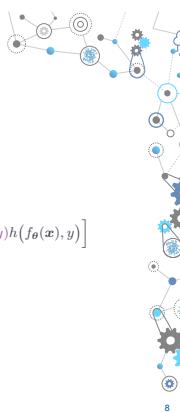
7

What about penalties?

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

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])
- ...

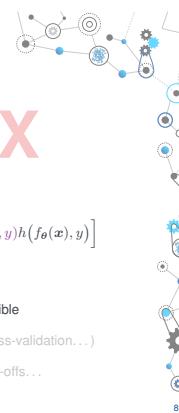


8

What about penalties?

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

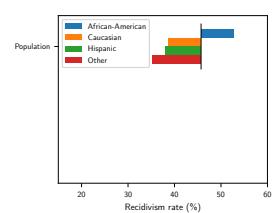
- ✖ There need not exist (λ, μ) for which the penalized solution is optimal *and* feasible
- ✖ Even if such (λ, μ) exist, they are not easy to find (hyperparameter search, cross-validation...)
- ✓ Constrained learning yields better guarantees, better performance, better trade-offs...



8

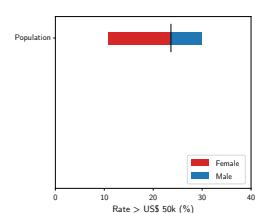
Fairness

Problem
Predict whether an individual will recidivate



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

Problem
Predict whether an individual makes > \$50k

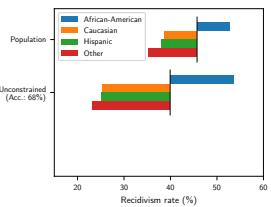


10

Fairness

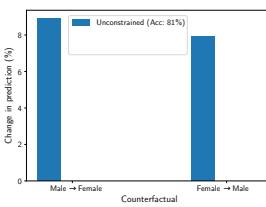
Problem

Predict whether an individual will recidivate



Problem

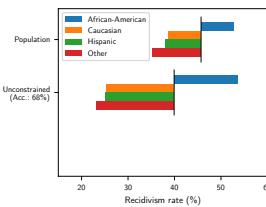
Predict whether an individual makes > \$50k



Fairness

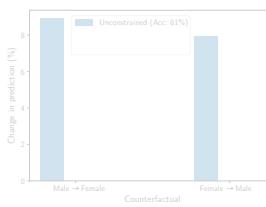
Problem

Predict whether an individual will recidivate



Problem

Predict whether an individual makes > \$50k



*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]

Fairness: “Equality” of odds

Problem

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

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

10

Problem

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

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

11

*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]

Fairness: “Equality” of odds

Problem

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

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

11

Fairness: “Equality” of odds

Problem

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

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

11

*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]

Fairness: “Equality” of odds

Problem

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

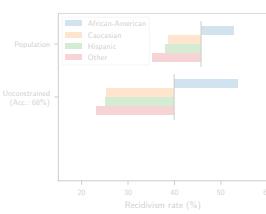
$$\begin{aligned} \min_{\theta} & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(\mathbf{x}_n), y_n) \\ \text{subject to} & \frac{1}{N} \sum_{n=1}^N \mathbb{E}[f_{\theta}(\mathbf{x}_n) = 1 \mid \text{Race}] \leq \frac{1}{N} \sum_{n=1}^N \mathbb{E}[f_{\theta}(\mathbf{x}_n) = 1] + c, \\ & \text{for Race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

11

Fairness

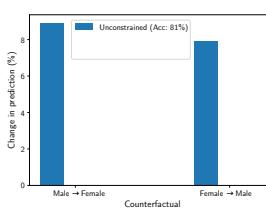
Problem

Predict whether an individual will recidivate



Problem

Predict whether an individual makes > \$50k



12

*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]

*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]

Counterfactual fairness

Problem

Predict whether an individual makes > \$50k while being invariant to gender

$$\begin{aligned} \min_{\theta} & \text{ Prediction error} \\ \text{subject to} & \text{ Change in prediction } (\rho x) \leq c \text{ a.e.} \\ & (\rho : \text{Male} \leftrightarrow \text{Female}) \end{aligned}$$

*We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Chamom and Ribeiro, NeurIPS'20]

Counterfactual fairness

Problem

Predict whether an individual makes > \$50k while being invariant to gender

$$\begin{aligned} \min_{\theta} & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} & \text{ Change in prediction } (\rho x) \leq c \text{ a.e.} \\ & (\rho : \text{Male} \leftrightarrow \text{Female}) \end{aligned}$$

*We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Chamom and Ribeiro, NeurIPS'20]

13

13

Counterfactual fairness

Problem

Predict whether an individual makes > \$50k while being invariant to gender

$$\begin{aligned} \min_{\theta} & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} & \frac{1}{N} \sum_{n=1}^N D_{\text{KL}}(f_{\theta}(x_n) \| f_{\theta}(\rho x_n)) \leq c, \text{ for all } n \\ & (\rho : \text{Male} \leftrightarrow \text{Female}) \end{aligned}$$

*We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
[Chamom and Ribeiro, NeurIPS'20]

Applications

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

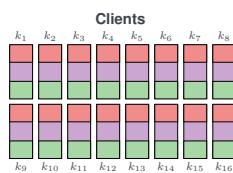
14

Federated learning

Problem

Learn a common model using data using data distributed among K clients

$$\min_{\theta} \text{ Average loss across clients}$$



13

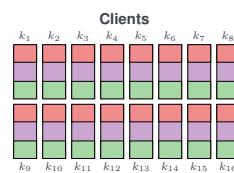
14

Federated learning

Problem

Learn a common model using data using data distributed among K clients

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



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

[Shen et al., ICRL'22]

15

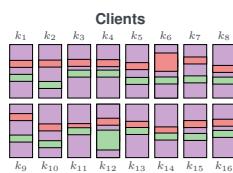
15

Heterogeneous federated learning

Problem

Learn a common model using data using data distributed among K clients

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



15

15

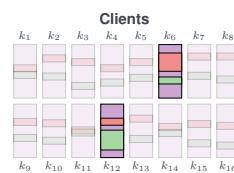
[Shen et al., ICRL'22]

Federated learning

Problem

Learn a common model using data using data distributed among K clients

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



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

[Shen et al., ICRL'22]

15

15

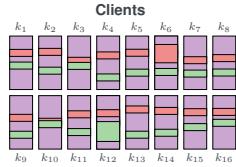
Federated learning

Problem

Learn a common model using data using data distributed among K clients

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

subject to Loss disparity (k -th client) $\leq c$,



15

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

[Shen et al., ICRL22]

Federated learning

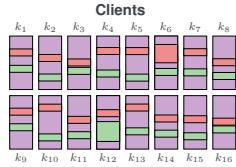
Problem

Learn a common model using data using data distributed among K clients

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

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

$k = 1, \dots, K$

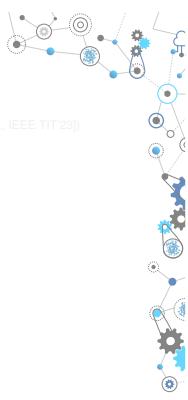


- k -th client loss: $\text{Loss}_k(\phi) = \frac{1}{N_k} \sum_{n_k=1}^{N_k} \text{Loss}(f_{\phi}(x_{n_k}), y_{n_k})$

[Shen et al., ICRL22]

Applications

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

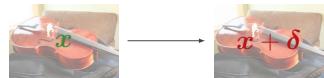


16

Robustness

Problem

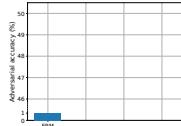
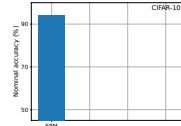
Learn a classifier that is robust to input perturbations



Cello



Hammer



17

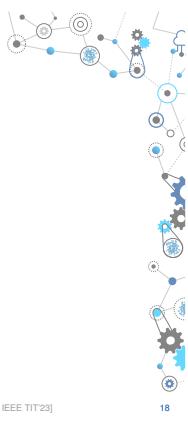
Robustness

Problem

Learn a classifier that is robust to input perturbations

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

subject to Adversarial loss $\leq c$



18

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

Robustness

Problem

Learn a 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 Adversarial loss $\leq c$



18

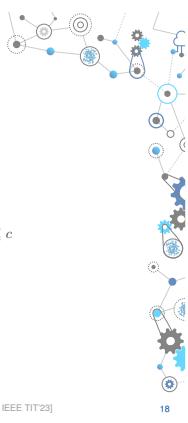
Robustness

Problem

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



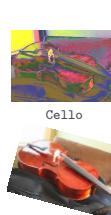
18

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

Invariance

Problem

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



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

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

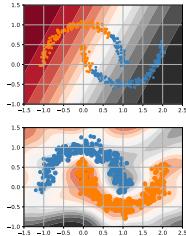
19

[Hounie et al., ICML'23]

Smoothness

Problem

Learn a classifier **that is smooth**, i.e., Lipschitz continuous (on a manifold)



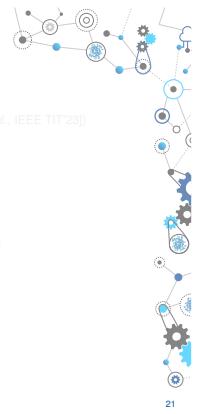
[Cerviño et al., ICML'23]

$$\begin{aligned} \min_{\theta} & \quad \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} & \quad \max_{x \in \mathcal{M}} \|\nabla_x f_{\theta}(x)\|^2 \leq L \end{aligned}$$

20

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])
- ...

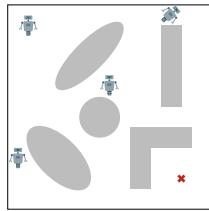


21

Safety

Problem

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

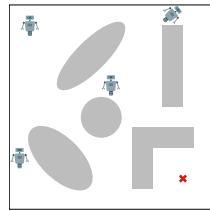


22

Safety

Problem

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



[Paternain et al., IEEE TAC'23]

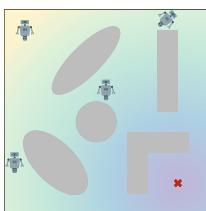
$$\begin{aligned} \text{maximize}_{\pi \in \mathcal{P}(\mathcal{S})} & \quad \text{Task reward} \\ \text{subject to} & \quad \Pr[\text{Colliding with } \mathcal{O}_i] \leq \delta, \\ & \quad \text{for } i = 1, 2, \dots \end{aligned}$$

23

Safety

Problem

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



[Paternain et al., IEEE TAC'23]

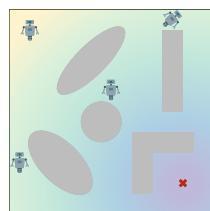
$$\begin{aligned} \text{maximize}_{\pi \in \mathcal{P}(\mathcal{S})} & \quad \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \\ \text{subject to} & \quad \Pr[\text{Colliding with } \mathcal{O}_i] \leq \delta, \\ & \quad \text{for } i = 1, 2, \dots \end{aligned}$$

23

Safety

Problem

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



[Paternain et al., IEEE TAC'23]

$$\begin{aligned} \text{maximize}_{\pi \in \mathcal{P}(\mathcal{S})} & \quad \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \\ \text{subject to} & \quad \Pr \left(\bigcap_{t=0}^{T-1} \{s_t \notin \mathcal{O}_i\} \mid \pi \right) \geq 1 - \delta_i, \\ & \quad \text{for } i = 1, 2, \dots \end{aligned}$$

23

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])
- Wireless resource allocation (e.g., [Eisen et al., IEEE TSP'19])
- Continual learning (e.g., [Peng et al., ICML'23])
- Active learning (e.g., [Elenter et al., NeurIPS'22])
- Semi-supervised learning (e.g., [Cerviño et al., ICML'23])
- Minimum norm interpolation, SVM...



24

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]
- $(x_n, y_n) \sim \mathcal{D}, (x_m, y_m) \sim \mathfrak{A}, (x_r, y_r) \sim \mathfrak{P}$ (i.i.d.)

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

26

What is (un)constrained learning?

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

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

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

27

Constrained learning challenges

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

Challenges

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

28

Constrained learning challenges

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

Challenges

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

28

Constrained learning challenges

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

Challenges

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

28

Agenda

Constrained learning theory

Constrained learning algorithms

Resilient constrained learning

29



Constrained learning challenges

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

Challenges

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

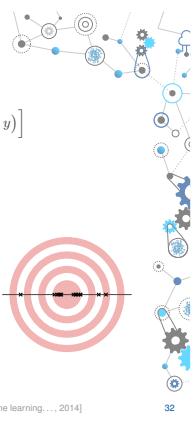
30

What classical learning theory says?

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

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

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



32

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

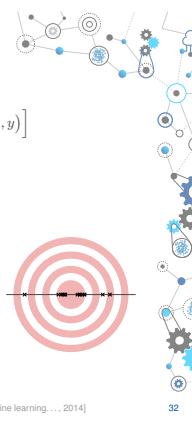
What classical learning theory says?

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

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

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

- ✗ Requirements?



32

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

- near-optimal

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



33

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

What's in a solution?

Definition (PACC learnability)

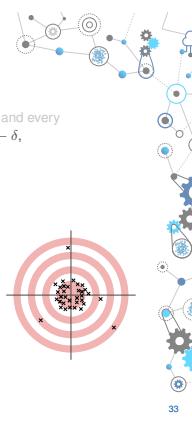
f_{θ} is a probably approximately correct constrained (PACC) learnable if for every ϵ, δ and every distributions $\mathcal{D}, \mathfrak{A}$, we can obtain f_{θ^*} from $N_f(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,

- near-optimal

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

- approximately feasible

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

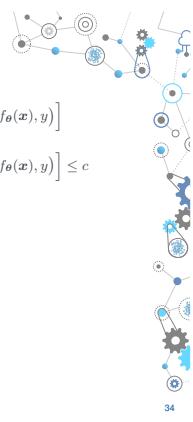


33

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

When is constrained learning possible?

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



34

Proposition

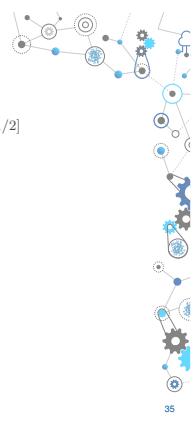
f_{θ} is PAC learnable $\Rightarrow f_{\theta}$ is PACC learnable

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

ECRM is not a PACC learner

Counter-example

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



35

ECRM is not a PACC learner

Counter-example

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

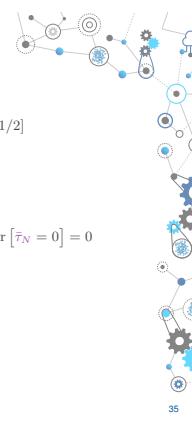
ECRM is not a PACC learner

Counter-example

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

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

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



35

ECRM is not a PACC learner

Counter-example

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

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

$$\begin{aligned} \hat{P}_r^* &= \min_{\theta \in \Theta} J(\theta) \\ \text{subject to} \quad &\theta_2 \bar{\tau}_N \leq \theta_1 - 1 + r_1 \\ &-\theta_1 \bar{\tau}_N \leq 1 - \theta_2 + r_2 \end{aligned}$$

$$\Pr[|\hat{P}_r^* - P^*| \leq 1/32] \leq 4e^{-0.001N},$$

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

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



Constrained learning challenges

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

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



Challenges

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

2) *Computational*: can we solve the constrained empirical problem?



Constrained learning challenges

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

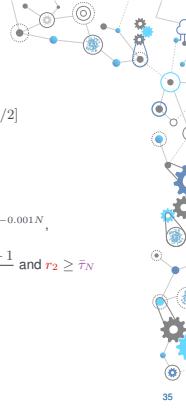
$$\begin{aligned} P^* &= \min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_\theta(\mathbf{x}), y)] \\ \text{subject to} \quad &\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{X}} [g(f_\theta(\mathbf{x}), y)] \leq c \\ &h(f_\theta(\mathbf{x}), y) \leq u \text{ a.e.} \end{aligned}$$



Duality

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

DUAL

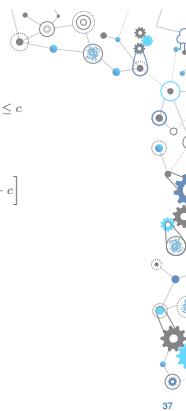


PRIMAL
↔
DUAL



Duality

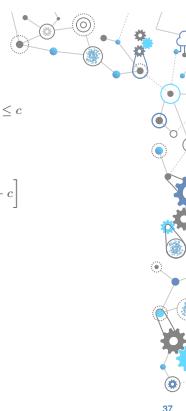
$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta(\mathbf{x}_n), y_n) \text{ subject to } \frac{1}{N} \sum_{m=1}^N g(f_\theta(\mathbf{x}_m), y_m) \leq c \\ &\updownarrow \\ \hat{D}^* &= \max_{\lambda \geq 0} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta(\mathbf{x}_n),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



Duality

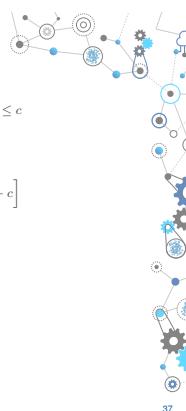
$$\begin{aligned} \hat{P}^* &= \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta(\mathbf{x}_n), y_n) \text{ subject to } \frac{1}{N} \sum_{m=1}^N g(f_\theta(\mathbf{x}_m), y_m) \leq c \\ &\updownarrow \end{aligned}$$

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



Duality

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



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

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

35

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

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

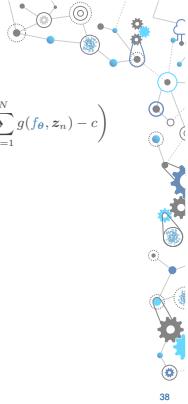
36

37

An alternative path

$$\begin{aligned} \hat{P}^* &= \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) \\ \text{s. to } &\frac{1}{N} \sum_{n=1}^N g(f_\theta, z_n) \leq c \\ P^* &= \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_\theta, z)] \\ \text{s. to } &\mathbb{E}_z [g(f_\theta, z)] \leq c \end{aligned}$$

↑ PAC

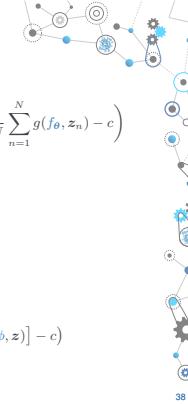


[Chamon and Ribeiro, NeurIPS'20; Chamon et al., IEEE TIT'23]

An alternative path

$$\begin{aligned} \hat{P}^* &= \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) \\ \text{s. to } &\frac{1}{N} \sum_{n=1}^N g(f_\theta, z_n) \leq c \\ P^* &= \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_\theta, z)] \\ \text{s. to } &\mathbb{E}_z [g(f_\theta, z)] \leq c \\ \downarrow &\mathcal{H}_\theta \subset \mathcal{H} \\ \hat{P}^* &= \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] \\ \text{s. to } &\mathbb{E}_z [g(\phi, z)] \leq c \end{aligned}$$

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

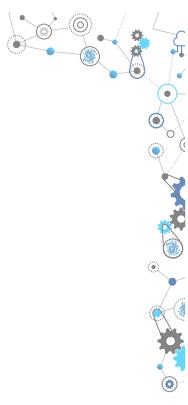


[Chamon and Ribeiro, NeurIPS'20; Chamon et al., IEEE TIT'23]

Non-convex variational duality

Convex optimization: Primal \longleftrightarrow Dual

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



39

Non-convex variational duality

Convex optimization: Primal \longleftrightarrow Dual

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

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



[Chamon et al., IEEE TSP'20]

Sparse logistic regression

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

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



40

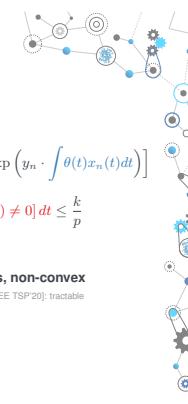
Sparse logistic regression

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

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

$$\begin{aligned} \min_{\theta \in L_2} & - \sum_{n=1}^N \log \left[1 + \exp \left(y_n \cdot \int \theta(t) x_n(t) dt \right) \right] \\ \text{s. to } & \|\theta\|_{L_0} = \int \mathbb{I}[\theta(t) \neq 0] dt \leq \frac{k}{p} \end{aligned}$$

Continuous, non-convex
[Chamon et al., IEEE TSP'20]: tractable



40

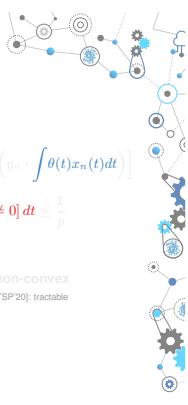
Sparse logistic regression

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

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

$$\begin{aligned} \min_{\theta \in L_2} & - \sum_{n=1}^N \log \left[1 + \exp \left(y_n \cdot \int \theta(t) x_n(t) dt \right) \right] \\ \text{s. to } & \|\theta\|_{L_0} = \int \mathbb{I}[\theta(t) \neq 0] dt \leq \frac{k}{p} \end{aligned}$$

Continuous, non-convex
[Chamon et al., IEEE TSP'20]: tractable

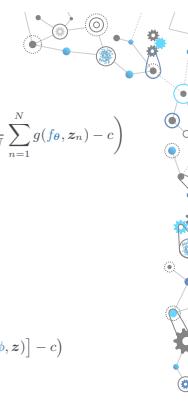


40

An alternative path

$$\begin{aligned} \hat{P}^* &= \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) \\ \text{s. to } &\frac{1}{N} \sum_{n=1}^N g(f_\theta, z_n) \leq c \\ P^* &= \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_\theta, z)] \\ \text{s. to } &\mathbb{E}_z [g(f_\theta, z)] \leq c \\ \hat{P}^* &= \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] \\ \text{s. to } &\mathbb{E}_z [g(\phi, z)] \leq c \end{aligned}$$

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



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

41

An alternative path

$$\begin{aligned} P^* &= \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) \\ \text{s. to } &\frac{1}{N} \sum_{n=1}^N g(f_\theta, z_n) \leq c \\ &\uparrow \text{PAC} \\ P^* &= \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_\theta, z)] \quad \xleftarrow{\epsilon_0} D^* = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) + \lambda \left(\frac{1}{N} \sum_{n=1}^N g(f_\theta, z_n) - c \right) \\ \text{s. to } &\mathbb{E}_z [g(f_\theta, z)] \leq c \\ &\uparrow \epsilon_0 \\ \tilde{P}^* &= \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] \quad \xleftarrow{\epsilon_0} \tilde{D}^* = \max_{\lambda \geq 0} \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] + \lambda (\mathbb{E}_z [g(\phi, z)] - c) \\ \text{s. to } &\mathbb{E}_z [g(\phi, z)] \leq c \end{aligned}$$

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

41

An alternative path

$$\begin{aligned} P^* &= \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) \\ \text{s. to } &\frac{1}{N} \sum_{n=1}^N g(f_\theta, z_n) \leq c \\ &\uparrow \text{PAC} \\ P^* &= \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_\theta, z)] \quad \xleftarrow{\epsilon_0} D^* = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta, z_n) + \lambda (\mathbb{E}_z [g(f_\theta, z)] - c) \\ \text{s. to } &\mathbb{E}_z [g(f_\theta, z)] \leq c \\ &\uparrow \epsilon_0 \\ \tilde{P}^* &= \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] \quad \xleftarrow{\epsilon_0} \tilde{D}^* = \max_{\lambda \geq 0} \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] + \lambda (\mathbb{E}_z [g(\phi, z)] - c) \\ \text{s. to } &\mathbb{E}_z [g(\phi, z)] \leq c \end{aligned}$$

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

41

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

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

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

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

42

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

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

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

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

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

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

(mild conditions apply)

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

42

Dual (near-)PACC learning

Theorem

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

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

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

$$\epsilon_0 = M\nu \quad \epsilon = B \sqrt{\frac{1}{N} \left[1 + \log \left(\frac{4m(2N)^{d_{VC}}}{\delta} \right) \right]}$$

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

Sources of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

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

43

Dual (near-)PACC learning

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

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

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

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

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

(mild conditions apply)

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

42

Dual (near-)PACC learning

Theorem

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

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

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

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

Sources of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

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

43

Dual (near-)PACC learning

Dual (near-)PACC learning

Theorem

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

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

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

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

Sources of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

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

43

Dual (near-)PACC learning

Theorem

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

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

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

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

Sources of error

parametrization richness (ν)

sample size (N)

requirements difficulty (λ^*)

43

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

Dual (near-)PACC learning

Theorem

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

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

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

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

Sources of error

parametrization richness (ν)

sample size (N)

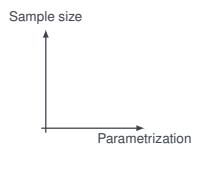
requirements difficulty (λ^*)

43

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

Dual learning trade-offs

- Unconstrained learning
parametrization \times sample size

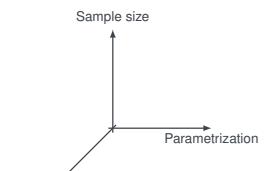


44

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

Dual learning trade-offs

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



44

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

When is constrained learning possible?

Corollary

f_θ is PAC learnable $\approx^* f_\theta$ is PACC learnable

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

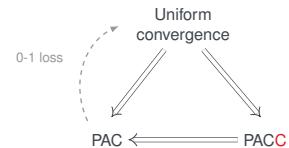


45

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

When is constrained learning possible?

Corollary



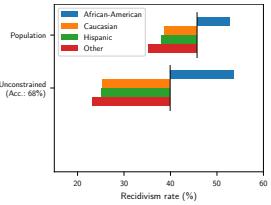
45

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

Fairness

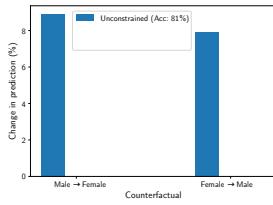
Problem

Predict whether an individual will recidivate



Problem

Predict whether an individual makes > \$50k



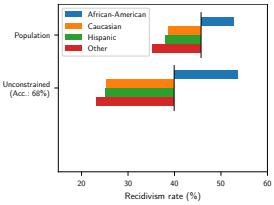
46

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

Fairness

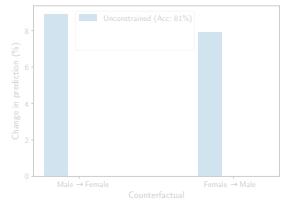
Problem

Predict whether an individual will recidivate



Problem

Predict whether an individual makes > \$50k



46

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

Fairness: “Equality” of odds

Problem

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

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

*We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Cotter et al., JMLR19; Chamon et al., IEEE TIT’23]

47

Fairness: “Equality” of odds

Problem

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

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

*We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Cotter et al., JMLR19; Chamon et al., IEEE TIT’23]

47

Fairness: “Equality” of odds

Problem

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

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

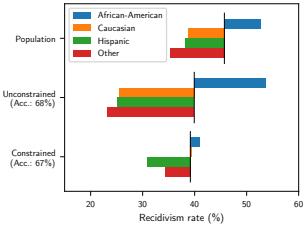
*We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Cotter et al., JMLR19; Chamon et al., IEEE TIT’23]

47

Fairness: “Equality” of odds

Problem

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



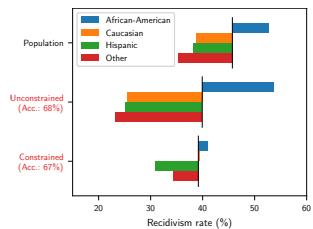
*We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Chamon et al., IEEE TIT’23]

48

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 et al., IEEE TIT’23]

48

Fairness: “Equality” of odds

Prediction

		True label	
		0	1
African-American	0	31% (16%)	16% (36%)
	1	16% (37%)	37% (23%)
Caucasian	0	52% (9%)	9% (44%)
	1	23% (16%)	16% (23%)
		Unconstrained	Constrained

*We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Chamon et al., IEEE TIT’23]

49

Fairness: “Equality” of odds

Fairness: “Equality” of odds

Prediction

		True label	
		0	1
African-American	0	31% (16%)	16% (36%)
	1	16% (37%)	37% (23%)
Caucasian	0	52% (9%)	9% (44%)
	1	23% (16%)	16% (23%)
		Unconstrained	Constrained

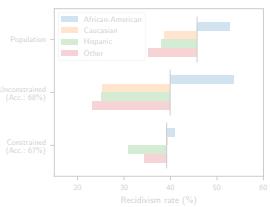
*We say “Race” to follow the terminology used during the data collection of the COMPAS dataset.
[Chamon et al., IEEE TIT’23]

49

Fairness

Problem

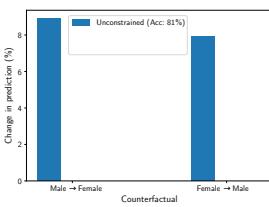
Predict whether an individual will recidivate.



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

Problem

Predict whether an individual makes > \$50k



50

Counterfactual fairness

Problem

Predict whether an individual makes > \$50k while being invariant to gender

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

$$\text{subject to } \text{DKL}(f_{\theta}(x_n) \| f_{\theta}(\rho x_n)) \leq c, \text{ for all } n$$

$$(\rho : \text{Male} \leftrightarrow \text{Female})$$

[Chamon and Ribeiro, NeurIPS'20]

51

Counterfactual fairness

Problem

Predict whether an individual makes > \$50k while being invariant to gender

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} \quad & \frac{1}{N} \sum_{n=1}^N \text{DKL}(f_{\theta}(x_n) \| f_{\theta}(\rho x_n)) \leq c, \text{ for all } n \\ & (\rho : \text{Male} \leftrightarrow \text{Female}) \end{aligned}$$

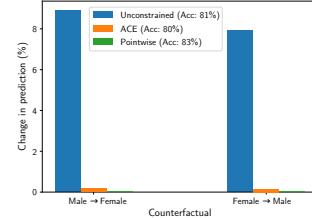
51

[Chamon and Ribeiro, NeurIPS'20]

Counterfactual fairness

Problem

Predict whether an individual makes > \$50k while being invariant to gender



[Chamon and Ribeiro, NeurIPS'20]

52

Counterfactual fairness

Problem

Predict whether an individual makes > \$50k while being invariant to gender

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} \quad & \text{DKL}(f_{\theta}(x_n) \| f_{\theta}(\rho x_n)) \leq c, \text{ for all } n \\ & (\rho : \text{Male} \leftrightarrow \text{Female}) \\ \max_{\lambda \geq 0} \quad & \min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) + \sum_{n=1}^N \lambda_n [\text{DKL}(f_{\theta}(x_n) \| f_{\theta}(\rho x_n)) - c] \end{aligned}$$

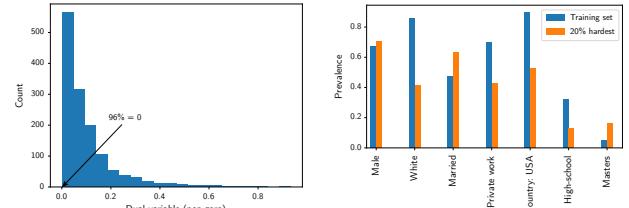
53

[Chamon and Ribeiro, NeurIPS'20]

Counterfactual fairness

Problem

Predict whether an individual makes > \$50k while being invariant to gender



[Chamon and Ribeiro, NeurIPS'20]

54

Agenda

Constrained learning theory

Constrained learning algorithms

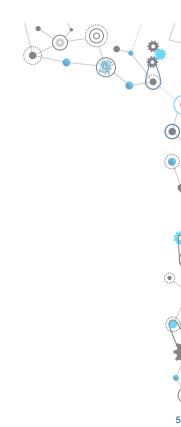
Resilient constrained learning



Constrained optimization methods

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

55



56

Constrained optimization methods

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

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

56

Constrained optimization methods

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

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

56

Constrained optimization methods

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

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

56

Dual learning algorithm

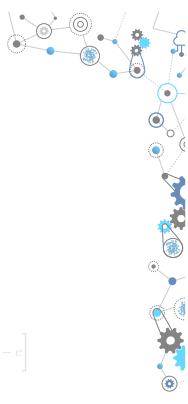
$$\hat{D}^* = \max_{\lambda \geq 0} \min_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{n=1}^N \ell(f_\theta(\mathbf{x}_n), y_n) + \lambda \left[\frac{1}{N} \sum_{m=1}^N g(f_\theta(\mathbf{x}_m), y_m) - c \right]$$

57

Dual learning algorithm

- Minimize the primal (\equiv ERM)

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



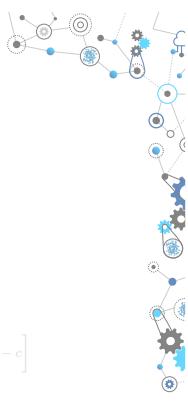
57

Dual learning algorithm

- Minimize the primal (\equiv ERM)

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

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



57

Dual learning algorithm

- Minimize the primal (\equiv ERM)

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

- Update the dual

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

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

A (near-)PACC learner

Theorem

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

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

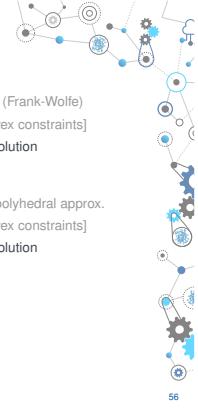
Then, after $T = \left\lceil \frac{\|\lambda^*\|^2}{2\eta M^2} \right\rceil + 1$ dual iterations with step size $\eta \leq \frac{2\epsilon}{mB^2}$,

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

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

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

[Chamon et al., IEEE TIT'23]



58

In practice...

- Minimize the primal (\equiv ERM)

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

- Update the dual

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

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

59



In practice...

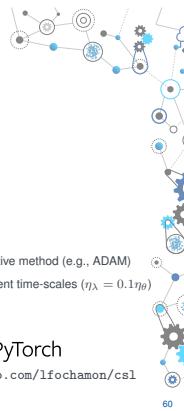
```

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:      $\beta_{n+1} \leftarrow \beta_n - \eta_{\theta} \nabla_{\beta} [\ell(f_{\beta_n}(x_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(x_n), y_n)]$ 
6:   end
7:    $\theta_t \leftarrow \beta_{N+1}$ 
8:    $\lambda_t = \left[ \lambda_{t-1} + \eta_{\lambda} \left( \frac{1}{N} \sum_{m=1}^N g(f_{\theta_t}(x_m), y_m) - c \right) \right]_+$ 
9: end
10: Output:  $\theta_T, \lambda_T$ 

```

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

60



In practice...

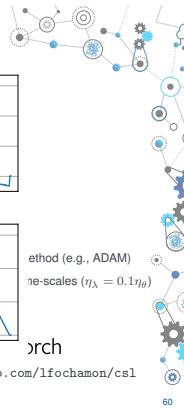
```

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:      $\beta_{n+1} \leftarrow \beta_n - \eta_{\theta} \nabla_{\beta} [\ell(f_{\beta_n}(x_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(x_n), y_n)]$ 
6:   end
7:    $\theta_t \leftarrow \beta_{N+1}$ 
8:    $\lambda_t = \left[ \lambda_{t-1} + \eta_{\lambda} \left( \frac{1}{N} \sum_{m=1}^N g(f_{\theta_t}(x_m), y_m) - c \right) \right]_+$ 
9: end
10: Output:  $\theta_T, \lambda_T$ 

```

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

60



In practice...

```

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

```

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

60

In practice...

- Minimize the primal (\equiv ERM)

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

- Update the dual

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

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

59



In practice...

```

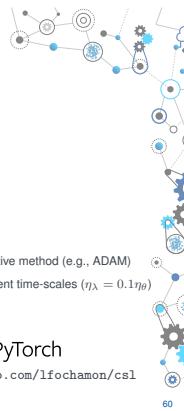
1: Initialize:  $\theta_0, \lambda_0$ 
2: for  $t = 1, \dots, T$ 
3:    $\beta_1 \leftarrow \theta_{t-1}$ 
4:   for  $n = 1, \dots, N$ 
5:      $\beta_{n+1} \leftarrow \beta_n - \eta_{\theta} \nabla_{\beta} [\ell(f_{\beta_n}(x_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(x_n), y_n)]$ 
6:   end
7:    $\theta_t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

```

PyTorch

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

60



In practice...

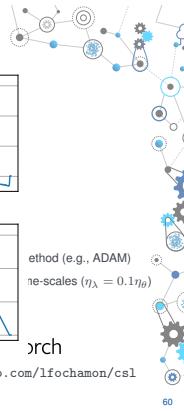
```

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:      $\beta_{n+1} \leftarrow \beta_n - \eta_{\theta} \nabla_{\beta} [\ell(f_{\beta_n}(x_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(x_n), y_n)]$ 
6:   end
7:    $\theta_t \leftarrow \beta_{N+1}$ 
8:    $\lambda_t = \left[ \lambda_{t-1} + \eta_{\lambda} \left( \frac{1}{N} \sum_{m=1}^N g(f_{\theta_t}(x_m), y_m) - c \right) \right]_+$ 
9: end
10: Output:  $\theta_T, \lambda_T$ 

```

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

60



In practice...

Penalty-based vs. dual learning

Penalty-based learning

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

Dual learning

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

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

- Parameter: λ (data-dependent)

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

- Parameter: c (requirement-dependent)

- Generalizes with respect to Loss and $\text{Penalty} \leq c$

61

Agenda

Constrained learning theory



Constrained learning algorithms

Resilient constrained learning

Heterogeneous federated learning

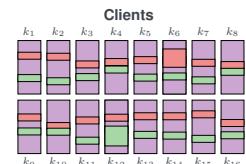
Problem

Learn a common model using data distributed among K clients

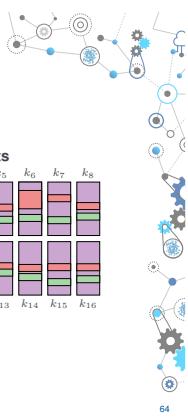
$$\min_{\theta} \frac{1}{K} \sum_{k=1}^K \text{Loss}_k(f_\theta)$$

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

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



62



Heterogeneous federated learning

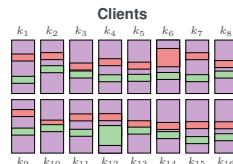
Problem

Learn a common model using data distributed among K clients

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

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

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



63

Heterogeneous federated learning

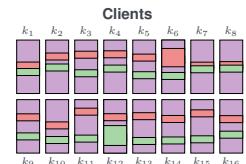
Problem

Learn a common model using data distributed among K clients

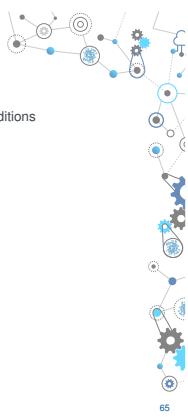
$$\min_{\theta} \frac{1}{K} \sum_{k=1}^K \text{Loss}_k(f_\theta)$$

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

- 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

Resilient constrained learning

Definition (Resilience)

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

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

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

65

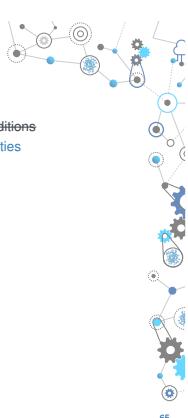
Resilient constrained learning

Definition (Resilience)

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

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

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



65



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

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

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

- 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

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

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

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

66

Resilient constrained learning

Definition (Resilience)

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

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

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

- 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!

66

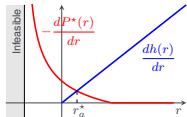
Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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



67

[Hounie et al., NeurIPS'23]

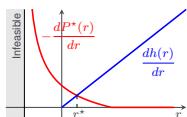
Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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



67

[Hounie et al., NeurIPS'23]

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

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

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

- After relaxing, $\lambda^*(\mathbf{r}^*)$ is smaller than $\lambda^*(0)$
 \Rightarrow Resilient constrained learning "generalizes better" (lower sample complexity)

68

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

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

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

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

[Hounie et al., NeurIPS'23]



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

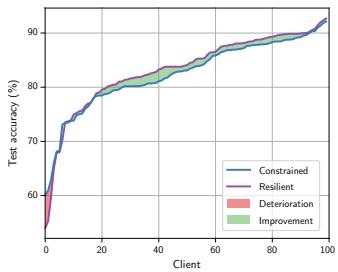
$$\begin{aligned} P^*(\mathbf{r}^*) &= \min_{\theta, \mathbf{r}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\text{Loss}(f_\theta(\mathbf{x}), y)] + h(\mathbf{r}) \\ \text{subject to } \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{X}_i} [g_i(f_\theta(\mathbf{x}_m), y_m)] &\leq c_i + r_i \end{aligned}$$

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

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

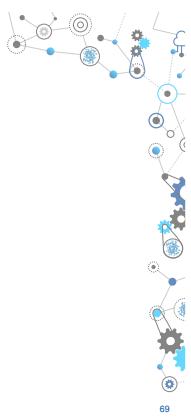
[Hounie et al., NeurIPS'23]

Heterogeneous federated learning

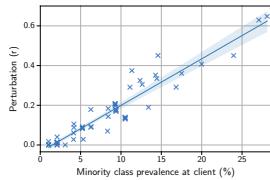


[Hounie et al., NeurIPS'23]

68

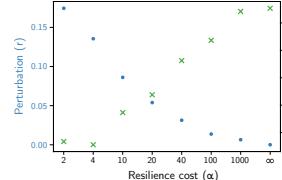


Heterogeneous federated learning



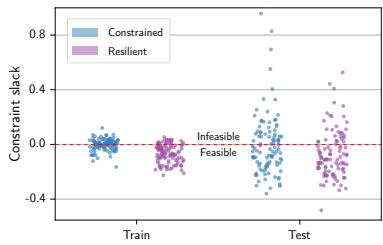
[Hounie et al., NeurIPS'23]

69



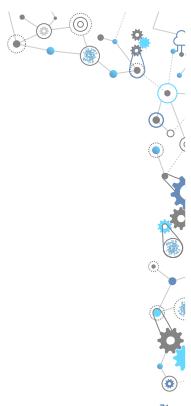
70

Heterogeneous federated learning



[Hounie et al., NeurIPS'23]

71



Summary

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

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?

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?

72



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., ICLR'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 [Houle et al., NeurIPS'23]

72

Robustness constraints



Agenda

Adversarially robust learning

Semi-infinite learning

Probabilistic robustness

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

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

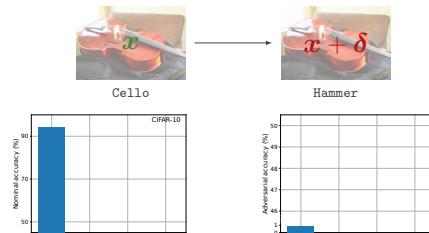
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \longrightarrow \min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

74

Robust learning

Problem

Learn an image classifier that is robust to input perturbations



75

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

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

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

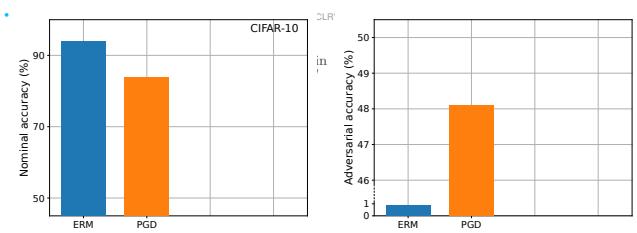
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) + \lambda \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

76

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations

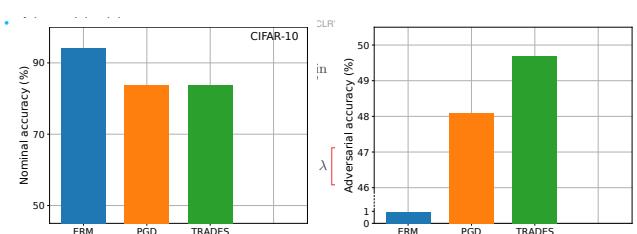


76

Adversarial training

Problem

Learn an image classifier that is robust to input perturbations



77

[Zhang et al., ICML'19]

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

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

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

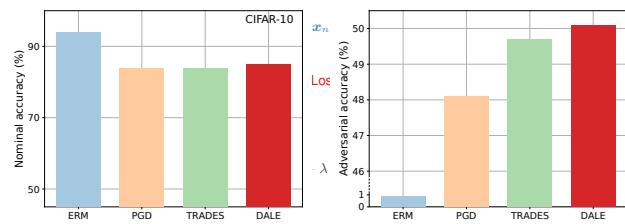


78

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations



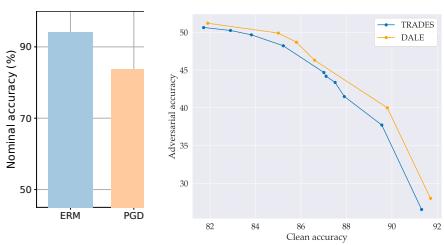
[Chamon and Ribeiro, NeurIPS'20; Robey et al., NeurIPS'21; Chamon et al., IEEE TIT'23]

78

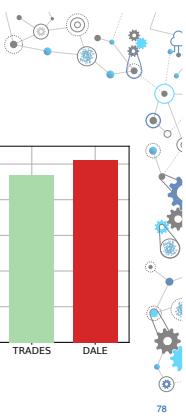
Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations



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



78

Penalty-based vs. dual learning

Penalty-based learning

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

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

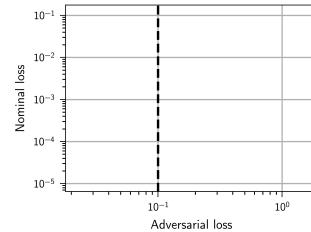
Dual learning

$$\begin{aligned} \theta^* &\in \operatorname{argmin}_{\theta} \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta) \\ \lambda^+ &= [\lambda + \eta(\text{Penalty}(\theta^*) - c)]_+ \end{aligned}$$

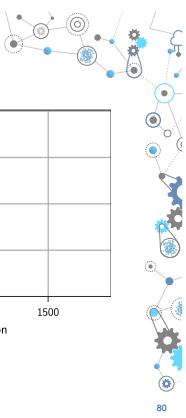
- Parameter: c (requirement-dependent)
- Generalizes with respect to Loss and $\text{Penalty} \leq c$

79

Constrained learning for robustness

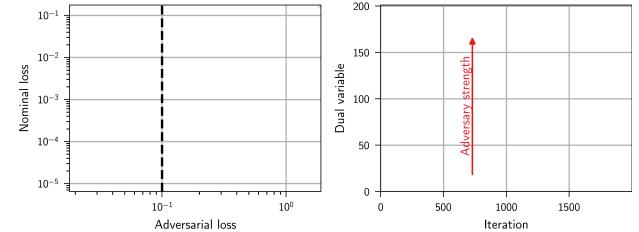


[Chamon et al., IEEE TIT'23]

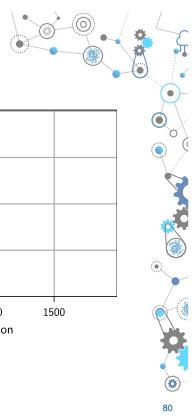


80

Constrained learning for robustness

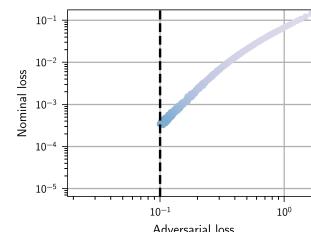


[Chamon et al., IEEE TIT'23]

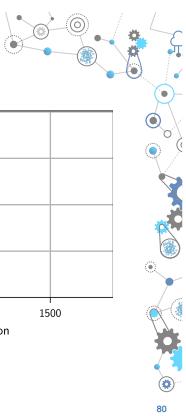


80

Constrained learning for robustness

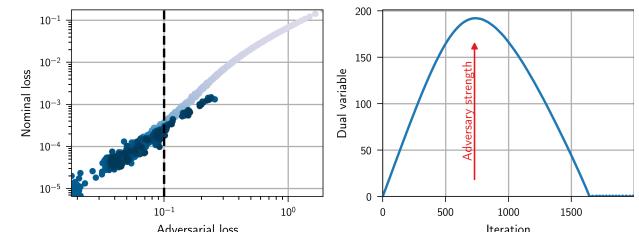


[Chamon et al., IEEE TIT'23]

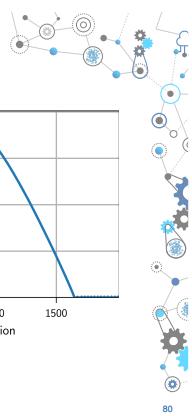


80

Constrained learning for robustness

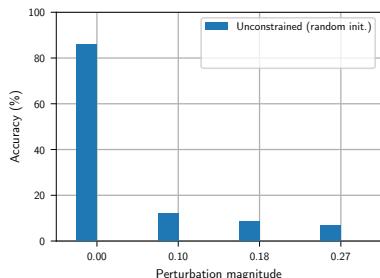


Empirical observations: [Zhang et al., ICML'20; Sitawarin, arXiv'20]



80

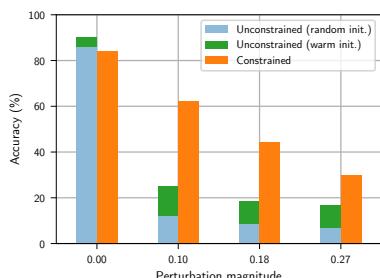
Constrained learning for robustness



[Chamon et al., IEEE TIT'23]

81

Constrained learning for robustness



[Chamon et al., IEEE TIT'23]

81

Constrained learning for robustness

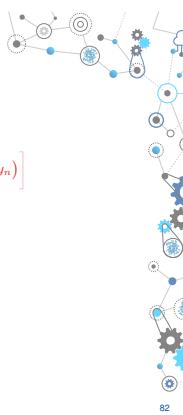
Problem

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



82

Adversarial training

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

✗ "PGD" [Madry et al., ICLR'18]

- 1: $\delta^1 \leftarrow \delta_{t-1}$
- 2: **for** $k = 1, \dots, K$
- 3: $\delta^{k+1} \leftarrow \text{proj}_{\Delta} \left[\delta^k + \eta \text{sign} \left(\nabla_{\delta} \text{Loss}(f_{\theta^k}(x + \delta^k), y) \right) \right]$
- 4: **end**
- 5: $\delta_t \leftarrow \delta^{K+1}$
- 6: $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \text{Loss}(f_{\theta}(x + \delta_t), y)$

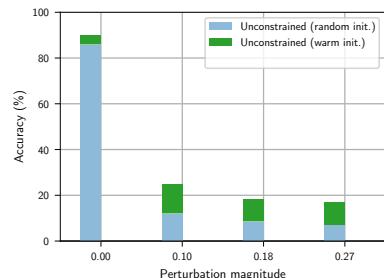
- Random initialization
- Restarts
- Pruning
- Adaptive step size



83

[Dhillon et al., ICLR'18; Carmon et al., NeurIPS'19; Wu et al., NeurIPS'20; Cheng et al., IJCAI'22]

Constrained learning for robustness



[Chamon et al., IEEE TIT'23]

81

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



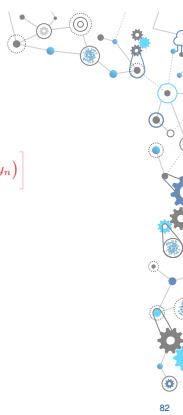
81

Adversarial training

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

• "PGD" [Madry et al., ICLR'18]

- 1: $\delta^1 \leftarrow \delta_{t-1}$
- 2: **for** $k = 1, \dots, K$
- 3: $\delta^{k+1} \leftarrow \text{proj}_{\Delta} \left[\delta^k + \eta \text{sign} \left(\nabla_{\delta} \text{Loss}(f_{\theta^k}(x + \delta^k), y) \right) \right]$
- 4: **end**
- 5: $\delta_t \leftarrow \delta^{K+1}$
- 6: $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \text{Loss}(f_{\theta}(x + \delta_t), y)$



83

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



84

Agenda

Adversarially robust learning



85

Semi-infinite learning

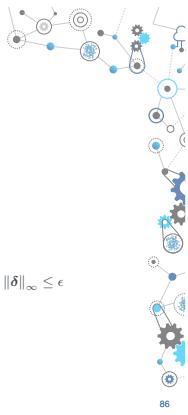
Probabilistic robustness

Semi-infinite constrained learning

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



86



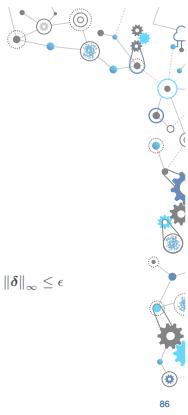
86

Semi-infinite constrained learning

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

- Epigraph formulation:

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



86

Semi-infinite constrained learning

$$\begin{aligned} & \min_{\theta} \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \\ \text{subject to } & \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_e), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_s), y_n) \leq t(x_n, y_n) \end{aligned}$$

- Epigraph formulation:

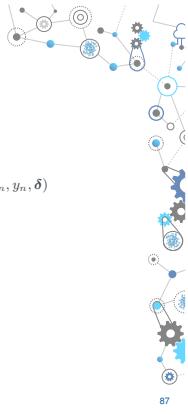
$$\max_{\|\delta\|_\infty \leq \epsilon} \text{Loss}(f_{\theta}(x + \delta), y) \leq t \iff \text{Loss}(f_{\theta}(x + \delta), y) \leq t, \text{ for all } \|\delta\|_\infty \leq \epsilon$$

- Semi-infinite program

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



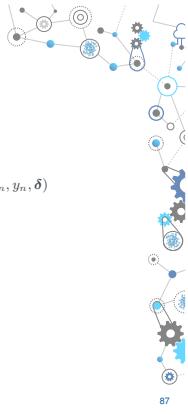
86



87

Duality

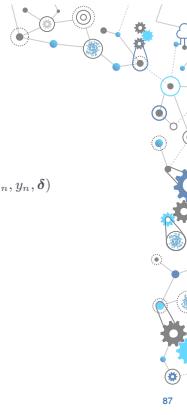
$$\begin{aligned} & \min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right] \\ & \uparrow = \\ & \min_{\theta} \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \text{ s. to } \text{Loss}(f_{\theta}(x_n + \delta), y_n) \leq t(x_n, y_n), \forall (x_n,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



87

Duality

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



87



88

From optimization to sampling

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



88

Proposition

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

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

From optimization to sampling

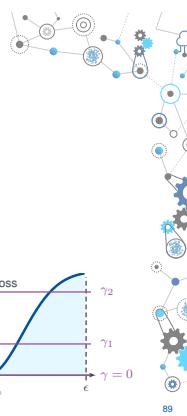
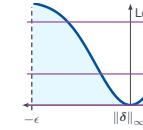
$$\begin{aligned} & \min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right] \\ & \uparrow \approx \\ & \min_{\theta} \sup_{\mu \in \mathcal{P}^2} \underbrace{\frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\delta \sim \mu(\cdot | x_n, y_n)} [\text{Loss}(f_{\theta}(x_n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

Proposition

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

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

[Robey et al., NeurIPS'21]



89

From optimization to sampling

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

$\uparrow \approx$

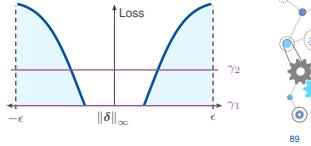
$$\min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\delta \sim \mu_{\gamma}} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]$$

$\underbrace{\quad}_{L(\theta, \mu)}$

Proposition

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

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



89

[Robey et al., NeurIPS'21]

From optimization to sampling

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

$\uparrow \approx$

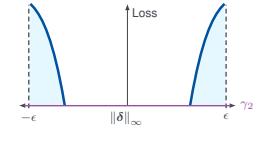
$$\min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\delta \sim \mu_{\gamma}} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]$$

$\underbrace{\quad}_{L(\theta, \mu)}$

Proposition

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

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



89

[Robey et al., NeurIPS'21]

From optimization to sampling

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

$\uparrow =$

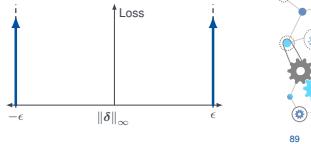
$$\min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\delta \sim \mu_{\gamma}} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]$$

$\underbrace{\quad}_{L(\theta, \mu)}$

Proposition

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

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



89

[Robey et al., NeurIPS'21]

From optimization to sampling

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

$\uparrow \approx$

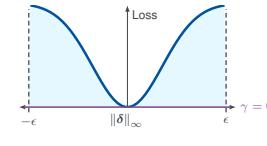
$$\min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\delta \sim \mu_{\gamma}} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]$$

$\underbrace{\quad}_{L(\theta, \mu)}$

Proposition

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

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



89

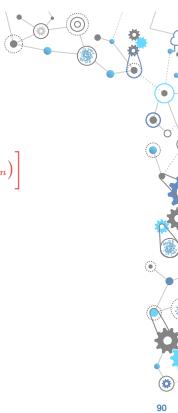
[Robey et al., NeurIPS'21]

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

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



90

• Balancing nominal accuracy and robustness \Rightarrow Dual constrained learning

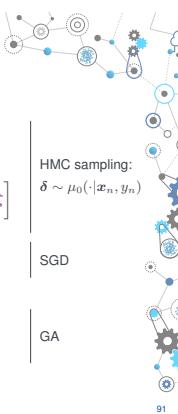
- ✖ Computing the worst-case perturbations
 - gradient ascent \rightarrow non-convex, underparametrized

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

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



91

• Balancing nominal accuracy and robustness \Rightarrow Dual constrained learning

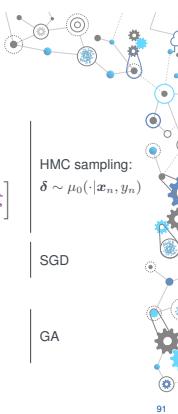
- ✖ Computing the worst-case perturbations
 - gradient ascent \rightarrow non-convex, underparametrized \Rightarrow sampling

Dual Adversarial LEarning

```

1: for n = 1, ..., N:
2:   δ_n ~ Random(Δ)
3:   for k = 1, ..., K:
4:     ζ ~ Laplace(0, I)
5:     δ_n ← proj_Δ [δ_n + η sign [∇_δ log (Loss(f_{θ_t}(x_n + δ_n), y_n))] + √{2ηT}ζ]
6:   end
7:   θ ← θ - η ∇_θ [Loss(f_{θ}(x_n), y_n) + λ Loss(f_{θ}(x_n + δ_n), y_n)]
8: end
9: λ ← [λ + η (1/N ∑_{n=1}^N Loss(f_{θ}(x_n + δ_n), y_n) - c)]_+

```



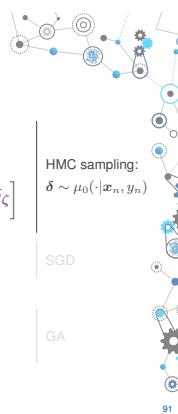
91

Dual Adversarial LEarning

```

1: for n = 1, ..., N:
2:   δ_n ~ Random(Δ)
3:   for k = 1, ..., K:
4:     ζ ~ Laplace(0, I)
5:     δ_n ← proj_Δ [δ_n + η sign [∇_δ log (Loss(f_{θ_t}(x_n + δ_n), y_n))] + √{2ηT}ζ]
6:   end
7:   θ ← θ - η ∇_θ [Loss(f_{θ}(x_n), y_n) + λ Loss(f_{θ}(x_n + δ_n), y_n)]
8: end
9: λ ← [λ + η (1/N ∑_{n=1}^N Loss(f_{θ}(x_n + δ_n), y_n) - c)]_+

```



91

[Robey et al., NeurIPS'21]

[Robey et al., NeurIPS'21]

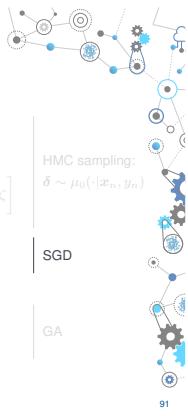
Dual Adversarial LEarning

```

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:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss} \left( f_{\theta_t}(\mathbf{x}_n + \delta_n), y_n \right) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss} \left( f_{\theta}(\mathbf{x}_n), y_n \right) + \lambda \text{Loss} \left( f_{\theta}(\mathbf{x}_n + \delta_n), y_n \right) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss} \left( f_{\theta}(\mathbf{x}_n + \delta_n), y_n \right) - c \right) \right]_+$ 

```

[Robey et al., NeurIPS'21]



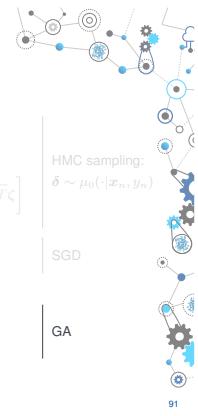
Dual Adversarial LEarning

```

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:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss} \left( f_{\theta_t}(\mathbf{x}_n + \delta_n), y_n \right) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss} \left( f_{\theta}(\mathbf{x}_n), y_n \right) + \lambda \text{Loss} \left( f_{\theta}(\mathbf{x}_n + \delta_n), y_n \right) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss} \left( f_{\theta}(\mathbf{x}_n + \delta_n), y_n \right) - c \right) \right]_+$ 

```

[Robey et al., NeurIPS'21]



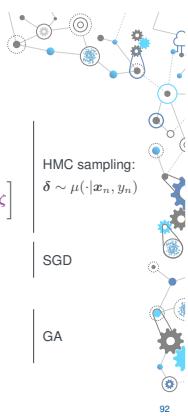
Dual Adversarial LEarning

```

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:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss} \left( f_{\theta_t}(\mathbf{x}_n + \delta_n), y_n \right) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss} \left( f_{\theta}(\mathbf{x}_n), y_n \right) + \lambda \text{Loss} \left( f_{\theta}(\mathbf{x}_n + \delta_n), y_n \right) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss} \left( f_{\theta}(\mathbf{x}_n + \delta_n), y_n \right) - c \right) \right]_+$ 

```

[Robey et al., NeurIPS'21]



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:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss} \left( f_{\theta_t}(\mathbf{x}_n + \delta_n), y_n \right) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss} \left( f_{\theta}(\mathbf{x}_n), y_n \right) + \lambda \text{Loss} \left( f_{\theta}(\mathbf{x}_n + \delta_n), y_n \right) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss} \left( f_{\theta}(\mathbf{x}_n + \delta_n), y_n \right) - c \right) \right]_+$ 

```

[Robey et al., NeurIPS'21]



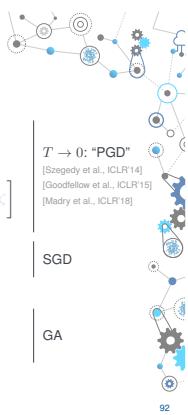
Dual Adversarial LEarning

```

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:      $\delta_n \leftarrow \text{proj}_{\Delta} \left[ \delta_n + \eta \text{sign} \left[ \nabla_{\delta} \log \left( \text{Loss} \left( f_{\theta_t}(\mathbf{x}_n + \delta_n), y_n \right) \right) \right] + \sqrt{2\eta T} \zeta \right]$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss} \left( f_{\theta}(\mathbf{x}_n), y_n \right) + \lambda \text{Loss} \left( f_{\theta}(\mathbf{x}_n + \delta_n), y_n \right) \right]$ 
8: end
9:  $\lambda \leftarrow \left[ \lambda + \eta \left( \frac{1}{N} \sum_{n=1}^N \text{Loss} \left( f_{\theta}(\mathbf{x}_n + \delta_n), y_n \right) - c \right) \right]_+$ 

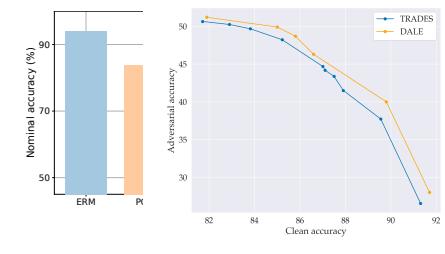
```

[Robey et al., NeurIPS'21]

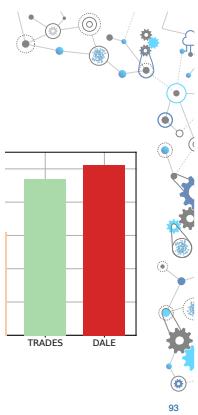


Dual Adversarial LEarning

Problem
Learn an image classifier that

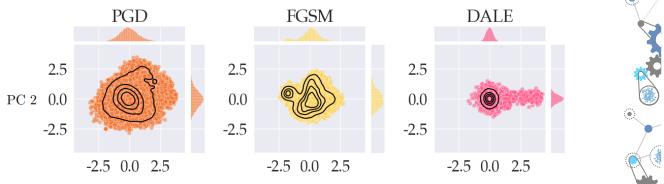


[Robey et al., NeurIPS'21]



Dual Adversarial LEarning

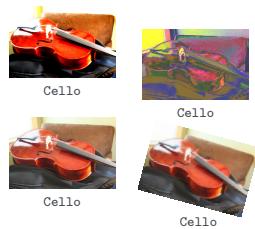
Problem
Learn an image classifier that is robust to input perturbations



94

Invariance

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



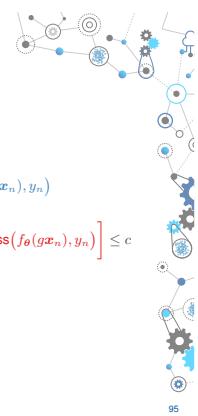
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss} \left(f_{\theta}(\mathbf{x}_n), y_n \right)$$

subject to

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

[Robey et al., NeurIPS'21]

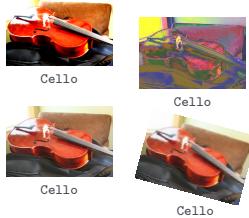
[Hounie et al., ICML'23]



Invariance

Problem

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



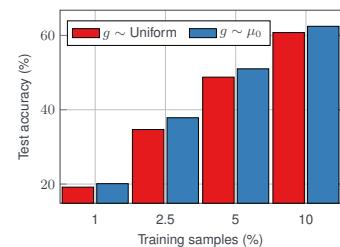
[Hounie et al., ICML'23]

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

- No differentiability required (e.g., Metropolis-Hastings)

95

Training on a subset of ImageNet-100



[Hounie et al., ICML'23]

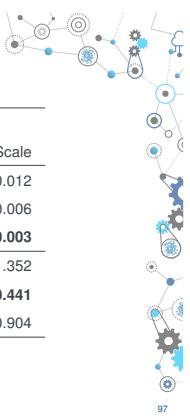


96

“Identifying” invariances

		Synthetic Invariance		
Dataset	Dual variable (λ)	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 et al., ICML'23]



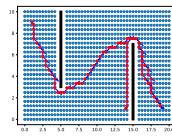
97

(Manifold) smoothness

Problem

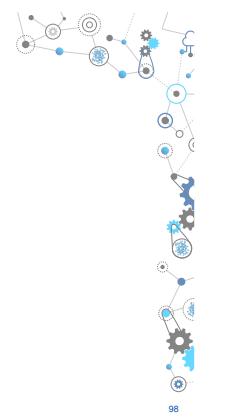
Leveraging unlabeled data

- Labeled data $\{\langle \text{Position}, \text{Action} \rangle\}$



Dataset

[Cerviño et al., ICML'23]



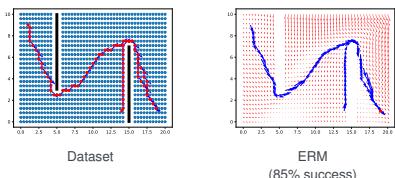
98

(Manifold) smoothness

Problem

Leveraging unlabeled data

- Labeled data $\{\langle \text{Position}, \text{Action} \rangle\}$



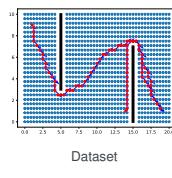
98

(Manifold) smoothness

Problem

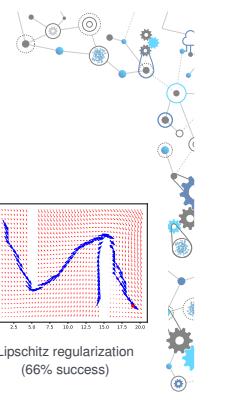
Leveraging unlabeled data

- Labeled data $\{\langle \text{Position}, \text{Action} \rangle\}$



Dataset

[Cerviño et al., ICML'23]



98

(Manifold) smoothness

Problem

Leveraging unlabeled data

- Labeled data $\{\langle \text{Position}, \text{Action} \rangle\}$ and unlabeled data $\{\langle \text{Position} \rangle\}$
- Use $\{\text{Position}\}$ to estimate a data manifold \mathcal{M}

$$\begin{aligned} & \min_{\theta} \frac{1}{N} \sum_{n=1}^N \|f_{\theta}(\text{Position}_n) - \text{Action}_n\|^2 \\ \text{subject to} \quad & \max_{\mathbf{x}} \|\nabla_{\mathcal{M}} f_{\theta}(\mathbf{x})\|^2 \leq c \end{aligned}$$

99

(Manifold) smoothness

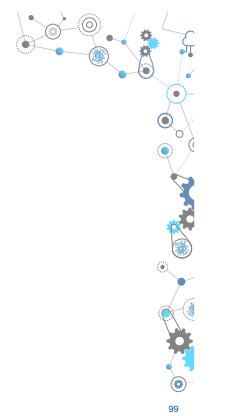
Problem

Leveraging unlabeled data

- Labeled data $\{\langle \text{Position}, \text{Action} \rangle\}$ and unlabeled data $\{\langle \text{Position} \rangle\}$
- Use $\{\text{Position}\}$ to estimate a data manifold \mathcal{M}

$$\begin{aligned} & \min_{\theta} \frac{1}{N} \sum_{n=1}^N \|f_{\theta}(\text{Position}_n) - \text{Action}_n\|^2 \\ \text{subject to} \quad & \max_{\mathbf{x} \sim \mu_0} \|\nabla_{\mathcal{M}} f_{\theta}(\mathbf{x})\|^2 \leq c \end{aligned}$$

[Cerviño et al., ICML'23]

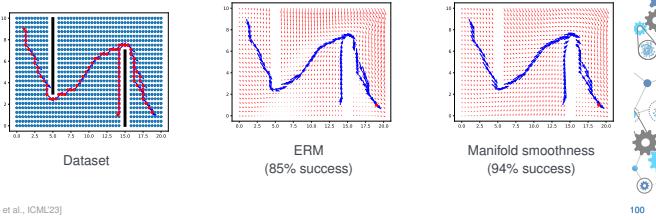


99

(Manifold) smoothness

Problem
Leveraging unlabeled data

- Labeled data ({Position, Action}) **and** unlabeled data ({Position})



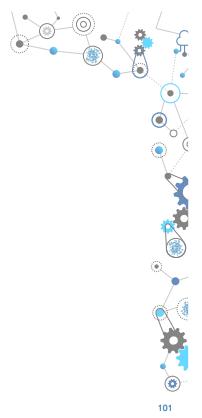
[Cerviño et al., ICML'23]

Agenda

Adversarially robust learning

Semi-infinite learning

Probabilistic robustness



Constrained learning challenges

$$\begin{aligned} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) &\xrightarrow{\text{PAC}} \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)] \\ \text{s. to } \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \ell(f_{\theta}(x_n + \delta), y_n) \right] \leq c &\xrightarrow{\text{PACC}} \text{s. to } \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\delta \in \Delta} \ell(f_{\theta}(x + \delta), y) \right] \leq c \end{aligned}$$

Challenges

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

102

Constrained learning challenges

$$\begin{aligned} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n) &\xrightarrow{\text{PAC}} \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)] \\ \text{s. to } \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \ell(f_{\theta}(x_n + \delta), y_n) \right] \leq c &\xrightarrow{\text{PACC?}} \text{s. to } \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\delta \in \Delta} \ell(f_{\theta}(x + \delta), y) \right] \leq c \end{aligned}$$

Challenges

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

102

Statistical complexity

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

- Is robust learning harder than non-robust learning? Do we need more samples?

A: YES AND NO

[Cullina, Bhagat, Mittal. PAC-learning in the presence of evasion adversaries, NeurIPS'18]

[Yin, Ramchandran, Bartlett. Rademacher Complexity for Adversarially Robust Generalization, ICML'19]

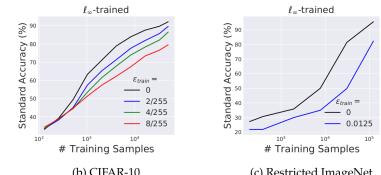
[Montasser, Hanneke, Srebro. VC Classes are Adversarially Robust Learnable, but Only Improperly, COLT'19]

[Awasthi, Frank, Mohri. Adversarial Learning Guarantees for Linear Hypotheses and Neural Networks, ICML'20]

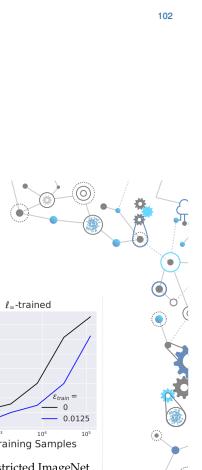
[Montasser, Hanneke, Srebro. Adversarially robust learning: A generic minimax optimal learner & characterization, NeurIPS'22]

103

Nominal performance of robust models



[Tsipras et al., ICLR'19]



“Softer” robustness

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

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

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

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

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

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

105

“Softer” robustness

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

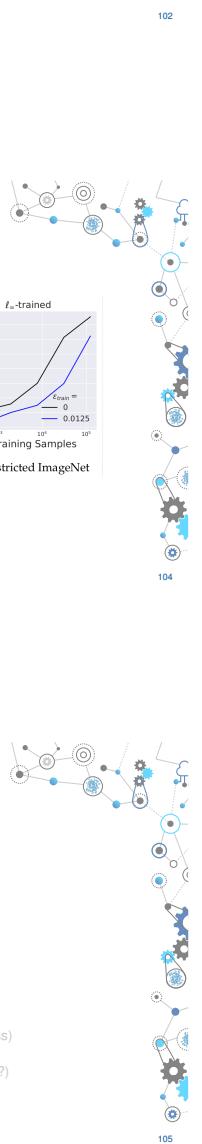
$$\min_{\theta} \mathbb{E}_{(x,y)} \left[\frac{1}{\tau} \log \left(\mathbb{E}_{\delta \sim m} \left[e^{\tau \cdot \text{Loss}(f_{\theta}(x+\delta), y)} \right] \right) \right]$$

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

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

Computationally challenging (especially as $\tau \rightarrow \infty$, i.e., stronger robustness)

No guaranteed advantages (lower sample complexity? improved trade-offs?)



Towards probabilistic robustness

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \\ \text{subject to} \quad & \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_1), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_2), y_n) \leq t(x_n, y_n) \\ & \text{Epigraph formulation:} \\ & \max_{\|\delta\|_{\infty} \leq \epsilon} \frac{\text{Loss}(f_{\theta}(x_n + \delta), y_n)}{\text{Loss}(f_{\theta}(x_n + \delta), y_n)} \geq \frac{t(x_n, y_n)}{\text{Loss}(f_{\theta}(x_n + \delta), y_n)} \leq t, \text{ for all } \|\delta\|_{\infty} \leq \epsilon \\ & \text{Semi-infinite program} \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^2}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^4}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{25}), y_n) \leq t(x_n, y_n) \end{aligned}$$

106

Towards probabilistic robustness

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \\ \text{subject to} \quad & \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_1), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_2), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^2}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^4}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^8}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{25}), y_n) \leq t(x_n, y_n) \end{aligned}$$

106

Towards probabilistic robustness

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \\ \text{subject to} \quad & \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_1), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^2}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^4}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^8}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{16}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{32}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{64}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{128}}), y_n) \leq t(x_n, y_n) \end{aligned}$$

107

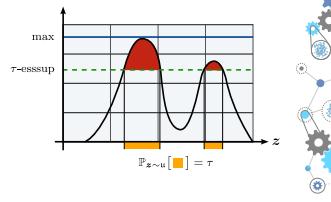
[Robey et al., ICML22 (spotlight)]

Probabilistic robustness

- Probabilistic robustness

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

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



108

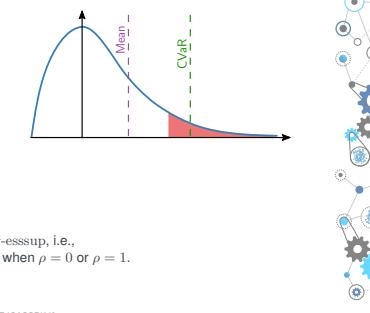
[Robey et al., ICML22 (spotlight)]

Probabilistic robustness and Risk

- Conditional value at risk:

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

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



110

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; Kalogerias et al., IEEE ICASSP'20]

Towards probabilistic robustness

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \\ \text{subject to} \quad & \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_1), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_2), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^2}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^4}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^8}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{16}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{32}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{64}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{128}}), y_n) \leq t(x_n, y_n) \end{aligned}$$

106

Towards probabilistic robustness

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \\ \text{subject to} \quad & \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_1), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_2), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^2}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^4}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^8}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{16}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{32}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{64}}), y_n) \leq t(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\pi^{128}}), y_n) \leq t(x_n, y_n) \end{aligned}$$

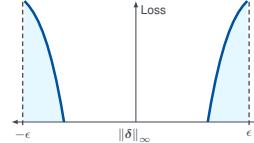
107

[Robey et al., ICML22 (spotlight)]

Probabilistic robustness

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

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



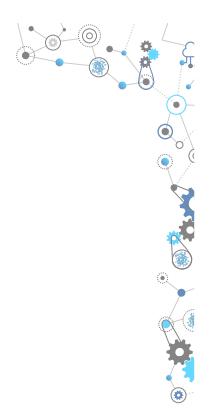
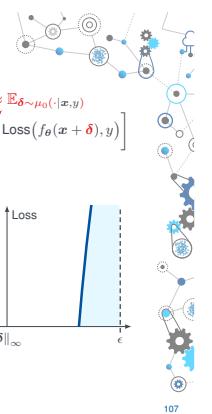
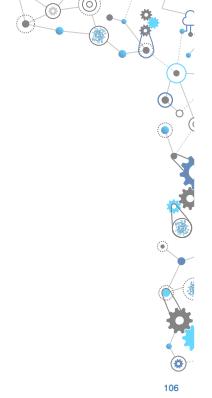
[Robey et al., ICML22 (spotlight)]

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:      $\alpha \leftarrow \alpha - \frac{\eta}{\tau} (\text{Loss}(f_{\theta}(x_n + \delta_t), y_n) - \alpha)$ 
6:   end
7:    $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[ \text{Loss}(f_{\theta}(x_n + \delta_T), y_n) - \alpha \right]_+$ 
      $\approx \text{CVaR}_{1-\tau} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]$ 
8: end
  
```

[Robey et al., ICML22 (spotlight)]

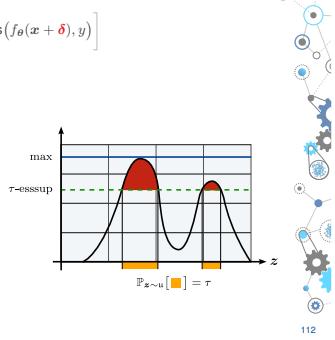


111

Probabilistic robustness

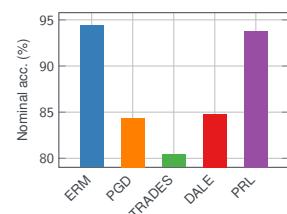
- Probabilistic robustness

$$\min_{\theta} \mathbb{E}_{(x,y)} \left[\underset{\delta \in \Delta}{\tau\text{-essup}} \text{Loss}(f_{\theta}(x + \delta), y) \right]$$
 - $\tau = 1/2$: classical learning (for symmetric m)
 - $\tau = 0$: adversarial robustness (ess sup)

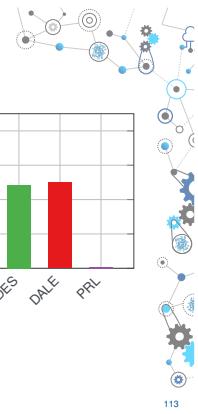


[Robey et al., ICML22 (spotlight)]

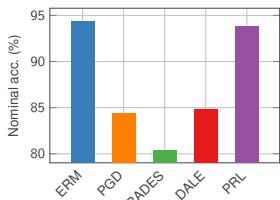
Probabilistically robust learning



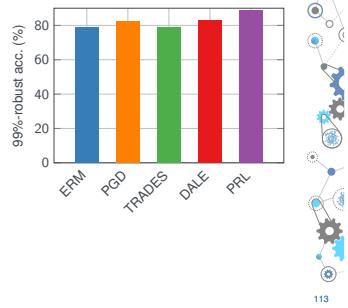
[Robey et al., ICML22 (spotlight)]



Probabilistically robust learning



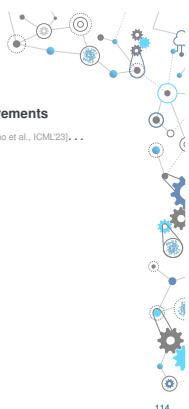
[Robey et al., ICML22 (spotlight)]



113

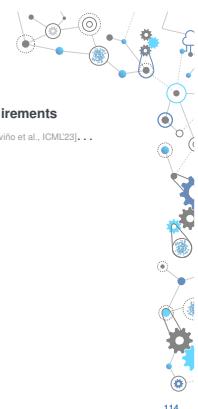
Summary

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



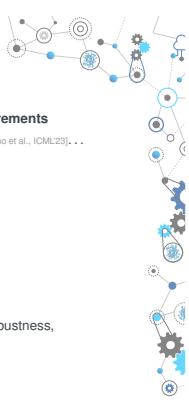
Summary

- Semi-infinite constrained learning is the a tool to enforce worst-case requirements
- 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 [Houne et al., ICML'23], smoothness [Cerviño 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]



Break