

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

EUSIPCO tutorial
Aug. 26, 2024

supervised and reinforcement learning under requirements

Agenda

- I. Constrained supervised learning
 - Constrained learning theory
 - Resilient constrained learning
 - Robust learning

Break (30 min)

- II. Constrained reinforcement learning
 - Constrained RL duality
 - Constrained RL algorithms

Q&A and discussions

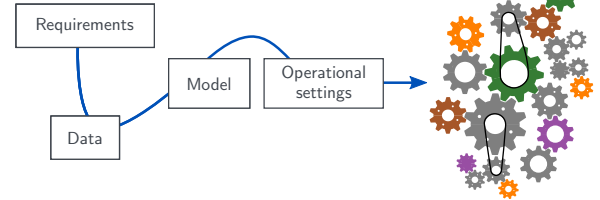


<https://luizchamon.com/eusipco>

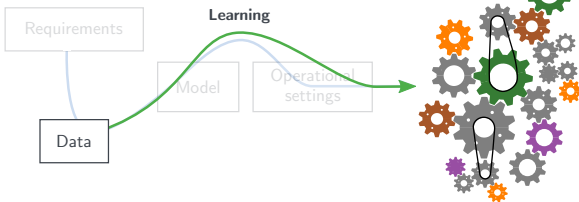
Why learning under requirements?



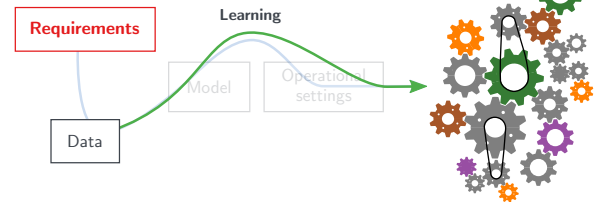
Why learning under requirements?



Why learning under requirements?

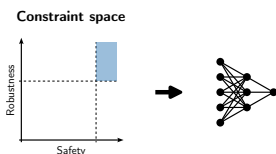


Why learning under requirements?



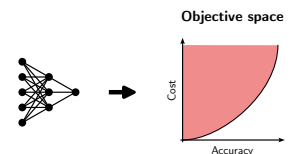
What is a requirements?

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



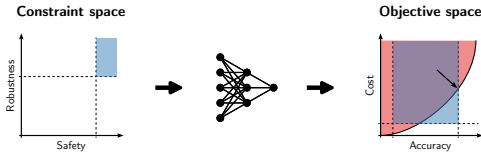
What is a requirements?

- 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



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]

3

What is (un)constrained learning?

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

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

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

4

What is (un)constrained learning?

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

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

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

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

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

4

What about penalties?

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

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

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

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

Applications

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

What about penalties?

NON-CONVEX

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

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

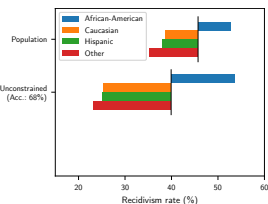
$$h(f_{\theta}(x), y) \leq u, \quad \mathbb{P}\text{-a.e.}$$

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

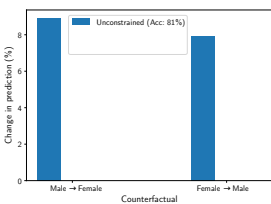
- There may not exist (λ, μ) such that the penalized solution is optimal and feasible
- Even if such (λ, μ) exist, they are not easy to find (hyperparameter search, cross-validation...)
- Constrained learning yields stronger guarantees, better performance, better trade-offs...

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.

7

Fairness: "Equality" of odds

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

$$\min_{\theta} \text{Prediction error}$$

$$\text{subject to } \text{Prediction rate disparity (Race)} \leq c,$$

$$\text{for Race} \in \{\text{African-American, Caucasian, Hispanic, Other}\}$$

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

8

Fairness: "Equality" of odds

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

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

subject to **Prediction rate disparity (Race) $\leq c$** ,
for Race $\in \{\text{African-American, Caucasian, Hispanic, Other}\}$

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset. [Goh et al., NeurIPS16; Kearns et al., ICML18; Cotter et al., JMLR19; Chamon et al., IEEE TIT23]

8

Fairness: "Equality" of odds

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

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

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

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset. [Goh et al., NeurIPS16; Kearns et al., ICML18; Cotter et al., JMLR19; Chamon et al., IEEE TIT23]

8

Counterfactual fairness

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

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

subject to **Change in prediction (ρx) $\leq c$ a.e.**
(ρ : Male \leftrightarrow Female)

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset. [Chamon and Ribeiro, NeurIPS20]

9

Counterfactual fairness

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

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

subject to $D_{\text{KL}}(f_{\theta}(x_n) \| f_{\theta}(\rho x_n)) \leq c$, for all n
(ρ : Male \leftrightarrow Female)

* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset. [Chamon and Ribeiro, NeurIPS20]

9

Applications

- Fairness (e.g., [Goh et al., NeurIPS16; Kearns et al., ICML18; Cotter et al., JMLR19; Chamon et al., IEEE TIT23])
- Federated learning (e.g., [Shen et al., ICLR22; Hounie et al., NeurIPS23])
- Adversarially robust learning (e.g., [Chamon et al., NeurIPS20; Robey et al., NeurIPS21; Chamon et al., IEEE TIT23])
- Safe learning (e.g., [Paternain et al., IEEE TAC23])
- Wireless resource allocation (e.g., [Eisen et al., IEEE TSP19; NaderiAlizadeh et al., IEEE TSP22; Chowdhury et al., Asilomar23])
- ...

10

Federated learning

Problem
Learn a common model using data from K clients

\min_{θ} Average loss across clients



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

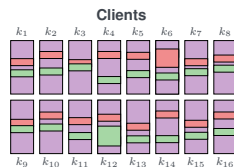
[Shen et al., ICLR22]

11

Federated learning

Problem
Learn a common model using data from K clients

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



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

[Shen et al., ICLR22]

11

Federated learning

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

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



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

[Shen et al., ICLR22]

11

Federated learning

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

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

subject to $\text{Loss disparity (} k\text{-th client)} \leq c,$
 $k = 1, \dots, K$



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

[Shen et al., ICRL22]

11

Federated learning

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

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

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



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

[Shen et al., ICRL22]

11

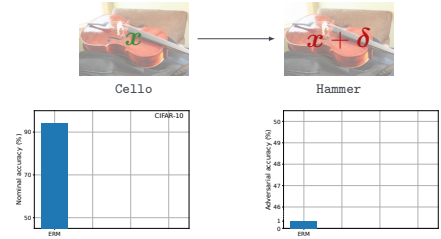
Applications

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

12

Robustness

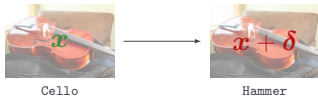
Problem
Learn an accurate classifier that is robust to input perturbations



13

Robustness

Problem
Learn an accurate classifier that is robust to input perturbations



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

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

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

13

Robustness

Problem
Learn an accurate classifier that is robust to input perturbations



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

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

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

13

Robustness

Problem
Learn an accurate classifier that is robust to input perturbations



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

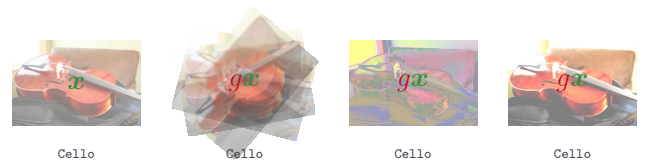
subject to $\frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right] \leq c$

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

13

Invariance

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



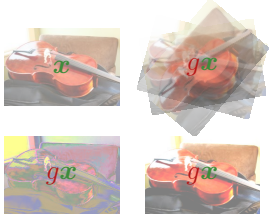
[Hounie, Chamon, Ribeiro, NeurIPS'23]

14

Invariance

Problem

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



$$\begin{aligned} & \min_{\theta} \text{Prediction error} \\ & \text{subject to } \text{Variance} \leq c \end{aligned}$$

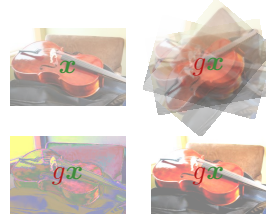
[Hounie, Chamon, Ribeiro, NeurIPS'23]

15

Invariance

Problem

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



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

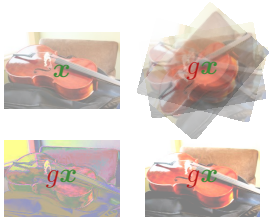
[Hounie, Chamon, Ribeiro, NeurIPS'23]

15

Invariance

Problem

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



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

[Hounie, Chamon, Ribeiro, NeurIPS'23]

15

Applications

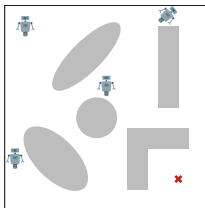
- Fairness (e.g., [Goh et al., NeurIPS'16; Kearns et al., ICML'18; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23])
- Federated learning (e.g., [Shen et al., ICLR'22; Hounie et al., NeurIPS'23])
- Adversarially robust learning (e.g., [Chamon et al., NeurIPS'20; Robey et al., NeurIPS'21; Chamon et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

16

Safety

Problem

Find a control policy that navigates the environment effectively and safely

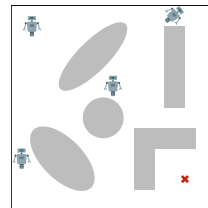


17

Safety

Problem

Find a control policy that navigates the environment effectively and safely



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

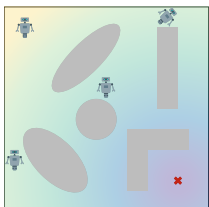
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

18

Safety

Problem

Find a control policy that navigates the environment effectively and safely



$$\begin{aligned} & \text{maximize}_{\pi \in \mathcal{P}(S)} \mathbb{E}_{s_0, a_0 \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \\ & \text{subject to } \mathbb{P}[\text{Colliding with } O_i] \leq \delta, \\ & \text{for } i = 1, 2, \dots \end{aligned}$$

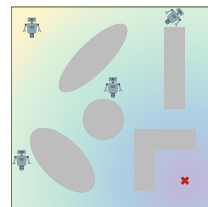
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

18

Safety

Problem

Find a control policy that navigates the environment effectively and safely



$$\begin{aligned} & \text{maximize}_{\pi \in \mathcal{P}(S)} \mathbb{E}_{s_0, a_0 \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \\ & \text{subject to } \mathbb{P} \left(\bigcap_{t=0}^{T-1} \{s_t \notin O_i\} \mid \pi \right) \geq 1 - \delta_i, \\ & \text{for } i = 1, 2, \dots \end{aligned}$$

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

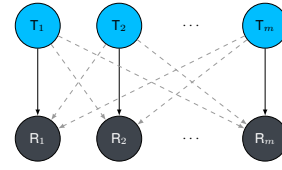
18

Applications

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

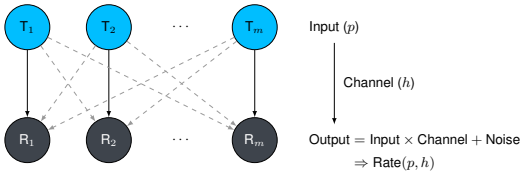
Wireless resource allocation

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



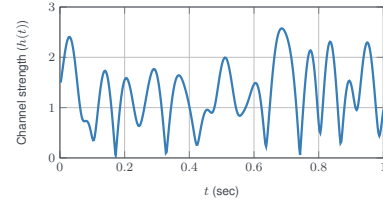
Wireless resource allocation

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



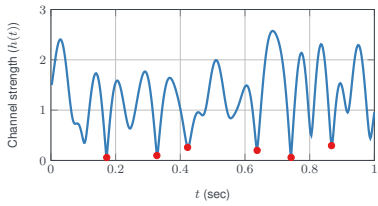
Wireless resource allocation

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



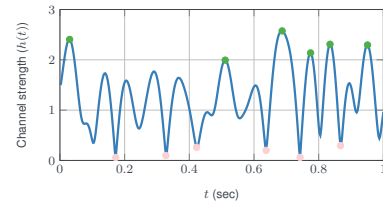
Wireless resource allocation

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



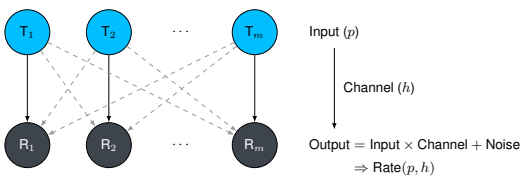
Wireless resource allocation

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



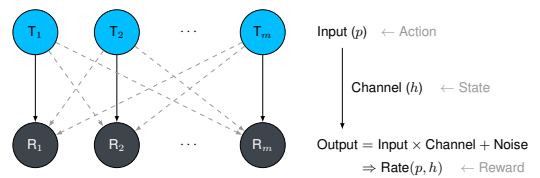
Wireless resource allocation

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



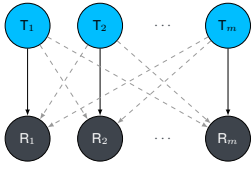
Wireless resource allocation

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



Wireless resource allocation

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



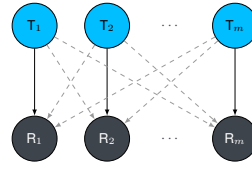
$$\begin{aligned} \min_{\pi \in \mathcal{P}(\mathcal{S})} & \text{Total transmit power} \\ \text{s. to} & \text{Rate } T_i \rightarrow R_i \geq c_i \end{aligned}$$

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

23

Wireless resource allocation

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



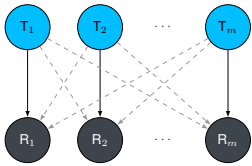
$$\begin{aligned} \min_{\pi \in \mathcal{P}(\mathcal{S})} & \sum_{i=1}^m \mathbb{E}_{h, p \sim \pi(h)} \left[\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^m p_{i,t} \right] \\ \text{s. to} & \text{Rate } T_i \rightarrow R_i \geq c_i \end{aligned}$$

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

23

Wireless resource allocation

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



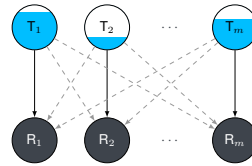
$$\begin{aligned} \min_{\pi \in \mathcal{P}(\mathcal{S})} & \sum_{i=1}^m \mathbb{E}_{h, p \sim \pi(h)} \left[\frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^m p_{i,t} \right] \\ \text{s. to} & \mathbb{E}_{h, p \sim \pi(h)} \left[\frac{1}{T} \sum_{t=0}^{T-1} \text{Rate}_i(p_t, h_t) \right] \geq c_i \end{aligned}$$

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

23

Wireless resource allocation

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



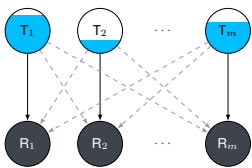
$$\begin{aligned} \min_p & \text{Total probability of depleting battery} \\ \text{s. to} & \mathbb{E}_{h, p \sim \pi(h, b)} \left[\frac{1}{T} \sum_{t=0}^{T-1} \text{Rate}_i(p_t, h_t) \right] \geq c_i \end{aligned}$$

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

23

Wireless resource allocation

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



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

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

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; Dwivedi et al., arXiv'24])
- Continual learning (e.g., [Peng et al., ICML'23])
- Active learning (e.g., [Elentner et al., NeurIPS'22])
- Semi-supervised learning (e.g., [Cervino et al., ICML'23])
- Minimum norm interpolation, SVM...

24

Constrained supervised learning

What is (un)constrained learning?

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

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

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

26

Constrained learning challenges

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

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

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

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

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

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

Challenges

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

27

Constrained learning challenges

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

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

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

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

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

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

Challenges

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

27

Constrained learning challenges

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

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

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

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

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

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

Challenges

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

27

Agenda

Constrained learning theory

Constrained learning algorithms

Resilient constrained learning

28

Constrained learning challenges

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

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

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

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

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

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

Challenges

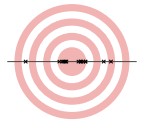
- 1) *Statistical*: does the solution of the constrained empirical problem generalize?
- 2) *Computational*: can we solve the constrained empirical problem?

29

What classical learning theory says?

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

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



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

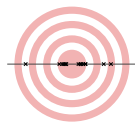
30

What classical learning theory says?

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

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

Requirements?



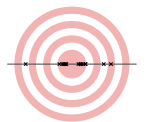
What's in a solution?

Definition (PAC learnability)

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

- near-optimal

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



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

30

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

31

What's in a solution?

Definition (PACC learnability)

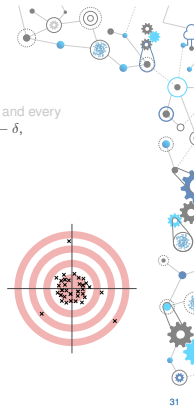
f_θ is a *probably approximately correct constrained (PACC)* learnable if for every ϵ, δ and every distributions \mathcal{D}, \mathcal{Q} , we can obtain f_{θ^*} from $N_{f_\theta}(\epsilon, \delta)$ samples such that, with prob. $1 - \delta$,

- near-optimal

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

- approximately feasible

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



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

31

When is constrained learning possible?

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

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

Proposition

f_θ is PAC learnable \Rightarrow f_θ is PACC learnable

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

32

ECRM is not a PACC learner

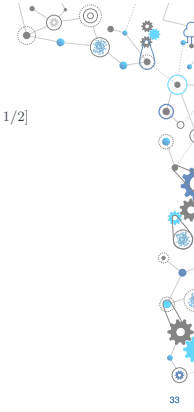
Counter-example

$$P^* = \min_{\theta \in \Theta} J(\theta)$$

subject to $\theta_2 \mathbb{E}_\tau[r] \leq \theta_1 - 1$
 $-\theta_1 \mathbb{E}_\tau[r] \leq \theta_2 - 1$

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

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



33

ECRM is not a PACC learner

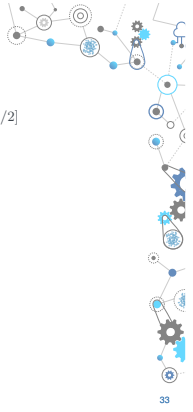
Counter-example

$$P^* = \min_{\theta \in \Theta} J(\theta) = \frac{1}{8}$$

subject to $\theta_2 \mathbb{E}_\tau[r] \leq \theta_1 - 1 \Rightarrow \theta_1 \geq 1$
 $-\theta_1 \mathbb{E}_\tau[r] \leq \theta_2 - 1 \Rightarrow \theta_2 \leq 1$

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

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



33

ECRM is not a PACC learner

Counter-example

$$P^* = \min_{\theta \in \Theta} J(\theta) = \frac{1}{8}$$

subject to $\theta_2 \mathbb{E}_\tau[r] \leq \theta_1 - 1 \Rightarrow \theta_1 \geq 1$
 $-\theta_1 \mathbb{E}_\tau[r] \leq \theta_2 - 1 \Rightarrow \theta_2 \leq 1$

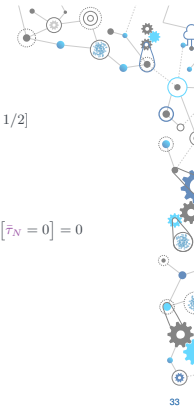
$$J(\theta) = \begin{cases} 1/16, & \theta = [1/2, 1/2] \\ 1/8, & \theta = [1, 1] \\ 1/4, & \theta = [1, 0] \end{cases}$$

$$\hat{P}^* = \min_{\theta \in \Theta} J(\theta)$$

subject to $\theta_2 \bar{\tau}_N \leq \theta_1 - 1$
 $-\theta_1 \bar{\tau}_N \leq \theta_2 - 1$

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

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



33

ECRM is not a PACC learner

Counter-example

$$P^* = \min_{\theta \in \Theta} J(\theta) = \frac{1}{8}$$

subject to $\theta_2 \mathbb{E}_\tau[r] \leq \theta_1 - 1 \Rightarrow \theta_1 \geq 1$
 $-\theta_1 \mathbb{E}_\tau[r] \leq \theta_2 - 1 \Rightarrow \theta_2 \leq 1$

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

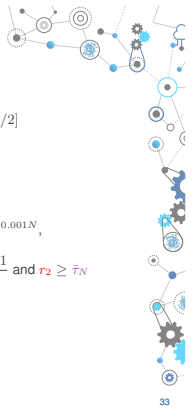
$$\hat{P}^* = \min_{\theta \in \Theta} J(\theta)$$

subject to $\theta_2 \bar{\tau}_N \leq \theta_1 - 1 + r_1$
 $-\theta_1 \bar{\tau}_N \leq \theta_2 - 1 + r_2$

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

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

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



33

Constrained learning challenges

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

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

$\xrightarrow{\text{PAC}}$

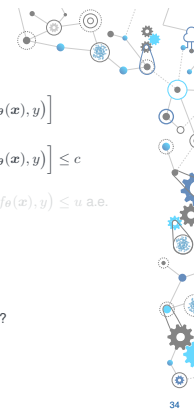
$$P^* = \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)]$$

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

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

Challenges

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



34

Constrained learning challenges

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

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

$\xrightarrow{\text{PAC}}$

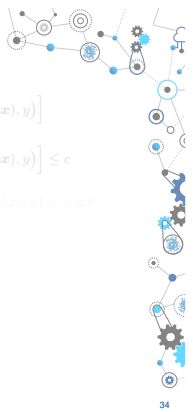
$$P^* = \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)]$$

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

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

Challenges

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



34

Duality

PRIMAL
↕
DUAL



35

Duality

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

↕

DUAL



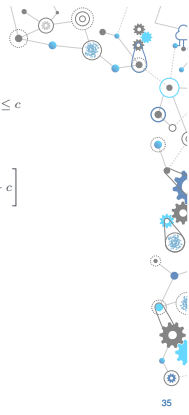
35

Duality

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

↕

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



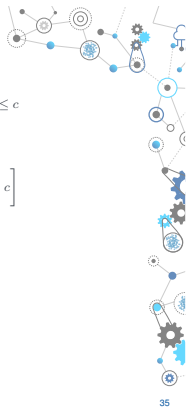
35

Duality

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

↕

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



35

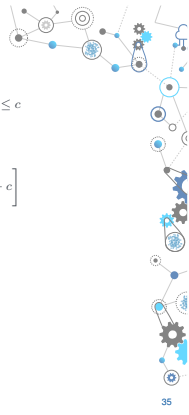
- In general, $\hat{D}^* \leq \hat{P}^*$
- But in some cases, $\hat{D}^* = \hat{P}^*$ (strong duality) [e.g., convex optimization]

Duality

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

↕

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



35

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

An alternative path

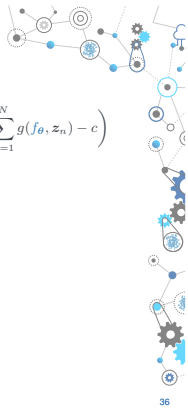
$$\hat{P}^* = \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}, z_n) \text{ s.t. } \frac{1}{N} \sum_{n=1}^N g(f_{\theta}, z_n) \leq c$$

↔

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

↑ PAC

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



36

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

An alternative path

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

↔

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

↑ PAC

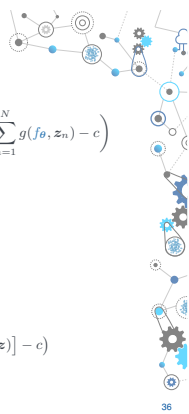
$$P^* = \min_{\theta \in \Theta} \mathbb{E}_z [\ell(f_{\theta}, z)] \text{ s.t. } \mathbb{E}_z [g(f_{\theta}, z)] \leq c$$

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

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

↔

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



36

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

Non-convex variational duality

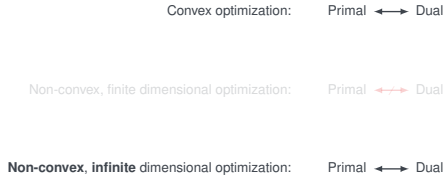
Convex optimization: Primal ↔ Dual

Non-convex, finite dimensional optimization: Primal ↔ Dual



37

Non-convex variational duality



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

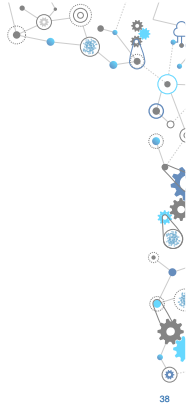
37

Sparse logistic regression

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

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

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



38

Sparse logistic regression

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

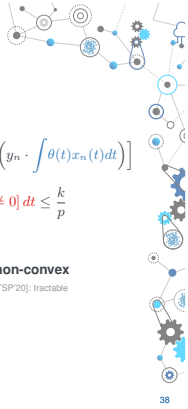
$$\text{s. to } \|\theta\|_0 = \sum_{i=1}^p \mathbb{I}[\theta_i \neq 0] \leq k$$

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

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

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

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



38

Sparse logistic regression

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

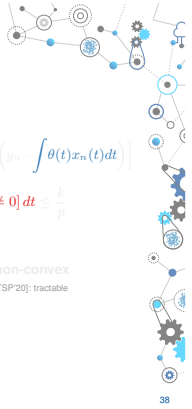
$$\text{s. to } \|\theta\|_0 = \sum_{i=1}^p \mathbb{I}[\theta_i \neq 0] \leq k$$

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

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

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

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



38

An alternative path

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

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

↓ PAC

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

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

↓ PAC

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

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

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

39

An alternative path

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

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

↓ PAC

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

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

↓ PAC

$$\hat{P}^* = \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] \quad \xleftarrow{\epsilon_0} \quad \hat{D}^* = \max_{\lambda \geq 0} \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] + \lambda (\mathbb{E}_z [g(\phi, z)] - c)$$

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

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

39

An alternative path

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

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

↓ PAC

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

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

↓ PAC

$$\hat{P}^* = \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] \quad \xleftarrow{\epsilon_0} \quad \hat{D}^* = \max_{\lambda \geq 0} \min_{\phi \in \mathcal{H}} \mathbb{E}_z [\ell(\phi, z)] + \lambda (\mathbb{E}_z [g(\phi, z)] - c)$$

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

↓ PAC

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

39

Dual (near-)PACC learning

Theorem

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

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

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



40

Dual (near-)PACC learning

Theorem

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

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

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

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

(mild additional conditions apply)

[Chamon and Ribeiro, NeurIPS'20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentor, Chamon, Ribeiro, ICLR'24]

40

Dual (near-)PACC learning

Theorem

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

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

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

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

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

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

(mild additional conditions apply)

[Chamon and Ribeiro, NeurIPS'20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentor, Chamon, Ribeiro, ICLR'24]

40

Dual (near-)PACC learning

Theorem

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

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

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

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

Sources of error

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

[Chamon and Ribeiro, NeurIPS'20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentor, Chamon, Ribeiro, ICLR'24]

41

Dual (near-)PACC learning

Theorem

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

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

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

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

Sources of error

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

[Chamon and Ribeiro, NeurIPS'20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentor, Chamon, Ribeiro, ICLR'24]

41

Dual (near-)PACC learning

Theorem

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

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

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

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

Sources of error

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

[Chamon and Ribeiro, NeurIPS'20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentor, Chamon, Ribeiro, ICLR'24]

41

Dual (near-)PACC learning

Theorem

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

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

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

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

Sources of error

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

[Chamon and Ribeiro, NeurIPS'20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentor, Chamon, Ribeiro, ICLR'24]

41

Dual (near-)PACC learning

Theorem

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

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

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

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

Sources of error

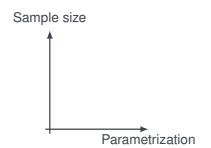
parametrization richness (ν) sample size (N) requirements difficulty (λ^*)

[Chamon and Ribeiro, NeurIPS'20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23; Elentor, Chamon, Ribeiro, ICLR'24]

41

Dual learning trade-offs

- Unconstrained learning
- parametrization \times sample size

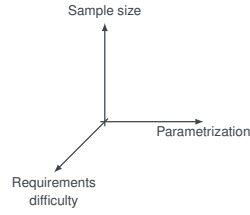


[Chamon and Ribeiro, NeurIPS'20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

42

Dual learning trade-offs

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



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

42

When is constrained learning possible?

Corollary

$$f_{\theta} \text{ is PAC learnable} \approx^* f_{\theta} \text{ is PACC learnable}$$

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

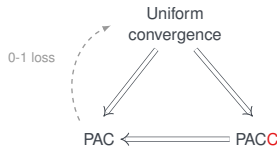
[mild conditions apply]

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

43

When is constrained learning possible?

Corollary



[mild conditions apply]

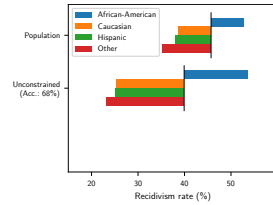
[Chamon and Ribeiro, NeurIPS'20; Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT'23]

43

Fairness

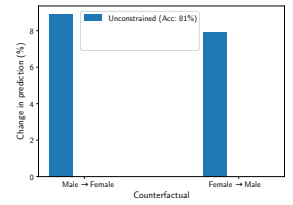
Problem

Predict whether an individual will recidivate



Problem

Predict whether an individual makes > \$50k



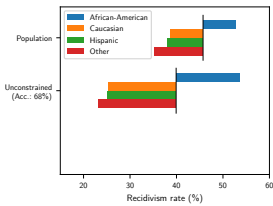
* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.

44

Fairness

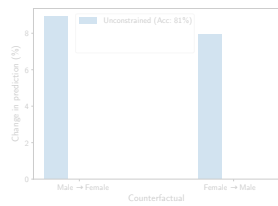
Problem

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.

44

Fairness: "Equality" of odds

Problem

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

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

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

45

Fairness: "Equality" of odds

Problem

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

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

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

45

Fairness: "Equality" of odds

Problem

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

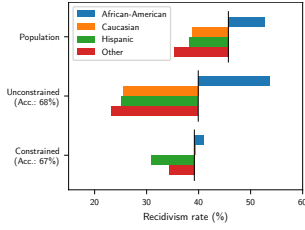
$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \\ \text{subject to} \quad & \frac{1}{N} \sum_{n=1}^N \sigma(f_{\theta}(x_n) - 0.5) \mathbb{1}[x_n \in \text{Race}] \leq \frac{1}{N} \sum_{n=1}^N \sigma(f_{\theta}(x_n) - 0.5) + c, \\ & \text{for Race}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

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

45

Fairness: "Equality" of odds

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

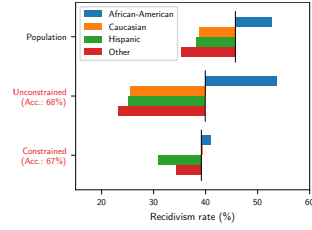


* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset. [Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT 23]

46

Fairness: "Equality" of odds

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



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

46

Fairness: "Equality" of odds

		Prediction	
		0	1
African-American	0	31%	16%
	1	16%	37%
Caucasian	0	52%	9%
	1	23%	16%

Unconstrained Constrained

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

47

Fairness: "Equality" of odds

		Prediction	
		0	1
African-American	0	31%	16%
	1	16%	37%
Caucasian	0	52%	9%
	1	23%	16%

Unconstrained Constrained

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

47

Fairness: "Equality" of odds

		Prediction	
		0	1
African-American	0	31%	16%
	1	16%	37%
Caucasian	0	52%	9%
	1	23%	16%

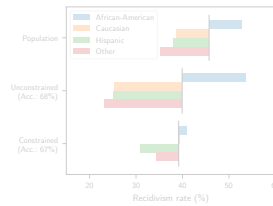
Unconstrained Constrained

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

47

Fairness

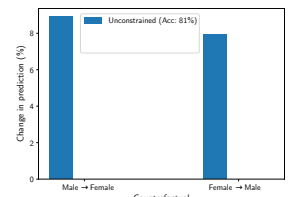
Problem
 Predict whether an individual will recidivate



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

48

Problem
 Predict whether an individual makes > \$50k



48

Counterfactual fairness

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

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

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

(ρ : Male ↔ Female)

[Chamon and Ribeiro, NeurIPS20]

49

Counterfactual fairness

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

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

subject to $\frac{1}{N} \sum_{n=1}^N D_{\text{KL}}(f_{\theta}(x_n) \| f_{\theta}(\rho x_n)) \leq c$, for all n

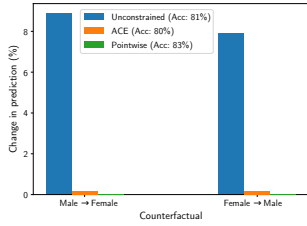
(ρ : Male ↔ Female)

[Chamon and Ribeiro, NeurIPS20]

49

Counterfactual fairness

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



[Chamion and Ribeiro, NeurIPS20]

50

Counterfactual fairness

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

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

subject to $D_{\text{KL}}(f_{\theta}(x_n) \| f_{\theta}(\rho x_n)) \leq c$, for all n
 (ρ : Male \leftrightarrow Female)

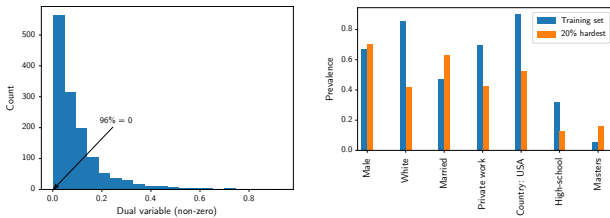
$$\max_{\lambda_n \geq 0} \min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) + \sum_{n=1}^N \lambda_n [D_{\text{KL}}(f_{\theta}(x_n) \| f_{\theta}(\rho x_n)) - c]$$

[Chamion and Ribeiro, NeurIPS20]

51

Counterfactual fairness

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



[Chamion and Ribeiro, NeurIPS20]

52

Agenda

Constrained learning theory

Constrained learning algorithms

Resilient constrained learning

53

Constrained optimization methods

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

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

[Chamion and Ribeiro, NeurIPS20]

54

Constrained optimization methods

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

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

[Chamion and Ribeiro, NeurIPS20]

54

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

Constrained optimization methods

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

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

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

[Chamion and Ribeiro, NeurIPS20]

54

Constrained optimization methods

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

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

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

[Chamion and Ribeiro, NeurIPS20]

54

Dual learning algorithm



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

55

Dual learning algorithm



- Minimize the primal (\equiv ERM)

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

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

55

Dual learning algorithm



- Minimize the primal (\equiv ERM)

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

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

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

55

Dual learning algorithm



- Minimize the primal (\equiv ERM)

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

- Update the dual

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

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

55

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

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

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

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

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

56

In practice...



- Minimize the primal (\equiv ERM)

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

- Update the dual

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

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

57

In practice...



- Minimize the primal (\equiv ERM)

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

- Update the dual

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

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

57

In practice...



- Initialize: θ_0, λ_0
- for $t = 1, \dots, T$
- $\beta_t \leftarrow \theta_{t-1}$
- for $n = 1, \dots, N$
- $\beta_{n+1} \leftarrow \beta_n - \eta \nabla_{\beta} \left[\ell(f_{\beta_n}(\mathbf{x}_n), y_n) + \lambda_{t-1} g(f_{\beta_n}(\mathbf{x}_n), y_n) \right]$
- end
- $\theta_t \leftarrow \beta_{N+1}$
- $\lambda_t = \left[\lambda_{t-1} + \eta \lambda \left(\frac{1}{N} \sum_{m=1}^N g(f_{\theta_t}(\mathbf{x}_m), y_m) - c \right) \right]_{+}$
- end
- Output: θ_T, λ_T

SGD

Dual update

PyTorch

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

58

In practice...

```

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

```

Use adaptive method (e.g., ADAM)

PyTorch

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

58

In practice...

```

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

```

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

PyTorch

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

58

In practice...

```

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

```

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

Use adaptive method (e.g., ADAM)

Use different time-scales ($\eta_{\lambda} = 0.1\eta$)

PyTorch

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

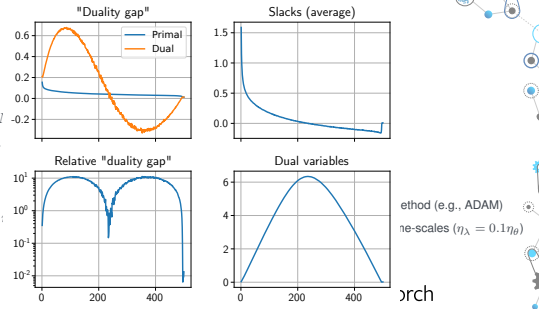
58

In practice...

```

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

```



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

PyTorch

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

58

Penalty-based vs. dual learning

Penalty-based learning

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

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

Dual learning

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

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

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

59

Agenda

Constrained learning theory

Constrained learning algorithms

Resilient constrained learning

60

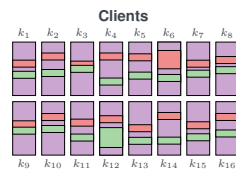
Heterogeneous federated learning

Problem

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

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

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



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

61

Heterogeneous federated learning

Problem

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

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

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



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

61

Resilient constrained learning

Definition (Resilience)

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

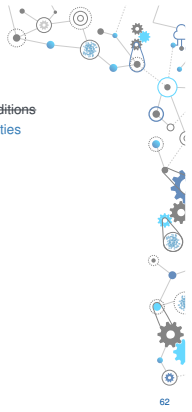


62

Resilient constrained learning

Definition (Resilience)

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



62

Resilient constrained learning

Definition (Resilience)

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

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

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



62

Resilient constrained learning

Definition (Resilience)

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

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

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



62

Resilient constrained learning

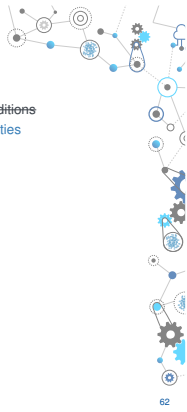
Definition (Resilience)

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

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

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

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



62

Resilient constrained learning

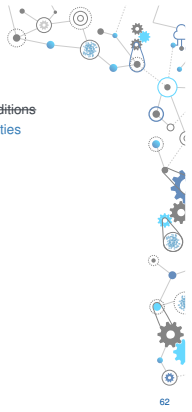
Definition (Resilience)

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

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

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

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



62

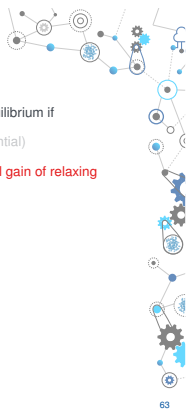
Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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



63

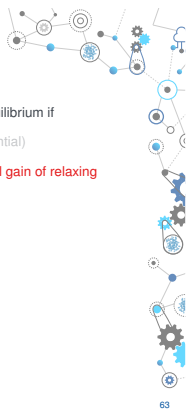
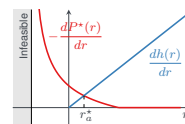
Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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



63

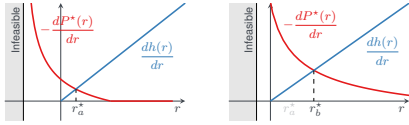
Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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



[Hounie, Chamon, Ribeiro, NeurIPS'23]

63

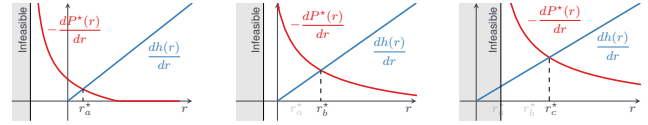
Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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



[Hounie, Chamon, Ribeiro, NeurIPS'23]

63

Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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

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

[Hounie, Chamon, Ribeiro, NeurIPS'23]

64

Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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

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

[Hounie, Chamon, Ribeiro, NeurIPS'23]

64

Resilient constrained learning

Definition (Resilient equilibrium)

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

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

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

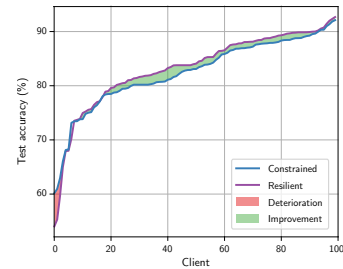
In words: at the resilient equilibrium the **marginal cost of relaxing** equals the **marginal gain of relaxing**

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

[Hounie, Chamon, Ribeiro, NeurIPS'23]

64

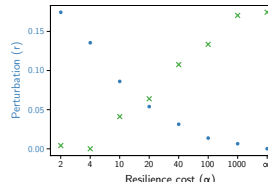
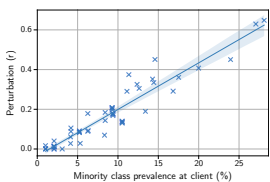
Heterogeneous federated learning



[Hounie, Chamon, Ribeiro, NeurIPS'23]

65

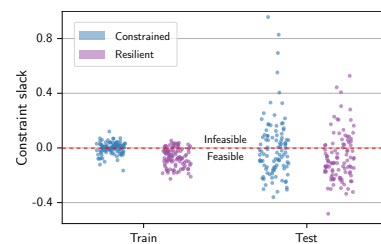
Heterogeneous federated learning



[Hounie, Chamon, Ribeiro, NeurIPS'23]

66

Heterogeneous federated learning



[Hounie, Chamon, Ribeiro, NeurIPS'23]

67

Summary

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



68

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



68

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?



68

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 [Hourie et al., NeurIPS'23]



68

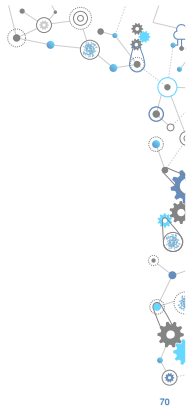
Robustness constraints

Agenda

Adversarially robust learning

Semi-infinite learning

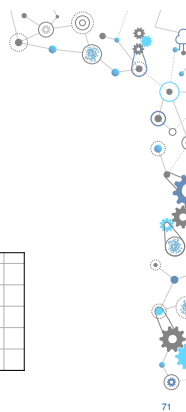
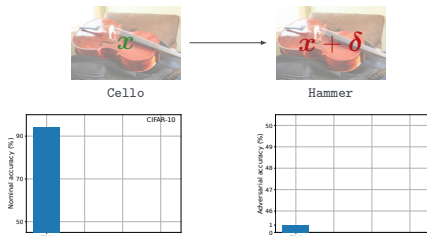
Probabilistic robustness



70

Robust learning

Problem
Learn an accurate classifier that is robust to input perturbations



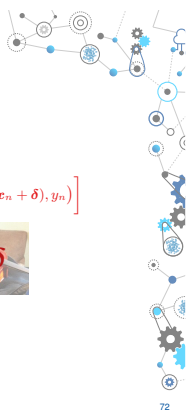
71

Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations

- Adversarial training [Szegedy et al., ICLR'14; Goodfellow et al., ICLR'15; Madry et al., ICLR'16; ...]

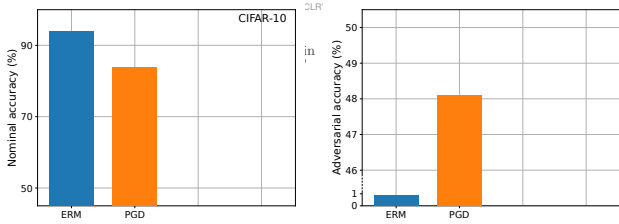
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \longrightarrow \min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$



72

Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations



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

72

Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations

- Adversarial training [Szegedy et al., ICLR14; Goodfellow et al., ICLR15; Madry et al., ICLR18; ...]

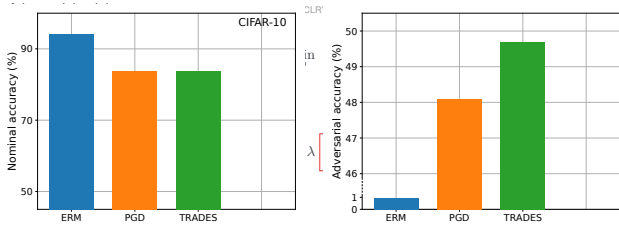
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n) \quad \min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

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

73

Adversarial training

Problem
Learn an accurate classifier that is robust to input perturbations



[Zhang et al., ICML19]

73

Constrained learning for robustness

Problem
Learn an accurate classifier that is robust to input perturbations

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

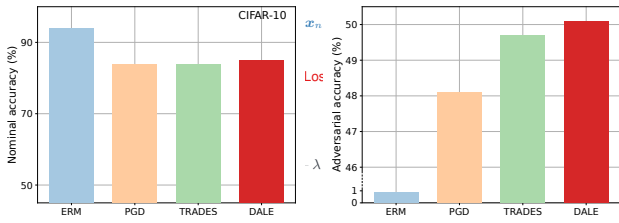
subject to $\frac{1}{N} \sum_{n=1}^N \left[\max_{\|\delta\|_{\infty} \leq \epsilon} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right] \leq c$

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

74

Constrained learning for robustness

Problem
Learn an accurate classifier that is robust to input perturbations

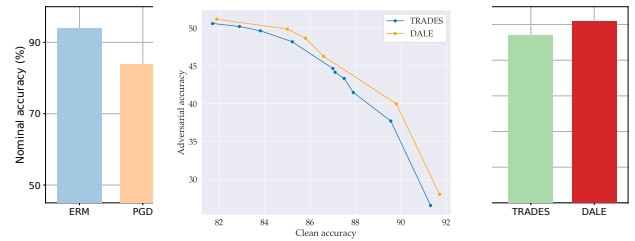


[Chamon and Ribeiro, NeurIPS20; Robey et al., NeurIPS21; Chamon et al., IEEE TIT23]

74

Constrained learning for robustness

Problem
Learn an accurate classifier that is robust to input perturbations



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

74

Penalty-based vs. dual learning

Penalty-based learning

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

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

Dual learning

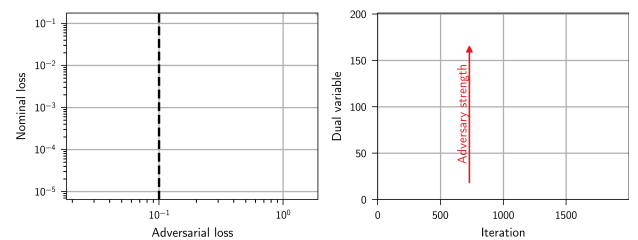
$$\theta^{\dagger} \in \text{argmin}_{\theta} \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$$

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

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

75

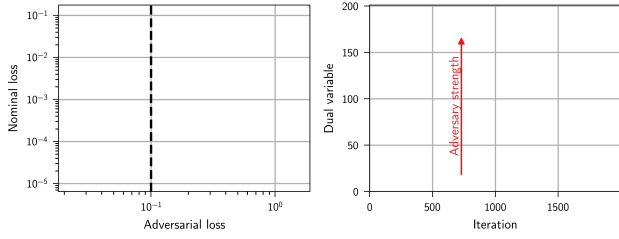
Constrained learning for robustness



[Chamon, Paternain, Calvo-Fullana, Ribeiro, IEEE TIT23]

76

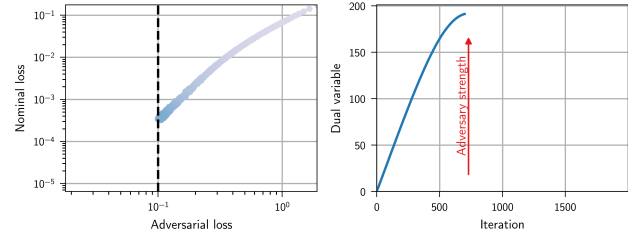
Constrained learning for robustness



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

76

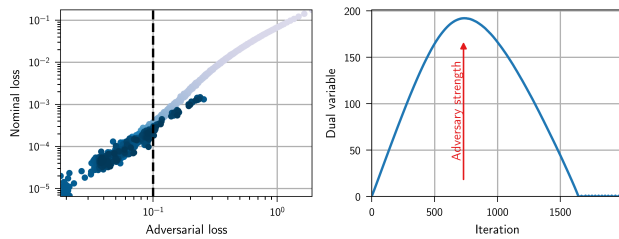
Constrained learning for robustness



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

76

Constrained learning for robustness



Empirical observations: [Zhang et al., ICML20; Sitawarin, arXiv'20]

76

Constrained learning for robustness

Problem

Learn an accurate 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

77

Constrained learning for robustness

Problem

Learn an accurate 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

77

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: $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \text{Loss}(f_{\theta}(x + \delta_t), y)$

78

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

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

78

Constrained learning for robustness

Problem

Learn an accurate 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

79

Agenda

Adversarially robust learning

Semi-infinite learning

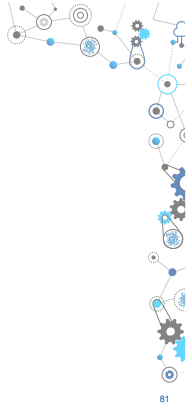
Probabilistic robustness



80

Semi-infinite constrained learning

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



81

Semi-infinite constrained learning

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

subject to $\text{Loss}(f_{\theta}(x_n + \delta), y_n) \leq t(x_n, y_n)$,
for all (x_n, y_n) and $\delta \in \Delta$

• Epigraph formulation:

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



81

Semi-infinite constrained learning

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

subject to $\text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq t(x_n, y_n)$
 $\text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n)$
 $\text{Loss}(f_{\theta}(x_n + \delta_{\epsilon}), y_n) \leq t(x_n, y_n)$
 $\text{Loss}(f_{\theta}(x_n + \delta_{\epsilon^*}), y_n) \leq t(x_n, y_n)$

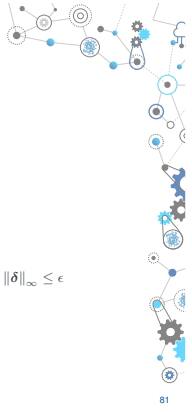
• Epigraph formulation:

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

• Semi-infinite program

$$\text{Loss}(f_{\theta}(x_n + \delta_{\epsilon^*}), y_n) \leq t(x_n, y_n)$$

$$\text{Loss}(f_{\theta}(x_n + \delta_{2\epsilon}), y_n) \leq t(x_n, y_n)$$



81

Duality

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

$$\iff \min_{\theta} \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \text{ s.t. } \text{Loss}(f_{\theta}(x_n + \delta), y_n) \leq t(x_n, y_n), \forall (x_n, y_n, \delta)$$

$$\iff \min_{\theta} \sup_{\mu \in \mathcal{P}} \frac{1}{N} \sum_{n=1}^N \int_{\Delta} \underbrace{\mu_n(\delta) \text{Loss}(f_{\theta}(x_n + \delta), y_n)}_{L(\theta, \mu)} d\delta$$



82

Duality

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

$$\iff \min_{\theta} \frac{1}{N} \sum_{n=1}^N [t(x_n, y_n)] \text{ s.t. } \text{Loss}(f_{\theta}(x_n + \delta), y_n) \leq t(x_n, y_n), \forall (x_n, y_n, \delta)$$

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



82

From optimization to sampling

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

$$\approx \min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbb{E}_{\delta \sim \mu_2} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]}_{L(\theta, \mu)}$$

Proposition

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

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



83

From optimization to sampling

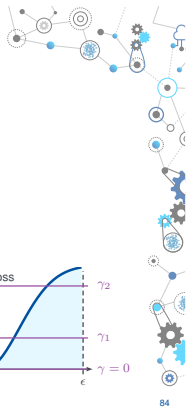
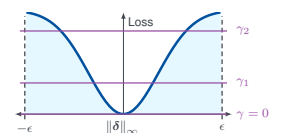
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

$$\approx \min_{\theta} \sup_{\mu \in \mathcal{P}^2} \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbb{E}_{\delta \sim \mu_2} [\text{Loss}(f_{\theta}(x_n + \delta), y_n)]}_{L(\theta, \mu)}$$

Proposition

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

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



84

From optimization to sampling

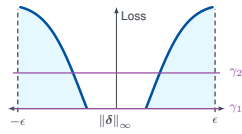
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

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

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

84

From optimization to sampling

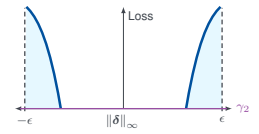
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

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

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

84

From optimization to sampling

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

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

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

84

From optimization to sampling

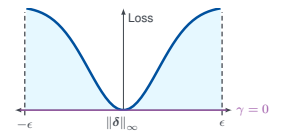
$$\min_{\theta} \frac{1}{N} \sum_{n=1}^N \left[\max_{\delta \in \Delta} \text{Loss}(f_{\theta}(x_n + \delta), y_n) \right]$$

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

Proposition

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

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



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

84

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

85

Constrained learning for robustness

Problem

Learn an image classifier that is robust to input perturbations

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

• Balancing nominal accuracy and robustness \Rightarrow Dual constrained learning

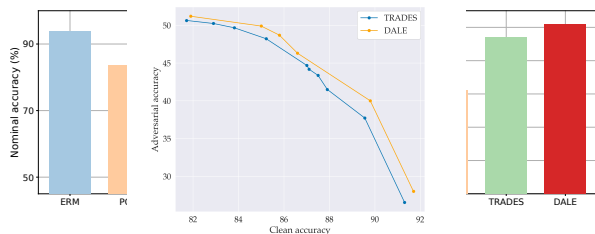
- Computing the worst-case perturbations
 - gradient ascent \rightarrow non-convex, underparametrized \Rightarrow sampling

85

Dual Adversarial Learning

Problem

Learn an image classifier th



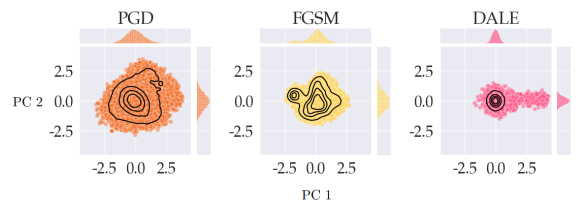
[Robey*, Chamon*, Pappas, Hassani, Ribeiro, NeurIPS21]

86

Dual Adversarial Learning

Problem

Learn an image classifier that is robust to input perturbations



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

87

Dual Adversarial Learning

```

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

```

HMC sampling:
 $\delta \sim \mu_{\theta}(\cdot | \mathbf{x}_n, y_n)$

SGD

GA

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

88

Dual Adversarial Learning

```

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

```

HMC sampling:
 $\delta \sim \mu_{\theta}(\cdot | \mathbf{x}_n, y_n)$

SGD

GA

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

88

Dual Adversarial Learning

```

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

```

HMC sampling:
 $\delta \sim \mu_{\theta}(\cdot | \mathbf{x}_n, y_n)$

SGD

GA

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

88

Dual Adversarial Learning

```

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

```

HMC sampling:
 $\delta \sim \mu_{\theta}(\cdot | \mathbf{x}_n, y_n)$

SGD

GA

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

88

Dual Adversarial Learning

```

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

```

HMC sampling:
 $\delta \sim \mu_{\theta}(\cdot | \mathbf{x}_n, y_n)$

SGD

GA

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

89

Dual Adversarial Learning

```

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

```

Gaussian
[Lopes et al., arXiv'19]

Patches
[Rusak et al., ECCV'20]

[Zhong et al., AAAI'20]

[Yun et al., ICCV'19]

...

SGD

GA

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

89

Dual Adversarial Learning

```

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

```

$T \rightarrow 0$: "PGD"
[Szegedy et al., ICLR14]
[Goodfellow et al., ICLR15]
[Madry et al., ICLR18]

SGD

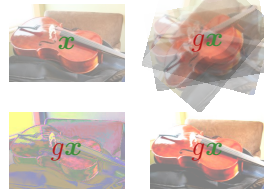
GA

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

89

Invariance

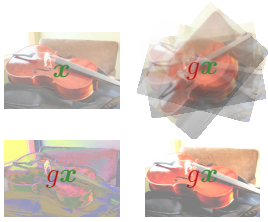
Problem
Learn a classifier that is invariant to transformation $g \in \mathcal{G}$



90

Invariance

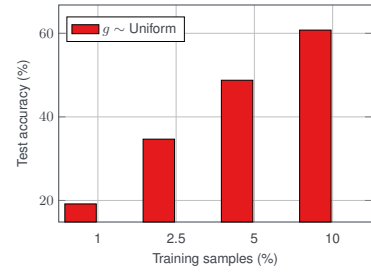
Problem
Learn a classifier that is invariant to transformation $g \in \mathcal{G}$



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

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

Training on a subset of ImageNet-100



Invariance

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

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

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

Invariance

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

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

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

Invariance

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

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

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

Invariance

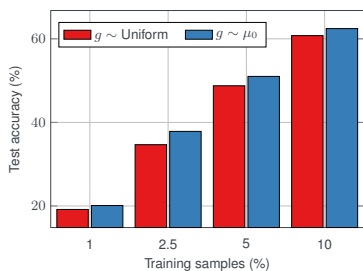
Problem
Learn a classifier that is invariant to transformation $g \in \mathcal{G}$

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

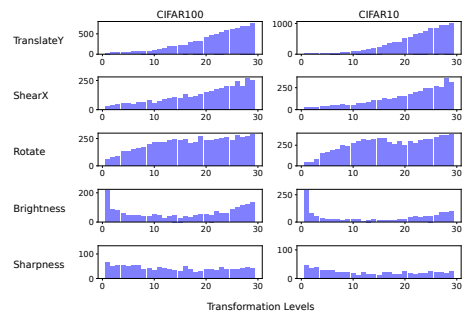
subject to $\frac{1}{N} \sum_{n=1}^N \mathbb{E}_{g \sim \mu_0(\cdot | x_n, y_n)} [\text{Loss}(f_{\theta}(g(x_n)), y_n)] \leq c$

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

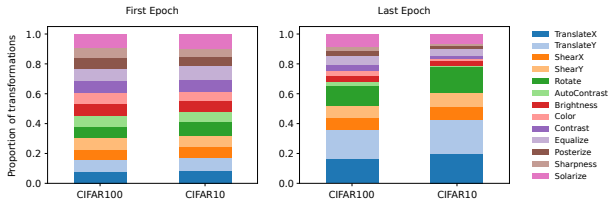
Training on a subset of ImageNet-100



Not all transformations are created equal



Not all transformations are created equal



95

“Identifying” invariances

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

[Hourie, Chamon, Ribeiro, ICML23]

96

Agenda

Adversarially robust learning

Semi-infinite learning

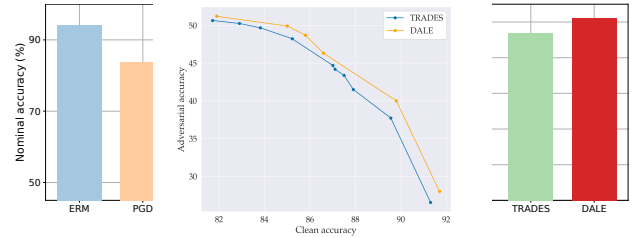
Probabilistic robustness

97

Constrained learning for robustness

Problem

Learn an **accurate** classifier



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

98

Constrained learning for robustness

Problem

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

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

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

99

“Softer” robustness

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

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

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

- L_p norms [Rice et al., NeurIPS21]

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

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

100

“Softer” robustness

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

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

- L_p norms [Rice et al., NeurIPS21]

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

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

100

Towards probabilistic robustness

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{N} \sum_{n=1}^N \ell(x_n, y_n) \\ \text{subject to} \quad & \text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_1), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\epsilon}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\epsilon/2}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\epsilon/4}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\epsilon/2^k}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{\epsilon^*}), y_n) \leq \ell(x_n, y_n) \\ & \text{Loss}(f_{\theta}(x_n + \delta_{2\epsilon^*}), y_n) \leq \ell(x_n, y_n) \end{aligned}$$

101

Towards probabilistic robustness

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

subject to

$$\text{Loss}(f_{\theta}(x_n + \delta_0), y_n) \leq \ell(x_n, y_n)$$

$$\text{Loss}(f_{\theta}(x_n + \delta_1), y_n) \leq \ell(x_n, y_n)$$

$$\text{Loss}(f_{\theta}(x_n + \delta_{\sqrt{2}}), y_n) \leq \ell(x_n, y_n)$$

$$\text{Loss}(f_{\theta}(x_n + \delta_{\pi}), y_n) \leq \ell(x_n, y_n)$$

$$\text{Loss}(f_{\theta}(x_n + \delta_{\pi/2}), y_n) \leq \ell(x_n, y_n)$$

$$\text{Loss}(f_{\theta}(x_n + \delta_{3\pi/4}), y_n) \leq \ell(x_n, y_n)$$

$$\text{Loss}(f_{\theta}(x_n + \delta_{5\pi/4}), y_n) \leq \ell(x_n, y_n)$$

$$\text{Loss}(f_{\theta}(x_n + \delta_{7\pi/4}), y_n) \leq \ell(x_n, y_n)$$

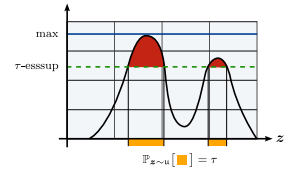
$$\text{Loss}(f_{\theta}(x_n + \delta_{3\pi/2}), y_n) \leq \ell(x_n, y_n)$$

101

Probabilistic robustness

- Probabilistic robustness

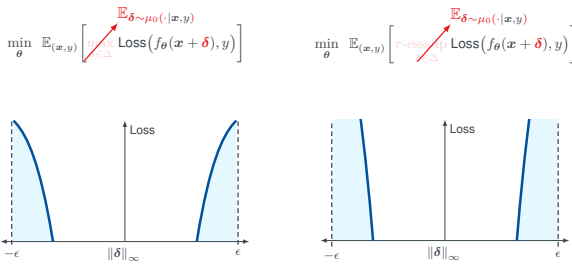
$$\min_{\theta} \mathbb{E}_{(x,y)} \left[\tau\text{-esssup}_{\delta \in \Delta} \text{Loss}(f_{\theta}(x + \delta), y) \right]$$
 - $\tau = 1/2$: classical learning (for symmetric m)
 - $\tau = 0$: adversarial robustness (ess sup)



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

102

Probabilistic robustness



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

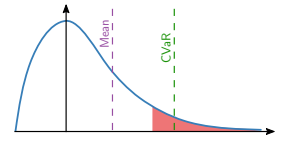
103

Probabilistic robustness and Risk

- Conditional value at risk:

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

$$= \inf_{\alpha \in \mathbb{R}} \alpha + \frac{\mathbb{E}_z [f(z) - \alpha]_+}{1 - \rho}$$
 - $\text{CVaR}_0(f) = \mathbb{E}_z [f(z)]$
 - $\text{CVaR}_1(f) = \text{ess sup}_z f(z)$



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

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

104

Probabilistically robust learning

- for $n = 1, \dots, N$:
 - $\alpha_0 = 0$
 - for $t = 1, \dots, T$:
 - $\delta_t \sim \text{Random}(\Delta)$
 - $\alpha \leftarrow \alpha - \frac{\eta}{T} \left[\tau - \mathbb{I} [\text{Loss}(f_{\theta}(x_n + \delta_t), y_n) \geq \alpha] \right]$
 - end
 - $\theta \leftarrow \theta - \eta \nabla_{\theta} \left[\underbrace{\text{Loss}(f_{\theta}(x_n + \delta_T), y_n) - \alpha}_+ \right]$
 - end
- SGD (CVaR)
- SGD (θ)

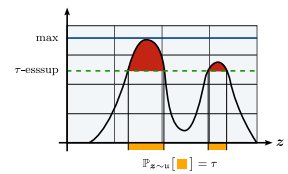
[Robey, Chamon, Pappas, Hassani, ICML'22 (spotlight)]

105

Probabilistic robustness

- Probabilistic robustness

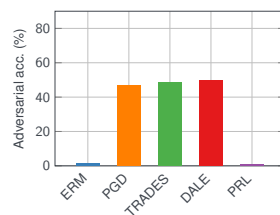
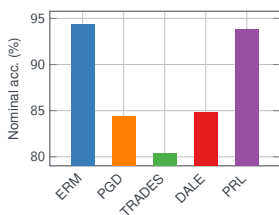
$$\min_{\theta} \mathbb{E}_{(x,y)} \left[\tau\text{-esssup}_{\delta \in \Delta} \text{Loss}(f_{\theta}(x + \delta), y) \right]$$
 - $\tau = 1/2$: classical learning (for symmetric m)
 - $\tau = 0$: adversarial robustness (ess sup)
- Potentially better sample complexity
 [Robey et al., ICML'22 (spotlight)]
 [Raman et al., NeurIPS ML Safety Workshop'22]
- Better performance trade-off
 [Robey et al., ICML'22 (spotlight)]



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

106

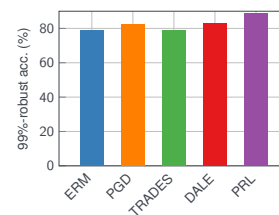
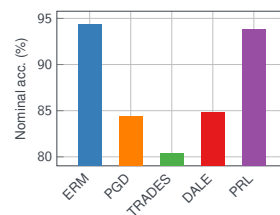
Probabilistically robust learning



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

107

Probabilistically robust learning



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

107

Summary

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

108

Summary

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

108

Summary

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

108

Summary

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

108

Agenda

- I. Constrained supervised learning
 - Constrained learning theory
 - Resilient constrained learning
 - Robust learning

Break (30 min)

- II. Constrained reinforcement learning
 - Constrained RL duality
 - Constrained RL algorithms

Q&A and discussions



<https://luizchamon.com/eusipco>

109