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

**L4DC tutorial
July 15, 2024**

**supervised and reinforcement learning under requirements**

---

## Agenda

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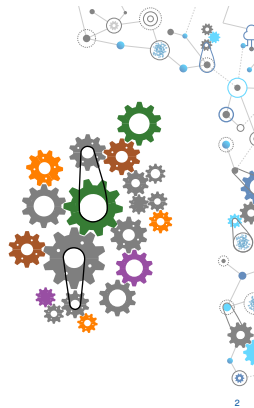- Resilient constrained learning
- Robust learning

Break (30 min)

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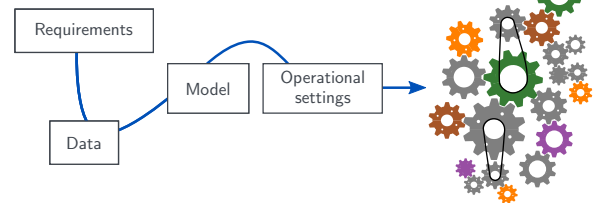- Constrained RL algorithms

https://luizchamon.com/l4dc
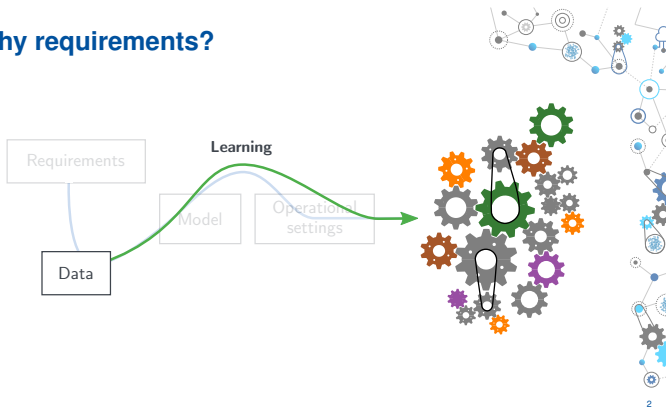
---

## Why requirements?

---

## Why requirements?



Requirements — Data — Model — Operational settings

---

## Why requirements?



Learning

Data

---

## Why requirements?



Requirements

Learning

Data

---

## What is a requirements?

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



Constraint space

Robustness / Safety

---

## 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)
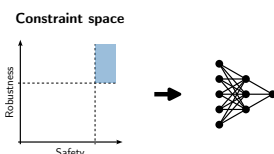  - *Objective space*: things the system achieves



Objective space

Cost / Accuracy

## What is a requirements?

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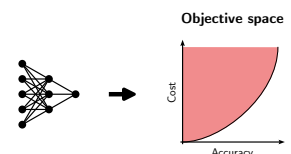  - *Constraint space*: things we decide

- **Goals are "should" statements:** express recommendations (once "shall" statements are satisfied)
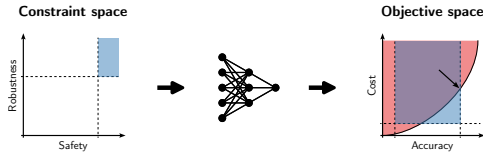  - *Objective space*: things the system achieves



Constraint space — Robustness vs Safety

Objective space — Cost vs Accuracy

[NASA. "Systems engineering handbook," 2019]

---

## What is (un)constrained learning?

$$P_U^\star = \min_\theta \quad \mathbb{E}_{(x,y)\sim\mathfrak{D}}\Big[\ell\big(f_\theta(x),y\big)\Big]$$

- $\ell, g$ are bounded, Lipschitz continuous (possibly non-convex) functions

- $f_\theta$ is a (possibly nonlinear) parametrization [e.g., logistic classifier, (G)(C)NN]

- $\mathfrak{D}, \mathfrak{A}, \mathfrak{P}$ unknown

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

---

## What is (un)constrained learning?

$$P^\star = \min_\theta \quad \mathbb{E}_{(x,y)\sim\mathfrak{D}}\Big[\ell\big(f_\theta(x),y\big)\Big]$$
$$\text{subject to} \quad \mathbb{E}_{(x,y)\sim\mathfrak{A}}\Big[g\big(f_\theta(x),y\big)\Big] \le c$$
$$h\big(f_\theta(x),y\big) \le u, \quad \mathfrak{P}\text{-a.e.}$$

- $\ell, g$ are bounded, Lipschitz continuous (possibly non-convex) functions

- $f_\theta$ is a (possibly nonlinear) parametrization [e.g., logistic classifier, (G)(C)NN]

- $\mathfrak{D}, \mathfrak{A}, \mathfrak{P}$ unknown

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

---

## What about penalties?

$$P^\star = \min_\theta \quad \mathbb{E}_{(x,y)\sim\mathfrak{D}}\Big[\ell\big(f_\theta(x),y\big)\Big]$$
$$\text{subject to} \quad \mathbb{E}_{(x,y)\sim\mathfrak{A}}\Big[g\big(f_\theta(x),y\big)\Big] \le c$$
$$h\big(f_\theta(x),y\big) \le u, \quad \mathfrak{P}\text{-a.e.}$$

$$\Updownarrow$$

$$\min_\theta \quad \mathbb{E}_{(x,y)\sim\mathfrak{D}}\Big[\ell\big(f_\theta(x),y\big)\Big] + \lambda\,\mathbb{E}_{(x,y)\sim\mathfrak{A}}\Big[g\big(f_\theta(x),y\big)\Big] + \mathbb{E}_{(x,y)\sim\mathfrak{P}}\Big[\mu(x,y)h\big(f_\theta(x),y\big)\Big]$$

---

## What about penalties?

$$P^\star = \min_\theta \quad \mathbb{E}_{(x,y)\sim\mathfrak{D}}\Big[\ell\big(f_\theta(x),y\big)\Big]$$
$$\text{subject to} \quad \mathbb{E}_{(x,y)\sim\mathfrak{A}}\Big[g\big(f_\theta(x),y\big)\Big] \le c$$
$$h\big(f_\theta(x),y\big) \le u, \quad \mathfrak{P}\text{-a.e.}$$

**NON-CONVEX**

$$\min_\theta \quad \mathbb{E}_{(x,y)\sim\mathfrak{D}}\Big[\ell\big(f_\theta(x),y\big)\Big] + \lambda\,\mathbb{E}_{(x,y)\sim\mathfrak{A}}\Big[g\big(f_\theta(x),y\big)\Big] + \mathbb{E}_{(x,y)\sim\mathfrak{P}}\Big[\mu(x,y)h\big(f_\theta(x),y\big)\Big]$$

- ❌ There may not exist $(\lambda, \mu)$ such that the penalized solution is optimal *and* feasible

- ❌ Even if such $(\lambda, \mu)$ exist, they are not easy to find (hyperparameter search, cross-validation...)

- ✅ Constrained learning yields better guarantees, better performance, better trade-offs...
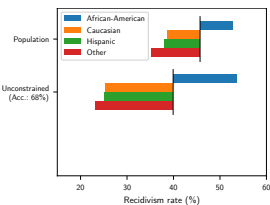
---

## Applications

- **Fairness**
  (e.g., [Goh et al., NeurIPS'16; Kearns et al., ICML'18; Cotter et al., JMLR'19; Chamon et al., IEEE TIT'23])

- **Federated learning**
  (e.g., [Shen et al., ICLR'22; Hounie et al., NeurIPS'23])

- **Adversarially robust learning**
  (e.g., [Chamon et al., NeurIPS'20; Robey et al., NeurIPS'21; Chamon et al., IEEE TIT'23])

- **Safe learning**
  (e.g., [Paternain et al., IEEE TAC'23])

- **Wireless resource allocation**
  (e.g., [Eisen et al., IEEE TSP'19; NaderiAlizadeh et al., IEEE TSP'22; Chowdhury et al., Asilomar'23])
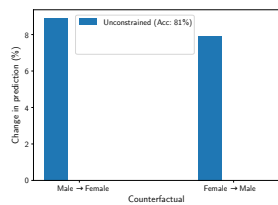
- ...

---

## 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.

---

## Fairness: "Equality" of odds

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

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

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

## Fairness: "Equality" of odds

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

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

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

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

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

9

## Fairness: "Equality" of odds

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

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

$$\text{subject to} \quad \frac{1}{N}\sum_{n=1}^{N} \mathbb{I}\left[f_{\theta}(\boldsymbol{x}_n) = 1 \mid \text{Race}\right] \leq \frac{1}{N}\sum_{n=1}^{N} \mathbb{I}\left[f_{\theta}(\boldsymbol{x}_n) = 1\right] + c,$$

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

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

9

## Counterfactual fairness

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

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

$$\text{subject to} \quad \text{Change in prediction } (\rho\boldsymbol{x}) \leq c \quad \text{a.e.}$$

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

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

10

## Counterfactual fairness

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

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

$$\text{subject to} \quad \text{D}_{\text{KL}}\big(f_{\theta}(\boldsymbol{x}_n)\|f_{\theta}(\rho\boldsymbol{x}_n)\big) \leq c, \quad \text{for all } n$$

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

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

10

## 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])

- ...

11

## Federated learning

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



**Clients**
$k_1 \; k_2 \; k_3 \; k_4 \; k_5 \; k_6 \; k_7 \; k_8$
$k_9 \; k_{10} \; k_{11} \; k_{12} \; k_{13} \; k_{14} \; k_{15} \; k_{16}$

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

- $k$-th client loss: $\text{Loss}_k(f_{\theta}) = \frac{1}{N_k}\sum_{n_k=1}^{N_k} \text{Loss}\big(f_{\theta}(\boldsymbol{x}_{n_k}), y_{n_k}\big)$
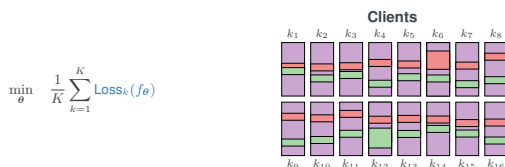
[Shen et al., ICRL'22]

12

## Federated learning

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



**Clients**
$k_1 \; k_2 \; k_3 \; k_4 \; k_5 \; k_6 \; k_7 \; k_8$
$k_9 \; k_{10} \; k_{11} \; k_{12} \; k_{13} \; k_{14} \; k_{15} \; k_{16}$

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

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

[Shen et al., ICRL'22]

12

## Federated learning

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



**Clients**
$k_1 \; k_2 \; k_3 \; k_4 \; k_5 \; k_6 \; k_7 \; k_8$
$k_9 \; k_{10} \; k_{11} \; k_{12} \; k_{13} \; k_{14} \; k_{15} \; k_{16}$

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

- $k$-th client loss: $\text{Loss}_k(f_{\theta}) = \frac{1}{N_k}\sum_{n_k=1}^{N_k} \text{Loss}\big(f_{\theta}(\boldsymbol{x}_{n_k}), y_{n_k}\big)$
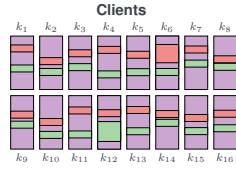
[Shen et al., ICRL'22]

12

## Federated learning

**Problem**
Learn a common model using data using data distributed among $K$ clients

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

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

**Clients**



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

---

## Federated learning

**Problem**
Learn a common model using data using data distributed among $K$ clients

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

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

**Clients**



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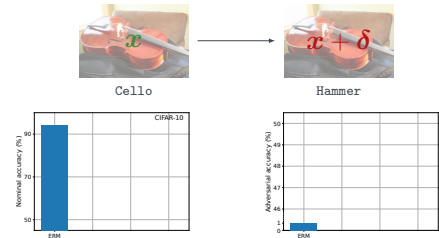
---

## 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])

- ...

---

## Robustness

**Problem**
Learn a classifier that is robust to input perturbations



Cello          Hammer

---

## Robustness

**Problem**
Learn a classifier that is robust to input perturbations



Cello          Hammer

$$\min_{\theta} \quad \text{Nominal accuracy}$$
$$\text{subject to} \quad \text{Robustness} \leq c$$

---

## Robustness

**Problem**
Learn a classifier that is robust to input perturbations



Cello          Hammer

$$\min_{\theta} \quad \frac{1}{N} \sum_{n=1}^{N} \text{Loss}\big(f_{\theta}(\boldsymbol{x}_n), y_n\big)$$
$$\text{subject to} \quad \frac{1}{N} \sum_{n=1}^{N} \text{Robustness} \leq c$$

---

## Robustness

**Problem**
Learn a classifier that is robust to input perturbations



Cello          Hammer

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

---

## (Manifold) smoothness

**Problem**
Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories

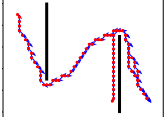# (Manifold) smoothness

**Problem**
Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories

$$\min_{\theta} \quad \text{Imitation error}$$
$$\text{subject to} \quad \text{Smoothness in free space} \leq L$$
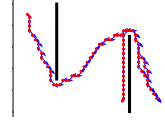
# (Manifold) smoothness

**Problem**
Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories

$$\min_{\theta} \quad \frac{1}{N}\sum_{n=1}^{N} \text{Loss}\big(f_{\theta}(\boldsymbol{x}_n), \boldsymbol{u}_n\big)$$
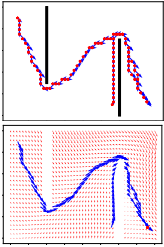$$\text{subject to} \quad \text{Smoothness in free space} \leq L$$

# (Manifold) smoothness

**Problem**
Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories

$$\min_{\theta} \quad \frac{1}{N}\sum_{n=1}^{N} \text{Loss}\big(f_{\theta}(\boldsymbol{x}_n), \boldsymbol{u}_n\big)$$
$$\text{subject to} \quad \max_{\boldsymbol{x}\in\mathcal{M}} \|\nabla_{\mathcal{M}} f_{\theta}(\boldsymbol{x})\|^2 \leq L$$
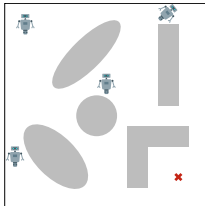
# Applications

- 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])

- …

# Safety

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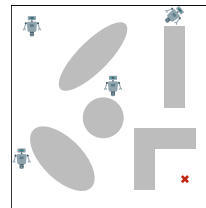
# Safety

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

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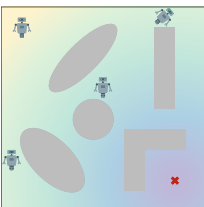$$\text{for } i = 1, 2, \dots$$

# Safety

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

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

# Safety

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

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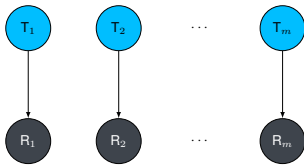$$\text{for } i = 1, 2, \dots$$

## 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])
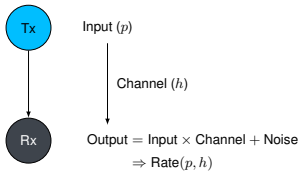
- …

## Wireless resource allocation

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

## Wireless resource allocation

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



Input ($p$)

Channel ($h$)

Output = Input × Channel + Noise
$\Rightarrow$ Rate$(p, h)$

## Wireless resource allocation

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

## Wireless resource allocation

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

## Wireless resource allocation

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

## Wireless resource allocation

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



Input ($p$)  ← Action

Channel ($h$)  ← State

Output = Input × Channel + Noise
$\Rightarrow$ Rate$(p, h)$  ← Reward

## Wireless resource allocation

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



$$\min_{p} \quad \text{Total transmit power}$$
$$\text{s. to} \quad \text{Rate } T_i \to R_i \geq c_i$$

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

## Wireless resource allocation

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

$$\min_{p} \quad \sum_{i=1}^{m} \mathbb{E}\left[\frac{1}{T}\sum_{t=0}^{T-1} p_i(h_t)\right]$$
$$\text{s. to} \quad \text{Rate } T_i \to R_i \geq c_i$$

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

---

## Wireless resource allocation

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

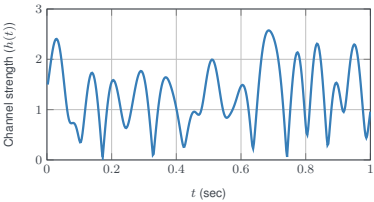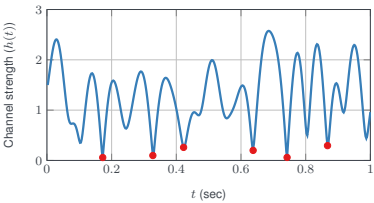$$\min_{p} \quad \sum_{i=1}^{m} \mathbb{E}\left[\frac{1}{T}\sum_{t=0}^{T-1} p_i(h_t)\right]$$
$$\text{s. to} \quad \mathbb{E}\left[\frac{1}{T}\sum_{t=0}^{T-1} \text{Rate}_i\big(\boldsymbol{p}(h_t), h_t\big)\right] \geq c_i$$

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

---

## Wireless resource allocation

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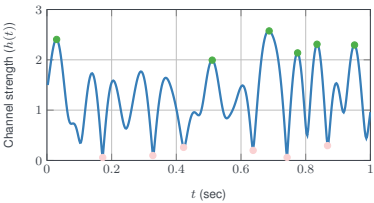Allocate the least transmit power to $m$ device pairs to achieve a communication rate

$$\min_{p} \quad \text{Total probability of depleting battery}$$
$$\text{s. to} \quad \mathbb{E}\left[\frac{1}{T}\sum_{t=0}^{T-1} \text{Rate}_i\big(\boldsymbol{p}(h_t), h_t\big)\right] \geq c_i$$

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

---

## Wireless resource allocation

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

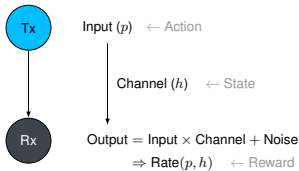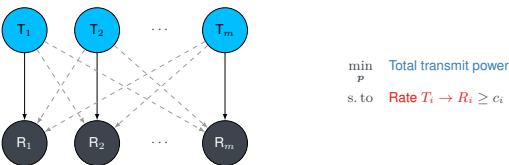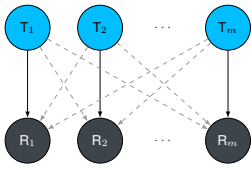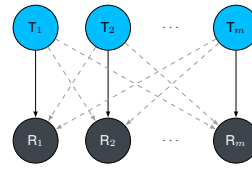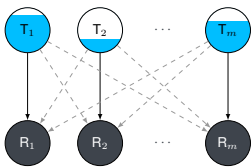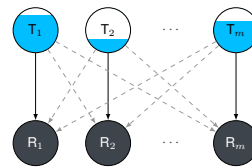$$\min_{p} \quad \sum_{i=1}^{m} \Pr\left[\bigcap_{t=0}^{T-1} \{b_{i,t}=0\}\right]$$
$$\text{s. to} \quad \mathbb{E}\left[\frac{1}{T}\sum_{t=0}^{T-1} \text{Rate}_i\big(\boldsymbol{p}(h_t), h_t\big)\right] \geq c_i$$

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

- Continual learning (e.g., [Peng et al., ICML'23])

- Active learning (e.g., [Elenter et al., NeurIPS'22])

- Data augmentation (e.g., [Hounie et al., ICML'23])

- Semi-supervised learning (e.g., [Cerviño et al., ICML'23])

- Minimum norm interpolation, SVM...

---

## Constrained supervised learning

---

## What is (un)constrained learning?

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

- $\ell, g$ are bounded, Lipschitz continuous (possibly non-convex) functions

- $f_{\boldsymbol{\theta}}$ is a (possibly nonlinear) parametrization [e.g., logistic classifier, (G)(C)NN]

- $(\boldsymbol{x}_n, y_n) \sim \mathfrak{D}, \ (\boldsymbol{x}_m, y_m) \sim \mathfrak{A}, \ (\boldsymbol{x}_r, y_r) \sim \mathfrak{P}$ (i.i.d.)

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

---

## What is (un)constrained learning?

$$P^\star = \min_{\boldsymbol{\theta}} \quad \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\left[\ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\right]$$
$$\text{subject to} \quad \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{A}}\left[g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\right] \leq c$$
$$h\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big) \leq u, \quad \mathfrak{P}\text{-a.e.}$$

- $\ell, g$ are bounded, Lipschitz continuous (possibly non-convex) functions

- $f_{\boldsymbol{\theta}}$ is a (possibly nonlinear) parametrization [e.g., logistic classifier, (G)(C)NN]

- $\mathfrak{D}, \mathfrak{A}, \mathfrak{P}$ unknown

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

## Constrained learning challenges

$$\hat{P}^\star = \min_{\boldsymbol{\theta}} \ \frac{1}{N}\sum_{n=1}^{N} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) \qquad \overset{?}{\longrightarrow} \qquad P^\star = \min_{\boldsymbol{\theta}} \ \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\Big[\ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big]$$

$$\text{subject to} \ \ \frac{1}{N}\sum_{m=1}^{N} g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big) \le c \qquad\qquad \text{subject to} \ \ \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{A}}\Big[g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big] \le c$$

$$h\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_r, y_r\big) \le u \qquad\qquad\qquad h\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big) \le u \ \text{a.e.}$$

**Challenges**

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

---

## Constrained learning challenges

$$\hat{P}^\star = \min_{\boldsymbol{\theta}} \ \frac{1}{N}\sum_{n=1}^{N} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) \qquad \overset{?}{\longrightarrow} \qquad P^\star = \min_{\boldsymbol{\theta}} \ \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\Big[\ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big]$$

$$\text{subject to} \ \ \frac{1}{N}\sum_{m=1}^{N} g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big) \le c \qquad\qquad \text{subject to} \ \ \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{A}}\Big[g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big] \le c$$

$$h\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_r, y_r\big) \le u \qquad\qquad\qquad h\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big) \le u \ \text{a.e.}$$

**Challenges**

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

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

---

## Constrained learning challenges

$$\hat{P}^\star = \min_{\boldsymbol{\theta}} \ \frac{1}{N}\sum_{n=1}^{N} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) \qquad \overset{?}{\longrightarrow} \qquad P^\star = \min_{\boldsymbol{\theta}} \ \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\Big[\ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big]$$

$$\text{subject to} \ \ \frac{1}{N}\sum_{m=1}^{N} g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big) \le c \qquad\qquad \text{subject to} \ \ \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{A}}\Big[g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big] \le c$$

$$h\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_r, y_r\big) \le u \qquad\qquad\qquad h\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big) \le u \ \text{a.e.}$$

**Challenges**

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

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

---

## Agenda

Constrained learning theory

Constrained learning algorithms

Resilient constrained learning

---

## Constrained learning challenges

$$\hat{P}^\star = \min_{\boldsymbol{\theta}} \ \frac{1}{N}\sum_{n=1}^{N} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) \qquad \overset{?}{\longrightarrow} \qquad P^\star = \min_{\boldsymbol{\theta}} \ \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\Big[\ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big]$$

$$\text{subject to} \ \ \frac{1}{N}\sum_{m=1}^{N} g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big) \le c \qquad\qquad \text{subject to} \ \ \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{A}}\Big[g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big] \le c$$

$$h\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_r, y_r\big) \le u \qquad\qquad\qquad h\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big) \le u \ \text{a.e.}$$
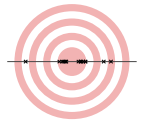
**Challenges**

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

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

---

## What classical learning theory says?

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

✅ $f_{\boldsymbol{\theta}}$ is *probably approximately correct (PAC)* learnable

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

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

---

## What classical learning theory says?

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

✅ $f_{\boldsymbol{\theta}}$ is *probably approximately correct (PAC)* learnable

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

❌ Requirements?

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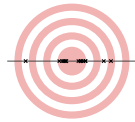
---

## What's in a solution?

**Definition (PAC learnability)**

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

- near-optimal

$$P^\star - \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\Big[\ell\big(f_{\boldsymbol{\theta}^\dagger}(\boldsymbol{x}), y\big)\Big] \le \epsilon$$

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

## What's in a solution?

**Definition (PACC learnability)**

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

- near-optimal

$$\left| P^\star - \mathbb{E}_{(x,y)\sim\mathfrak{D}}\left[ \ell\big(f_{\theta\dagger}(x), y\big)\right] \right| \le \epsilon$$

- approximately feasible

$$\mathbb{E}_{(x,y)\sim\mathfrak{A}}\left[ g\big(f_{\theta\dagger}(x), y\big)\right] \le c + \epsilon$$

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

---

## When is constrained learning possible?

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

$$\text{subject to} \quad \frac{1}{N}\sum_{m=1}^N g\big(f_\theta(x_m), y_m\big) \le c \qquad \text{subject to} \quad \mathbb{E}_{(x,y)\sim\mathfrak{A}}\left[g\big(f_\theta(x), y\big)\right] \le c$$

**Proposition**

$$f_\theta \text{ is PAC learnable} \;\not\Rightarrow\; f_\theta \text{ is PACC learnable}$$

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

---

## ECRM is not a PACC learner

**Counter-example**

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

$$\text{subject to} \quad \theta_2 \mathbb{E}_\tau[\tau] \le \theta_1 - 1$$
$$-\theta_1 \mathbb{E}_\tau[\tau] \le \theta_2 - 1$$

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

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

---

## ECRM is not a PACC learner

**Counter-example**

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

$$\text{subject to} \quad \theta_2 \mathbb{E}_\tau[\tau] \le \theta_1 - 1 \Rightarrow \theta_1 \ge 1$$
$$-\theta_1 \mathbb{E}_\tau[\tau] \le \theta_2 - 1 \Rightarrow \theta_2 \le 1$$

$$J(\theta) = \begin{cases} 1/16, & \theta = [1/2, 1/2] \\ 1/8, & \theta = [1, 1] \\ 1/4, & \theta = [1, 0] \end{cases}$$
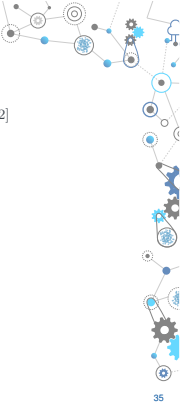
- $\tau \sim \text{Uniform}\big(-1/2, 1/2\big)$

---

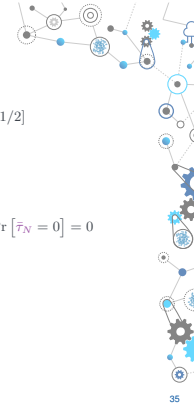## ECRM is not a PACC learner

**Counter-example**

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

$$\text{subject to} \quad \theta_2 \mathbb{E}_\tau[\tau] \le \theta_1 - 1 \Rightarrow \theta_1 \ge 1$$
$$-\theta_1 \mathbb{E}_\tau[\tau] \le \theta_2 - 1 \Rightarrow \theta_2 \le 1$$

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

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

$$\text{subject to} \quad \theta_2 \bar{\tau}_N \le \theta_1 - 1$$
$$-\theta_1 \bar{\tau}_N \le 1 - \theta_2$$

$$\Pr\left[|\hat{P}^\star - P^\star| \le 1/32\right] = \Pr\left[\bar{\tau}_N = 0\right] = 0$$

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

---

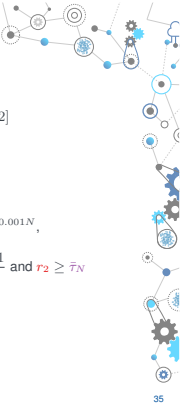## ECRM is not a PACC learner

**Counter-example**

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

$$\text{subject to} \quad \theta_2 \mathbb{E}_\tau[\tau] \le \theta_1 - 1 \Rightarrow \theta_1 \ge 1$$
$$-\theta_1 \mathbb{E}_\tau[\tau] \le \theta_2 - 1 \Rightarrow \theta_2 \le 1$$

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

$$\hat{P}_r^\star = \min_{\theta\in\Theta} \quad J(\theta)$$

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

$$\Pr\left[|\hat{P}_r^\star - P^\star| \le 1/32\right] \le 4e^{-0.001N},$$
$$\text{unless } \bar{\tau}_N \le r_1 < \frac{\bar{\tau}_N + 1}{2} \text{ and } r_2 \ge \bar{\tau}_N$$

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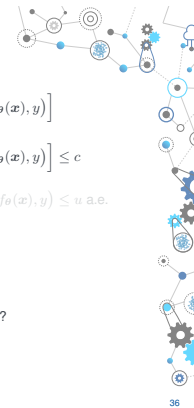
---

## Constrained learning challenges

$$\hat{P}^\star = \min_\theta \frac{1}{N}\sum_{n=1}^N \ell\big(f_\theta(x_n), y_n\big) \qquad P^\star = \min_\theta \mathbb{E}_{(x,y)\sim\mathfrak{D}}\left[\ell\big(f_\theta(x), y\big)\right]$$

$$\text{subject to} \quad \frac{1}{N}\sum_{m=1}^N g\big(f_\theta(x_m), y_m\big) \le c \quad \xrightarrow{\text{PAC}} \quad \text{subject to} \quad \mathbb{E}_{(x,y)\sim\mathfrak{A}}\left[g\big(f_\theta(x), y\big)\right] \le c$$

$$h\big(f_\theta(x_r, y_r\big) \le u \qquad\qquad h\big(f_\theta(x), y\big) \le u \text{ a.e.}$$

**Challenges**

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

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

---

## Constrained learning challenges

$$\hat{P}^\star = \min_\theta \frac{1}{N}\sum_{n=1}^N \ell\big(f_\theta(x_n), y_n\big) \qquad P^\star = \min_\theta \mathbb{E}_{(x,y)\sim\mathfrak{D}}\left[\ell\big(f_\theta(x), y\big)\right]$$

$$\text{subject to} \quad \frac{1}{N}\sum_{m=1}^N g\big(f_\theta(x_m), y_m\big) \le c \quad \xrightarrow{\text{PAC}} \quad \text{subject to} \quad \mathbb{E}_{(x,y)\sim\mathfrak{A}}\left[g\big(f_\theta(x), y\big)\right] \le c$$

$$h\big(f_\theta(x_r, y_r\big) \le u$$

**Challenges**

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

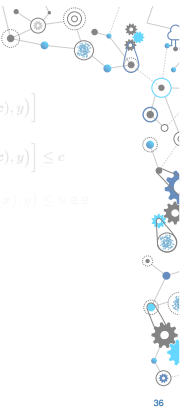2) *Computational*: can we solve the constrained empirical problem?

## Duality

PRIMAL

DUAL

---

## Duality

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

---

## Duality

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

$$\hat{D}^\star = \max_{\lambda \ge 0} \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^{N} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) + \lambda\Big[\frac{1}{N} \sum_{m=1}^{N} g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big) - c\Big]$$

---

## Duality

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

$$\hat{D}^\star = \max_{\lambda \ge 0} \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^{N} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) + \lambda\Big[\frac{1}{N} \sum_{m=1}^{N} g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big) - c\Big]$$

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

---

## Duality

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

$$\hat{D}^\star = \max_{\lambda \ge 0} \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^{N} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) + \lambda\Big[\frac{1}{N} \sum_{m=1}^{N} g\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big) - c\Big]$$

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

---

## An alternative path

$$\hat{P}^\star = \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^{N} \ell(f_{\boldsymbol{\theta}}, z_n)$$
$$\text{s. to } \frac{1}{N} \sum_{n=1}^{N} g(f_{\boldsymbol{\theta}}, z_n) \le c$$

$$\hat{D}^\star = \max_{\lambda \ge 0} \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{N} \sum_{n=1}^{N} \ell(f_{\boldsymbol{\theta}}, z_n) + \lambda\Big(\frac{1}{N} \sum_{n=1}^{N} g(f_{\boldsymbol{\theta}}, z_n) - c\Big)$$

PAC

$$P^\star = \min_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_{\boldsymbol{z}}\big[\ell(f_{\boldsymbol{\theta}}, \boldsymbol{z})\big]$$
$$\text{s. to } \mathbb{E}_{\boldsymbol{z}}\big[g(f_{\boldsymbol{\theta}}, \boldsymbol{z})\big] \le c$$

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

---

## An alternative path

$$\hat{P}^\star = \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^{N} \ell(f_{\boldsymbol{\theta}}, z_n)$$
$$\text{s. to } \frac{1}{N} \sum_{n=1}^{N} g(f_{\boldsymbol{\theta}}, z_n) \le c$$

$$\hat{D}^\star = \max_{\lambda \ge 0} \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{N} \sum_{n=1}^{N} \ell(f_{\boldsymbol{\theta}}, z_n) + \lambda\Big(\frac{1}{N} \sum_{n=1}^{N} g(f_{\boldsymbol{\theta}}, z_n) - c\Big)$$

PAC

$$P^\star = \min_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_{\boldsymbol{z}}\big[\ell(f_{\boldsymbol{\theta}}, \boldsymbol{z})\big]$$
$$\text{s. to } \mathbb{E}_{\boldsymbol{z}}\big[g(f_{\boldsymbol{\theta}}, \boldsymbol{z})\big] \le c$$

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

$$\tilde{P}^\star = \min_{\phi \in \mathcal{H}} \mathbb{E}_{\boldsymbol{z}}\big[\ell(\phi, \boldsymbol{z})\big]$$
$$\text{s. to } \mathbb{E}_{\boldsymbol{z}}\big[g(\phi, \boldsymbol{z})\big] \le c$$

?

$$\tilde{D}^\star = \max_{\lambda \ge 0} \min_{\phi \in \mathcal{H}} \mathbb{E}_{\boldsymbol{z}}\big[\ell(\phi, \boldsymbol{z})\big] + \lambda\big(\mathbb{E}_{\boldsymbol{z}}\big[g(\phi, \boldsymbol{z})\big] - c\big)$$

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

---

## Non-convex variational duality

Convex optimization:  Primal ⟷ Dual

Non-convex, finite dimensional optimization:  Primal ⟷ Dual

## Non-convex variational duality

Convex optimization:     Primal ⟷ Dual

Non-convex, finite dimensional optimization:     Primal ⟷ Dual

**Non-convex, infinite** dimensional optimization:     Primal ⟷ Dual

## Sparse logistic regression

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

$$\text{s. to } \|\boldsymbol{\theta}\|_0 = \sum_{t=1}^{P} \mathbb{I}\left[\boldsymbol{\theta}_t \neq 0\right] \leq k$$

**Discrete, non-convex**

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

## Sparse logistic regression

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

$$\text{s. to } \|\boldsymbol{\theta}\|_0 = \sum_{t=1}^{P} \mathbb{I}\left[\boldsymbol{\theta}_t \neq 0\right] \leq k$$

**Discrete, non-convex**

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

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

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

**Continuous, non-convex**

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

## Sparse logistic regression

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

$$\text{s. to } \|\boldsymbol{\theta}\|_0 = \sum_{t=1}^{P} \mathbb{I}\left[\boldsymbol{\theta}_t \neq 0\right] \leq k$$

**Discrete**, non-convex

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

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

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

**Continuous**, non-convex

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

## An alternative path

PRIMAL     DUAL

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

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

PAC

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

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

## An alternative path

PRIMAL     DUAL

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

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

PAC

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

$$\overset{\epsilon_0}{\longleftrightarrow} D^\star = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \mathbb{E}_z\left[\ell(f_\theta, z)\right] + \lambda\left(\mathbb{E}_z\left[g(f_\theta, z)\right] - c\right)$$

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

## An alternative path

PRIMAL     DUAL

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

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

PAC

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

$$\overset{\epsilon_0}{\longleftrightarrow} D^\star = \max_{\lambda \geq 0} \min_{\theta \in \Theta} \mathbb{E}_z\left[\ell(f_\theta, z)\right] + \lambda\left(\mathbb{E}_z\left[g(f_\theta, z)\right] - c\right)$$

$$O(\epsilon)$$

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

## Dual (near-)PACC learning

**Theorem**

Let $f$ be $\nu$-universal, i.e., for each $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2$, and $\gamma \in [0,1]$ there exists $\boldsymbol{\theta}$ such that

$$\mathbb{E}\left[\left|\gamma f_{\boldsymbol{\theta}_1}(\boldsymbol{x}) + (1-\gamma) f_{\boldsymbol{\theta}_2}(\boldsymbol{x}) - f_{\boldsymbol{\theta}}(\boldsymbol{x})\right|\right] \leq \nu$$

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

## Dual (near-)PACC learning

**Theorem**
Let $f$ be $\nu$-universal, i.e., for each $\theta_1$, $\theta_2$, and $\gamma \in [0,1]$ there exists $\theta$ 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 $\hat{D}^\star$ is a (near-)PACC learner, i.e., there exists a solution $\theta^\dagger$ that, with probability $1 - \delta$,

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

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

(mild conditions apply)

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

---

## Dual (near-)PACC learning

**Theorem**
Let $f$ be $\nu$-universal, i.e., for each $\theta_1$, $\theta_2$, and $\gamma \in [0,1]$ there exists $\theta$ 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 $\hat{D}^\star$ is a (near-)PACC learner, i.e., there exists a solution $\theta^\dagger$ that, with probability $1 - \delta$,

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

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

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

(mild conditions apply)

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

---

## Dual (near-)PACC learning

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

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

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

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

**Sources of error**
$\qquad$ parametrization richness ($\nu$) $\qquad$ sample size ($N$) $\qquad$ requirements difficulty ($\lambda^\star$)

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

---

## Dual (near-)PACC learning

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

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

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

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

**Sources of error**
$\qquad$ parametrization richness ($\nu$) $\qquad$ sample size ($N$) $\qquad$ requirements difficulty ($\lambda^\star$)

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

---

## Dual (near-)PACC learning

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

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

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

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

**Sources of error**
$\qquad$ parametrization richness ($\nu$) $\qquad$ sample size ($N$) $\qquad$ requirements difficulty ($\lambda^\star$)

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

---

## Dual (near-)PACC learning

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

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

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

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

**Sources of error**
$\qquad$ parametrization richness ($\nu$) $\qquad$ sample size ($N$) $\qquad$ requirements difficulty ($\lambda^\star$)

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

---

## Dual learning trade-offs

- Unconstrained learning
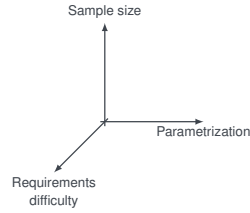
  parametrization $\times$ sample size



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

## Dual learning trade-offs

- Unconstrained learning

  parametrization $\times$ sample size

- Constrained learning

  parametrization $\times$ sample size $\times$ requirements


Sample size / Parametrization / Requirements difficulty

---

## When is constrained learning possible?

Corollary

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

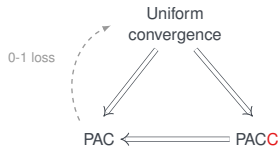Constrained learning is **essentially as hard as** unconstrained learning

[mild conditions apply]

---

## When is constrained learning possible?

Corollary


Uniform convergence / 0-1 loss / PAC ⟵ PACC

[mild conditions apply]

---

## Fairness

Problem
Predict whether an individual will recidivate



Problem
Predict whether an individual makes > $50k

---

## Fairness

Problem
Predict whether an individual will recidivate



Problem
Predict whether an individual makes > $50k

---

## Fairness: "Equality" of odds

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

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

$$\text{subject to} \quad \frac{1}{N}\sum_{n=1}^N \mathbb{I}\left[f_\theta(\boldsymbol{x}_n) = 1 \mid \text{Race}\right] \le \frac{1}{N}\sum_{n=1}^N \mathbb{I}\left[f_\theta(\boldsymbol{x}_n) = 1\right] + c,$$

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

---

## Fairness: "Equality" of odds

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

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

$$\text{subject to} \quad \frac{1}{N}\sum_{n=1}^N \mathbb{I}\left[f_\theta(\boldsymbol{x}_n) = 1 \mid \text{Race}\right] \le \frac{1}{N}\sum_{n=1}^N \mathbb{I}\left[f_\theta(\boldsymbol{x}_n) = 1\right] + c,$$

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

---

## Fairness: "Equality" of odds

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

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

$$\text{subject to} \quad \frac{1}{N}\sum_{n=1}^N \sigma\big(f_\theta(\boldsymbol{x}_n) - 0.5\big)\,\mathbb{I}\left[\boldsymbol{x}_n \in \text{Race}\right] \le \frac{1}{N}\sum_{n=1}^N \sigma\big(f_\theta(\boldsymbol{x}_n) - 0.5\big) + c,$$

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

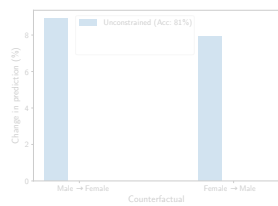# Fairness: "Equality" of odds

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



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

48

---

# Fairness: "Equality" of odds

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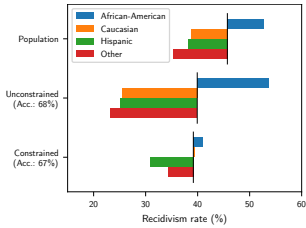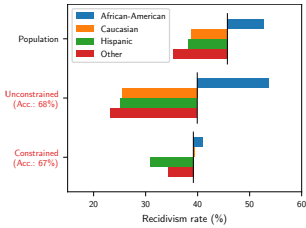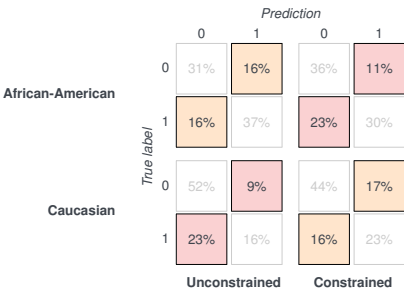Predict whether an individual will recidivate at the same rate across races



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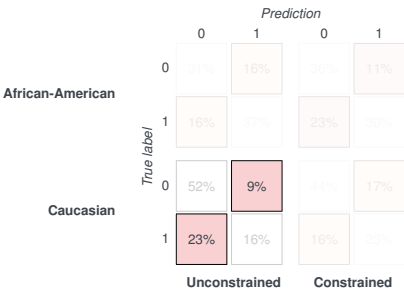[Chamon et al., IEEE TIT'23]

48

---

# Fairness: "Equality" of odds



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

49

---

# Fairness: "Equality" of odds



\* We say "Race" to follow the terminology used during the data collection of the COMPAS dataset.
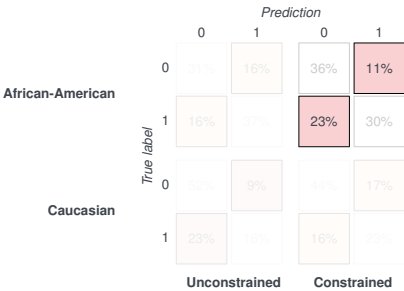[Chamon et al., IEEE TIT'23]
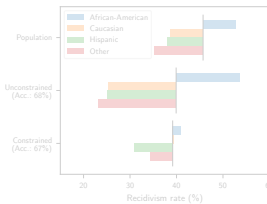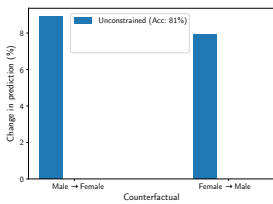
49

---

# Fairness: "Equality" of odds



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

49

---

# Fairness

Problem
Predict whether an individual will recidivate



Problem
Predict whether an individual makes > $50k

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

50

---

# Counterfactual fairness

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

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

$$\text{subject to} \quad \text{D}_{\text{KL}}\big(f_{\theta}(\boldsymbol{x}_n)\|f_{\theta}(\rho\boldsymbol{x}_n)\big) \leq c, \quad \text{for all } n$$

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

[Chamon and Ribeiro, NeurIPS'20]

51

---

# Counterfactual fairness

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

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

$$\text{subject to} \quad \frac{1}{N}\sum_{n=1}^{N} \text{D}_{\text{KL}}\big(f_{\theta}(\boldsymbol{x}_n)\|f_{\theta}(\rho\boldsymbol{x}_n)\big) \leq c, \quad \text{for all } n$$

$$(\rho : \text{Male} \leftrightarrow \text{Female})$$
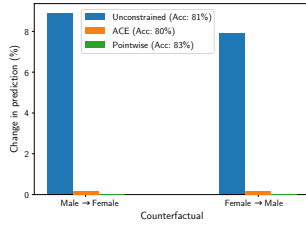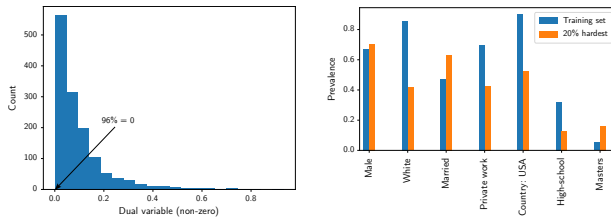
[Chamon and Ribeiro, NeurIPS'20]

51

# Counterfactual fairness

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



*Legend: Unconstrained (Acc: 81%), ACE (Acc: 80%), Pointwise (Acc: 83%)*

---

# Counterfactual fairness

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

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

$$\text{subject to} \quad \text{D}_{\text{KL}}\big(f_{\theta}(\boldsymbol{x}_n)\|f_{\theta}(\rho\boldsymbol{x}_n)\big) \le c, \quad \text{for all } n$$

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

$$\updownarrow$$

$$\max_{\lambda_n \ge 0}\ \min_{\theta}\ \frac{1}{N}\sum_{n=1}^{N}\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n), y_n\big) + \sum_{n=1}^{N}\lambda_n\Big[\text{D}_{\text{KL}}\big(f_{\theta}(\boldsymbol{x}_n)\|f_{\theta}(\rho\boldsymbol{x}_n)\big) - c\Big]$$

---

# Counterfactual fairness

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



*Legend: Training set, 20% hardest*

---

# Agenda

Constrained learning theory

Constrained learning algorithms

Resilient constrained learning

---

# Constrained optimization methods

$$\hat{P}^{\star} = \min_{\theta} \quad \frac{1}{N}\sum_{n=1}^{N}\ell\big(f_{\theta}(\boldsymbol{x}_n), y_n\big)$$

$$\text{subject to} \quad \frac{1}{N}\sum_{m=1}^{N}g\big(f_{\theta}(\boldsymbol{x}_m), y_m\big) \le c$$

$$h\big(f_{\theta}(\boldsymbol{x}_r), y_r\big) \le u$$

---

# Constrained optimization methods

$$\hat{P}^{\star} = \min_{\theta} \quad \frac{1}{N}\sum_{n=1}^{N}\ell\big(f_{\theta}(\boldsymbol{x}_n), y_n\big)$$

$$\text{subject to} \quad \frac{1}{N}\sum_{m=1}^{N}g\big(f_{\theta}(\boldsymbol{x}_m), y_m\big) \le c$$

$$h\big(f_{\theta}(\boldsymbol{x}_r), y_r\big) \le u$$

- Feasible update methods
  e.g., conditional gradients (Frank-Wolfe)

- Interior point methods
  e.g., barriers, projection, polyhedral approx.

---

# Constrained optimization methods

$$\hat{P}^{\star} = \min_{\theta} \quad \frac{1}{N}\sum_{n=1}^{N}\ell\big(f_{\theta}(\boldsymbol{x}_n), y_n\big)$$

$$\text{subject to} \quad \frac{1}{N}\sum_{m=1}^{N}g\big(f_{\theta}(\boldsymbol{x}_m), y_m\big) \le c$$

$$h\big(f_{\theta}(\boldsymbol{x}_r), y_r\big) \le u$$

- 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

---

# Constrained optimization methods

$$\hat{P}^{\star} = \min_{\theta} \quad \frac{1}{N}\sum_{n=1}^{N}\ell\big(f_{\theta}(\boldsymbol{x}_n), y_n\big)$$

$$\text{subject to} \quad \frac{1}{N}\sum_{m=1}^{N}g\big(f_{\theta}(\boldsymbol{x}_m), y_m\big) \le c$$

$$h\big(f_{\theta}(\boldsymbol{x}_r), y_r\big) \le u$$

- 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]

## Dual learning algorithm

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

---

## Dual learning algorithm

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

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

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

---

## Dual learning algorithm

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

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

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

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

---

## Dual learning algorithm

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

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

- Update the dual

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

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

---

## A (near-)PACC learner

**Theorem**

Suppose $\boldsymbol{\theta}^{\dagger}$ is a $\rho$-approximate solution of the regularized ERM:

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

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

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

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

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

[Chamon et al., IEEE TIT'23]

---

## In practice...

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

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

- Update the dual

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

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

---

## In practice...

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

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

- Update the dual

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

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

---

## In practice...

1: Initialize: $\boldsymbol{\theta}_0, \lambda_0$
2: **for** $t = 1, \ldots, T$
3:      $\boldsymbol{\beta}_1 \leftarrow \boldsymbol{\theta}_{t-1}$
4:      **for** $n = 1, \ldots, N$
5:          $\boldsymbol{\beta}_{n+1} \leftarrow \boldsymbol{\beta}_n - \eta_{\theta}\nabla_{\boldsymbol{\beta}}\left[\ell\left(f_{\boldsymbol{\beta}_n}(\boldsymbol{x}_n), y_n\right) + \lambda_{t-1} g\left(f_{\boldsymbol{\beta}_n}(\boldsymbol{x}_n), y_n\right)\right]$    SGD
6:      **end**
7:      $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\beta}_{N+1}$
8:      $\lambda_t = \left[\lambda_{t-1} + \eta_{\lambda}\left(\frac{1}{N}\sum_{m=1}^{N} g\left(f_{\boldsymbol{\theta}_t}(\boldsymbol{x}_m), y_n\right) - c\right)\right]_{+}$    Dual update
9: **end**
10: Output: $\boldsymbol{\theta}_T, \lambda_T$

**PyTorch**

https://github.com/lfochamon/csl

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

Use adaptive method (e.g., ADAM)

○ PyTorch

https://github.com/lfochamon/csl

60

---

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

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

○ PyTorch

https://github.com/lfochamon/csl

60

---

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

Check slack:
- feasibility: $s_t \leq 0$
- "duality gap": $\lambda_t s_t$

$$s_t = \frac{1}{N} \sum_{n=1}^{N} g\left(f_{\theta_t}(\boldsymbol{x}_n), y_n\right) - c$$

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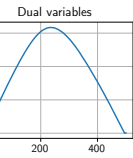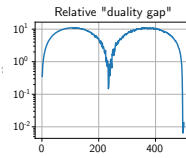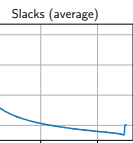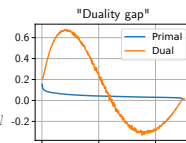Use different time-scales ($\eta_\lambda = 0.1\eta_\theta$)

○ PyTorch

https://github.com/lfochamon/csl

60

---

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

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

https://github.com/lfochamon/csl

60

---

## Penalty-based vs. dual learning

**Penalty-based learning**

$$\theta^\dagger \in \underset{\theta}{\arg\min} \; \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$$

**Dual learning**

$$\theta^\dagger \in \underset{\theta}{\arg\min} \; \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$$

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

- Parameter: $\lambda$ (data-dependent)

- Generalizes with respect to Loss + $\lambda$Penalty

- Parameter: $c$ (requirement-dependent)

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

61

---

## Agenda

Constrained learning theory

Constrained learning algorithms

Resilient constrained learning

62

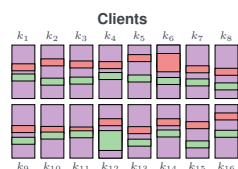---

## Heterogeneous federated learning

**Problem**
Learn a common model using data using data distributed among $K$ clients

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

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

$$k = 1, \ldots, K$$

**Clients**



$k_1$ $k_2$ $k_3$ $k_4$ $k_5$ $k_6$ $k_7$ $k_8$

$k_9$ $k_{10}$ $k_{11}$ $k_{12}$ $k_{13}$ $k_{14}$ $k_{15}$ $k_{16}$

- $k$-th client loss: $\text{Loss}_k(\phi) = \dfrac{1}{N_k} \sum_{n_k=1}^{N_k} \text{Loss}\left(f_\theta(\boldsymbol{x}_{n_k}), y_{n_k}\right)$

63

---

## Heterogeneous federated learning

**Problem**
Learn a common model using data using data distributed among $K$ clients

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

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

$$k = 1, \ldots, K$$

**Clients**



$k_1$ $k_2$ $k_3$ $k_4$ $k_5$ $k_6$ $k_7$ $k_8$

$k_9$ $k_{10}$ $k_{11}$ $k_{12}$ $k_{13}$ $k_{14}$ $k_{15}$ $k_{16}$

- $k$-th client loss: $\text{Loss}_k(\phi) = \dfrac{1}{N_k} \sum_{n_k=1}^{N_k} \text{Loss}\left(f_\theta(\boldsymbol{x}_{n_k}), y_{n_k}\right)$

64

## Heterogeneous federated learning

**Problem**
Learn a common model using data using data distributed among $K$ clients

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

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

$$k = 1, \ldots, K$$

**Clients**

$k_1$ $k_2$ $k_3$ $k_4$ $k_5$ $k_6$ $k_7$ $k_8$

$k_9$ $k_{10}$ $k_{11}$ $k_{12}$ $k_{13}$ $k_{14}$ $k_{15}$ $k_{16}$

- $k$-th client loss: $\mathsf{Loss}_k(\phi) = \dfrac{1}{N_k}\sum_{n_k=1}^{N_k}\mathsf{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_{n_k}), y_{n_k}\big)$

---

## Resilient constrained learning

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

---

## Resilient constrained learning

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

---

## Resilient constrained learning

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

$$P^{\star} = \min_{\boldsymbol{\theta}} \;\; \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\Big[\mathsf{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big]$$

$$\text{subject to} \;\; \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{A}_i}\Big[g_i\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big)\Big] \leq c_i$$

---

## Resilient constrained learning

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

$$P^{\star}(\boldsymbol{r}) = \min_{\boldsymbol{\theta}} \;\; \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\Big[\mathsf{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big]$$

$$\text{subject to} \;\; \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{A}_i}\Big[g_i\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big)\Big] \leq c_i + r_i$$

---

## Resilient constrained learning

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

$$P^{\star}(\boldsymbol{r}) = \min_{\boldsymbol{\theta}} \;\; \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\Big[\mathsf{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big]$$

$$\text{subject to} \;\; \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{A}_i}\Big[g_i\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big)\Big] \leq c_i + r_i$$

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

---

## Resilient constrained learning

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

$$P^{\star}(\boldsymbol{r}) = \min_{\boldsymbol{\theta}} \;\; \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\Big[\mathsf{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big]$$

$$\text{subject to} \;\; \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{A}_i}\Big[g_i\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big)\Big] \leq c_i + r_i$$

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

- Resilience is a compromise!

---

## Resilient constrained learning

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

$$P^{\star}(\boldsymbol{r}) = \min_{\boldsymbol{\theta}} \;\; \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{D}}\Big[\mathsf{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}), y\big)\Big]$$

$$\text{subject to} \;\; \mathbb{E}_{(\boldsymbol{x},y)\sim\mathfrak{A}_i}\Big[g_i\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_m), y_m\big)\Big] \leq c_i + r_i$$

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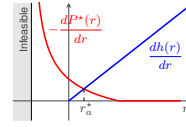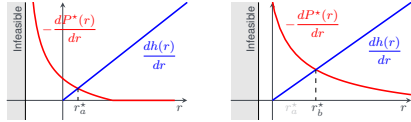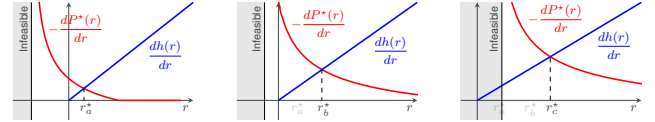
- Resilience is a compromise!

# Resilient constrained learning

**Definition (Resilient equilibrium)**

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

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

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

---

# Resilient constrained learning

**Definition (Resilient equilibrium)**

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

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

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

---

# Resilient constrained learning

**Definition (Resilient equilibrium)**

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

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

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

---

# Resilient constrained learning

**Definition (Resilient equilibrium)**

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

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

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

---

# Resilient constrained learning

**Definition (Resilient equilibrium)**

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

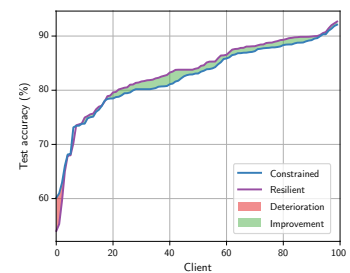$$\nabla h(r^\star) \in -\partial P^\star(r^\star) = \lambda^\star(r^\star)$$

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

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

---

# Resilient constrained learning

**Definition (Resilient equilibrium)**

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

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

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

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

✅ The resilient equilibrium *exists and is unique*  (because $h$ is strictly convex)

---

# Resilient constrained learning
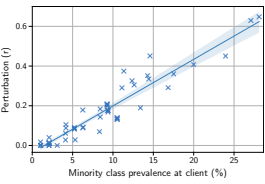
**Definition (Resilient equilibrium)**

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

$$P^\star(r^\star) = \min_{\theta, r} \quad \mathbb{E}_{(x,y)\sim\mathfrak{D}}\left[\mathsf{Loss}\big(f_\theta(x), y\big)\right] + h(r)$$
$$\text{subject to} \quad \mathbb{E}_{(x,y)\sim\mathfrak{A}_i}\left[g_i\big(f_\theta(x_m), y_m\big)\right] \leq c_i + r_i$$

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

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

✅ The resilient equilibrium *exists and is unique*  (because $h$ is strictly convex)

---

# Heterogeneous federated learning

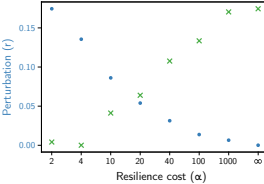## Heterogeneous federated learning

## Heterogeneous federated learning

## Summary

- **Constrained learning is ~~the~~ a tool to learn under requirements**

- **Constrained learning is hard...**

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

## Summary

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

- **Constrained learning is hard...**

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

## Summary

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

- **Constrained learning is hard...**
  Constrained, non-convex, statistical optimization problem

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

## Summary

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

- **Constrained learning is hard...**
  Constrained, non-convex, statistical optimization problem

- **...but possible. How?**
  We can learn under requirements (essentially) whenever we can learn at all by solving *(penalized) ERM problems*. Resilient learning can then be used to adapt the requirements to the task difficulty [Hounie et al., NeurIPS'23]

## Robustness constraints
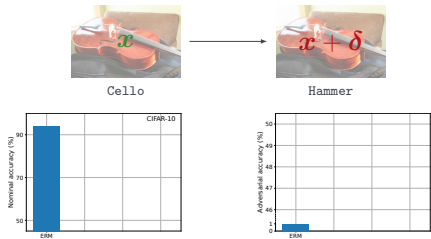
## Agenda

Adversarially robust learning

Semi-infinite learning

Probabilistic robustness

## Robust learning

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

## Adversarial training

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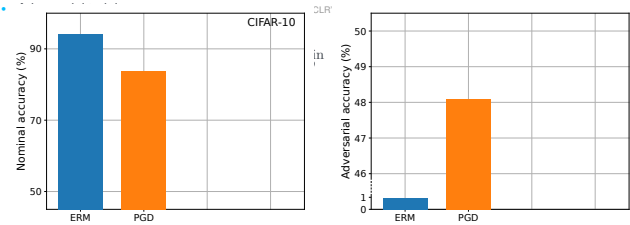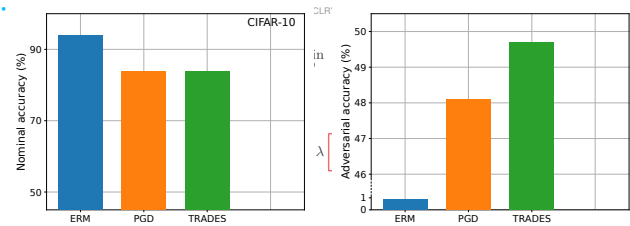Learn an image classifier that is robust to input perturbations

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

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

## Adversarial training

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



[Robey et al., NeurIPS'21]

## Adversarial training

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

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

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

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

## Adversarial training

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



[Zhang et al., ICML'19]

## Constrained learning for robustness

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

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

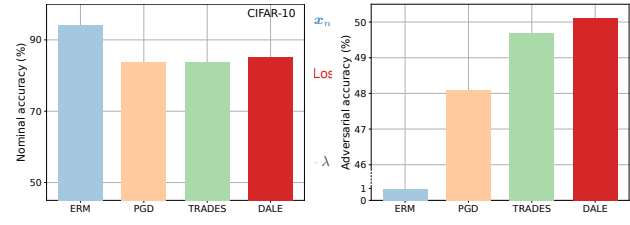$$\text{subject to}\ \frac{1}{N}\sum_{n=1}^{N}\left[\max_{\|\delta\|_{\infty}\le\epsilon}\text{Loss}\big(f_{\theta}(x_n+\delta), y_n\big)\right] \le c$$
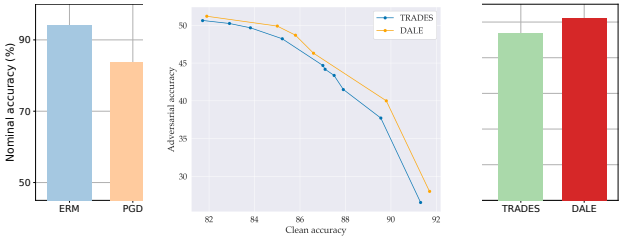
[Chamon and Ribeiro, NeurIPS'20; Robey et al., NeurIPS'21; Chamon et al., IEEE TIT'23]

## Constrained learning for robustness

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



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

## Constrained learning for robustness

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



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

## Penalty-based vs. dual learning

### Penalty-based learning

$$\theta^\dagger \in \arg\min_{\theta} \; \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$$
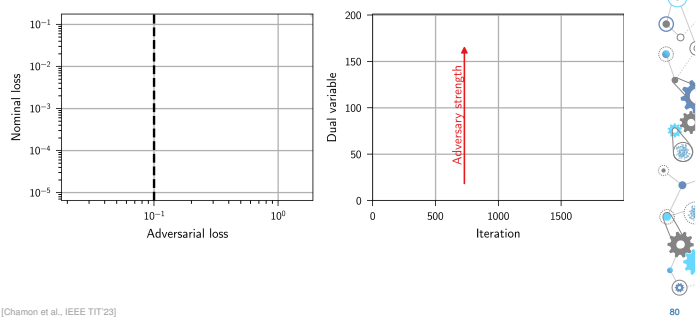
- Parameter: $\lambda$ (data-dependent)
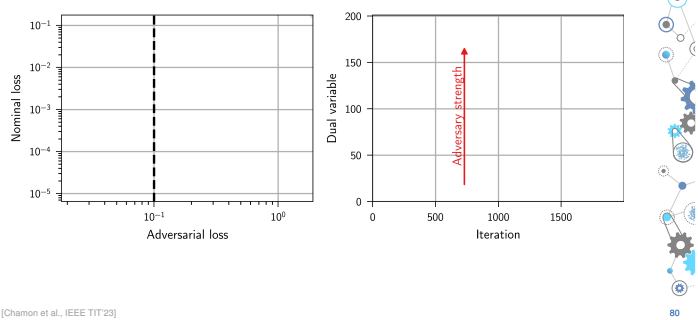- Generalizes with respect to $\text{Loss} + \lambda\text{Penalty}$

### Dual learning

$$\theta^\dagger \in \arg\min_{\theta} \; \text{Loss}(\theta) + \lambda \cdot \text{Penalty}(\theta)$$

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

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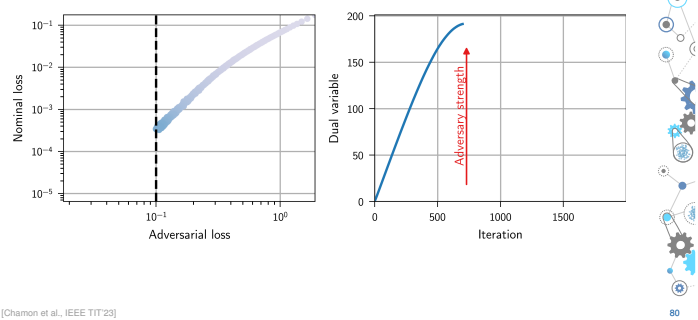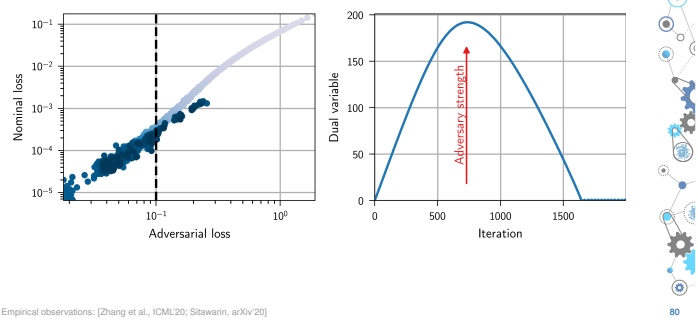- Generalizes with respect to $\text{Loss}$ and $\text{Penalty} \leq c$

79

## Constrained learning for robustness

80

## Constrained learning for robustness

80

## Constrained learning for robustness

80

## Constrained learning for robustness
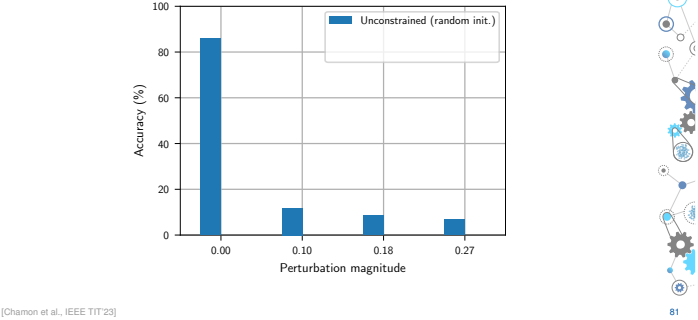
80

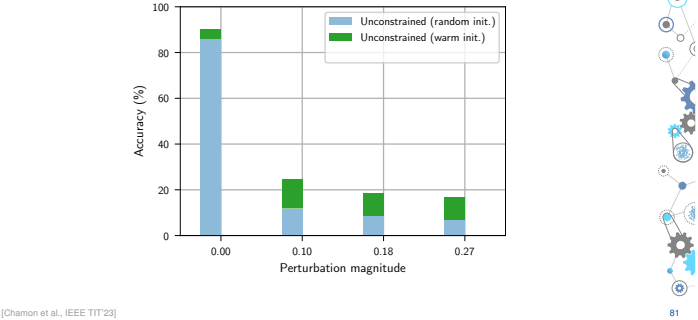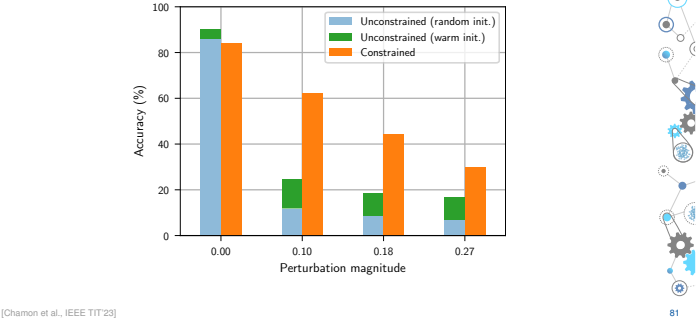## Constrained learning for robustness

81

## Constrained learning for robustness

81

## Constrained learning for robustness

81

## Constrained learning for robustness

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

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

- ✅ Balancing nominal accuracy and robustness ⇒ Dual constrained learning

---

## Constrained learning for robustness

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

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

- ✅ Balancing nominal accuracy and robustness ⇒ Dual constrained learning
- ❌ Computing the worst-case perturbations

---

## Adversarial training

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

- "PGD" [Mądry et al., ICLR'18]

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

---

## Adversarial training

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

- "PGD" [Mądry et al., ICLR'18]

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

- Random initialization
- Restarts
- Pruning
- Adaptive step size

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

---

## Constrained learning for robustness

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

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

- ✅ Balancing nominal accuracy and robustness ⇒ Dual constrained learning
- ❌ Computing the worst-case perturbations
    - ~~gradient ascent~~ → non-convex, underparametrized

---

## Agenda

Adversarially robust learning

Semi-infinite learning

Probabilistic robustness

---

## Semi-infinite constrained learning

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

---

## Semi-infinite constrained learning

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

$$\text{subject to} \quad \text{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n + \boldsymbol{\delta}), y_n\big) \leq t(\boldsymbol{x}_n, y_n),$$
$$\text{for all } (\boldsymbol{x}_n, y_n) \text{ and } \boldsymbol{\delta} \in \Delta$$

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

# Semi-infinite constrained learning

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

$$\text{subject to}\quad \text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}_0),y_n\big)\ \le t(\boldsymbol{x}_n,y_n)$$
$$\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}_{\sqrt 2}),y_n\big)\ \le t(\boldsymbol{x}_n,y_n)$$
$$\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}_e),y_n\big)\ \le t(\boldsymbol{x}_n,y_n)$$
$$\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}_{\pi}),y_n\big)\ \le t(\boldsymbol{x}_n,y_n)$$

- Epigraph formulation:
$$\max_{\|\boldsymbol{\delta}\|_{\infty}\le\epsilon}\text{Loss}\big(f_{\theta}(\boldsymbol{x}+\boldsymbol{\delta}),y\big)\le t \iff \text{Loss}\big(f_{\theta}(\boldsymbol{x}+\boldsymbol{\delta}),y\big)\le t,\ \text{for all }\|\boldsymbol{\delta}\|_{\infty}\le\epsilon$$

- Semi-infinite program

86

# Duality

$$\min_{\theta}\ \frac{1}{N}\sum_{n=1}^{N}\Big[\max_{\boldsymbol{\delta}\in\Delta}\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]$$

$$\updownarrow =$$

$$\min_{\theta}\ \frac{1}{N}\sum_{n=1}^{N}\big[t(\boldsymbol{x}_n,y_n)\big]\ \text{s. to}\ \text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\le t(\boldsymbol{x}_n,y_n),\ \forall(\boldsymbol{x}_n,y_n,\boldsymbol{\delta})$$

$$\updownarrow =$$

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

87

# Duality

$$\min_{\theta}\ \frac{1}{N}\sum_{n=1}^{N}\Big[\max_{\boldsymbol{\delta}\in\Delta}\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]$$

$$\updownarrow =$$

$$\min_{\theta}\ \frac{1}{N}\sum_{n=1}^{N}\big[t(\boldsymbol{x}_n,y_n)\big]\ \text{s. to}\ \text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\le t(\boldsymbol{x}_n,y_n),\ \forall(\boldsymbol{x}_n,y_n,\boldsymbol{\delta})$$

$$\updownarrow =$$

$$\min_{\theta}\ \sup_{\mu\in\mathcal{P}}\ \frac{1}{N}\sum_{n=1}^{N}\underbrace{\mathbb{E}_{\boldsymbol{\delta}\sim\mu(\cdot|\boldsymbol{x}_n,y_n)}\Big[\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]}_{L(\theta,\mu_n)}$$

87

# From optimization to sampling

$$\min_{\theta}\ \frac{1}{N}\sum_{n=1}^{N}\Big[\max_{\boldsymbol{\delta}\in\Delta}\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]$$

$$\updownarrow \approx$$

$$\min_{\theta}\ \sup_{\mu\in\mathcal{P}^2}\ \frac{1}{N}\sum_{n=1}^{N}\underbrace{\mathbb{E}_{\boldsymbol{\delta}\sim\mu_{\gamma}(\cdot|\boldsymbol{x}_n,y_n)}\Big[\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]}_{L(\theta,\mu)}$$

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

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

88

# From optimization to sampling

$$\min_{\theta}\ \frac{1}{N}\sum_{n=1}^{N}\Big[\max_{\boldsymbol{\delta}\in\Delta}\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]$$

$$\updownarrow \approx$$

$$\min_{\theta}\ \sup_{\mu\in\mathcal{P}^2}\ \frac{1}{N}\sum_{n=1}^{N}\underbrace{\mathbb{E}_{\boldsymbol{\delta}\sim\mu_{\gamma}(\cdot|\boldsymbol{x}_n,y_n)}\Big[\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]}_{L(\theta,\mu)}$$

**Proposition**
For any approximation error, $\exists\ \gamma(\boldsymbol{x},y)$ such that

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



[Robey et al., NeurIPS'21]

89

# From optimization to sampling

$$\min_{\theta}\ \frac{1}{N}\sum_{n=1}^{N}\Big[\max_{\boldsymbol{\delta}\in\Delta}\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]$$

$$\updownarrow \approx$$

$$\min_{\theta}\ \sup_{\mu\in\mathcal{P}^2}\ \frac{1}{N}\sum_{n=1}^{N}\underbrace{\mathbb{E}_{\boldsymbol{\delta}\sim\mu_{\gamma}(\cdot|\boldsymbol{x}_n,y_n)}\Big[\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]}_{L(\theta,\mu)}$$

**Proposition**
For any approximation error, $\exists\ \gamma(\boldsymbol{x},y)$ such that

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



[Robey et al., NeurIPS'21]

89

# From optimization to sampling

$$\min_{\theta}\ \frac{1}{N}\sum_{n=1}^{N}\Big[\max_{\boldsymbol{\delta}\in\Delta}\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]$$

$$\updownarrow \approx$$

$$\min_{\theta}\ \sup_{\mu\in\mathcal{P}^2}\ \frac{1}{N}\sum_{n=1}^{N}\underbrace{\mathbb{E}_{\boldsymbol{\delta}\sim\mu_{\gamma}(\cdot|\boldsymbol{x}_n,y_n)}\Big[\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]}_{L(\theta,\mu)}$$

**Proposition**
For any approximation error, $\exists\ \gamma(\boldsymbol{x},y)$ such that

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



[Robey et al., NeurIPS'21]

89

# From optimization to sampling

$$\min_{\theta}\ \frac{1}{N}\sum_{n=1}^{N}\Big[\max_{\boldsymbol{\delta}\in\Delta}\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]$$

$$\updownarrow =$$

$$\min_{\theta}\ \sup_{\mu\in\mathcal{P}^2}\ \frac{1}{N}\sum_{n=1}^{N}\underbrace{\mathbb{E}_{\boldsymbol{\delta}\sim\mu_{\gamma}(\cdot|\boldsymbol{x}_n,y_n)}\Big[\text{Loss}\big(f_{\theta}(\boldsymbol{x}_n+\boldsymbol{\delta}),y_n\big)\Big]}_{L(\theta,\mu)}$$

**Proposition**
For any approximation error, $\exists\ \gamma(\boldsymbol{x},y)$ such that

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



[Robey et al., NeurIPS'21]

89

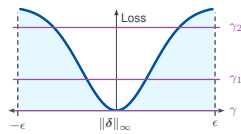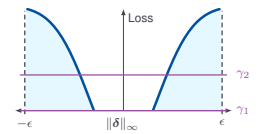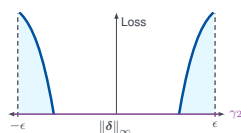## From optimization to sampling

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

$$\updownarrow \approx$$

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

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

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

[Robey et al., NeurIPS'21]

89

---

## Constrained learning for robustness

### Problem
Learn an image classifier that is robust to input perturbations

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

- ✓ Balancing nominal accuracy and robustness ⇒ Dual constrained learning

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

90

---

## Constrained learning for robustness

### Problem
Learn an image classifier that is robust to input perturbations

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

- ✓ Balancing nominal accuracy and robustness ⇒ Dual constrained learning

- ✓ Computing the worst-case perturbations
  - gradient ascent → non-convex, underparametrized ⇒ sampling

[Robey et al., NeurIPS'21]

90

---

## Dual Adversarial LEarning

```
1: for n = 1, ..., N:
2:     δₙ ~ Random(Δ)
3:     for k = 1, ..., K:
4:         ζ ~ Laplace(0, I)
5:         δₙ ← proj_Δ [δₙ + η sign [∇_δ log (Loss(f_θₜ(xₙ+δₙ), yₙ))] + √(2ηT)ζ]
6:     end
7:     θ ← θ - η∇_θ [Loss(f_θ(xₙ), yₙ) + λLoss(f_θ(xₙ+δₙ), yₙ)]
8: end
9: λ ← [λ + η (1/N Σₙ₌₁ᴺ Loss(f_θ(xₙ+δₙ), yₙ) - c)]₊
```

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

SGD

GA

[Robey et al., NeurIPS'21]

91

---

## Dual Adversarial LEarning

```
1: for n = 1, ..., N:
2:     δₙ ~ Random(Δ)
3:     for k = 1, ..., K:
4:         ζ ~ Laplace(0, I)
5:         δₙ ← proj_Δ [δₙ + η sign [∇_δ log (Loss(f_θₜ(xₙ+δₙ), yₙ))] + √(2ηT)ζ]
6:     end
7:     θ ← θ - η∇_θ [Loss(f_θ(xₙ), yₙ) + λLoss(f_θ(xₙ+δₙ), yₙ)]
8: end
9: λ ← [λ + η (1/N Σₙ₌₁ᴺ Loss(f_θ(xₙ+δₙ), yₙ) - c)]₊
```

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

SGD

GA

[Robey et al., NeurIPS'21]

91

---

## Dual Adversarial LEarning

```
1: for n = 1, ..., N:
2:     δₙ ~ Random(Δ)
3:     for k = 1, ..., K:
4:         ζ ~ Laplace(0, I)
5:         δₙ ← proj_Δ [δₙ + η sign [∇_δ log (Loss(f_θₜ(xₙ+δₙ), yₙ))] + √(2ηT)ζ]
6:     end
7:     θ ← θ - η∇_θ [Loss(f_θ(xₙ), yₙ) + λLoss(f_θ(xₙ+δₙ), yₙ)]
8: end
9: λ ← [λ + η (1/N Σₙ₌₁ᴺ Loss(f_θ(xₙ+δₙ), yₙ) - c)]₊
```

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

SGD

GA

[Robey et al., NeurIPS'21]

91

---

## Dual Adversarial LEarning

```
1: for n = 1, ..., N:
2:     δₙ ~ Random(Δ)
3:     for k = 1, ..., K:
4:         ζ ~ Laplace(0, I)
5:         δₙ ← proj_Δ [δₙ + η sign [∇_δ log (Loss(f_θₜ(xₙ+δₙ), yₙ))] + √(2ηT)ζ]
6:     end
7:     θ ← θ - η∇_θ [Loss(f_θ(xₙ), yₙ) + λLoss(f_θ(xₙ+δₙ), yₙ)]
8: end
9: λ ← [λ + η (1/N Σₙ₌₁ᴺ Loss(f_θ(xₙ+δₙ), yₙ) - c)]₊
```

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

SGD

GA

[Robey et al., NeurIPS'21]

92

## Dual Adversarial LEarning

1: **for** $n = 1, \dots, N$:

2:     $\boldsymbol{\delta}_n \sim \text{Random}(\Delta)$

3:     **for** $k = 1, \dots, K$:

4:       $\zeta \sim \text{Laplace}(0, I)$

5:       $\boldsymbol{\delta}_n \leftarrow \underset{\Delta}{\text{proj}} \left[ \boldsymbol{\delta}_n + \eta \text{sign} \left[ \nabla_{\boldsymbol{\delta}} \log \left( \text{Loss}\big(f_{\boldsymbol{\theta}_t}(\boldsymbol{x}_n + \boldsymbol{\delta}_n), y_n\big) \right) \right] + \sqrt{2\eta T}\zeta \right]$

6:     **end**

7:     $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \left[ \text{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) + \lambda \text{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n + \boldsymbol{\delta}_n), y_n\big) \right]$

8: **end**

9: $\lambda \leftarrow \left[ \lambda + \eta \left( \dfrac{1}{N} \sum\limits_{n=1}^{N} \text{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n + \boldsymbol{\delta}_n), y_n\big) - c \right) \right]_+$

Gaussian
[Lopes et al., arXiv'19]
[Rusak et al., ECCV'20]

Patches
[Zhong et al., AAAI'20]
[Yun et al., ICCV'19]

…

SGD

GA

[Robey et al., NeurIPS'21]
92

---

## Dual Adversarial LEarning

1: **for** $n = 1, \dots, N$:

2:     $\boldsymbol{\delta}_n \sim \text{Random}(\Delta)$

3:     **for** $k = 1, \dots, K$:

4:       $\zeta \sim \text{Laplace}(0, I)$

5:       $\boldsymbol{\delta}_n \leftarrow \underset{\Delta}{\text{proj}} \left[ \boldsymbol{\delta}_n + \eta \text{sign} \left[ \nabla_{\boldsymbol{\delta}} \log \left( \text{Loss}\big(f_{\boldsymbol{\theta}_t}(\boldsymbol{x}_n + \boldsymbol{\delta}_n), y_n\big) \right) \right] + \sqrt{2\eta T}\zeta \right]$

6:     **end**

7:     $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \left[ \text{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) + \lambda \text{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n + \boldsymbol{\delta}_n), y_n\big) \right]$

8: **end**

9: $\lambda \leftarrow \left[ \lambda + \eta \left( \dfrac{1}{N} \sum\limits_{n=1}^{N} \text{Loss}\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n + \boldsymbol{\delta}_n), y_n\big) - c \right) \right]_+$

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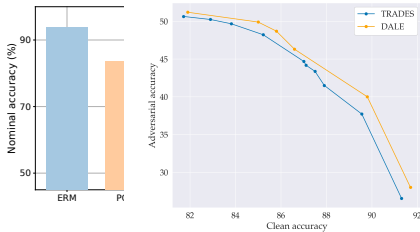[Goodfellow et al., ICLR'15]
[Madry et al., ICLR'18]

SGD

GA

[Robey et al., NeurIPS'21]
92

---

## Dual Adversarial LEarning

### Problem
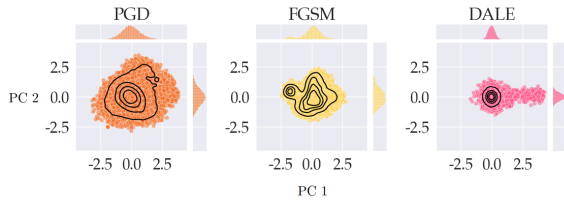Learn an image classifier that is robust to input perturbations



[Robey et al., NeurIPS'21]
93

---

## Dual Adversarial LEarning

### Problem
Learn an image classifier that is robust to input perturbations
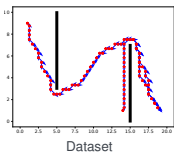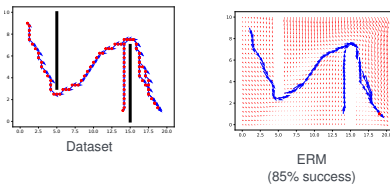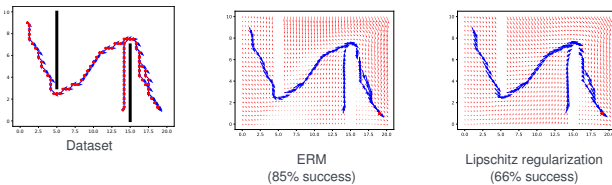


[Robey et al., NeurIPS'21]
94

---

## (Manifold) smoothness

### Problem
Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories

- Labeled data ({State, Action})


Dataset

[Cerviño et al., ICML'23]
95

---

## (Manifold) smoothness

### Problem
Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories

- Labeled data ({State, Action})


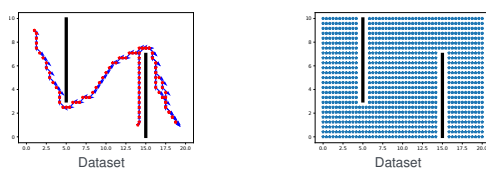Dataset     ERM (85% success)

[Cerviño et al., ICML'23]
95

---

## (Manifold) smoothness

### Problem
Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories

- Labeled data ({State, Action})


Dataset     ERM (85% success)     Lipschitz regularization (66% success)

[Cerviño et al., ICML'23]
95

---

## (Manifold) smoothness

### Problem
Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories

- Labeled data ({State, Action}) and unlabeled data ({State in free space})


Dataset     Dataset

[Cerviño et al., ICML'23]
96

## (Manifold) smoothness

Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories

- Labeled data ({State, Action}) and unlabeled data ({State in free space})

- Use {State in free space} to estimate a data manifold $\mathcal{M}$

$$\min_{\boldsymbol{\theta}} \quad \frac{1}{N} \sum_{n=1}^{N} \|f_{\boldsymbol{\theta}}(\boldsymbol{x}_n) - \boldsymbol{u}_n\|^2$$
$$\text{subject to} \quad \max_{\boldsymbol{x}} \|\nabla_{\mathcal{M}} f_{\boldsymbol{\theta}}(\boldsymbol{x})\|^2 \le c$$
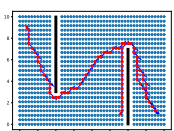
[Cerviño et al., ICML'23]

---

## (Manifold) smoothness

Problem
Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories

- Labeled data ({State, Action}) and unlabeled data ({State in free space})

- Use {State in free space} to estimate a data manifold $\mathcal{M}$

$$\min_{\boldsymbol{\theta}} \quad \frac{1}{N} \sum_{n=1}^{N} \|f_{\boldsymbol{\theta}}(\boldsymbol{x}_n) - \boldsymbol{u}_n\|^2$$
$$\text{subject to} \quad \max_{\boldsymbol{x}} \|\nabla_{\mathcal{M}} f_{\boldsymbol{\theta}}(\boldsymbol{x})\|^2 \le c$$
$$\mathbb{E}_{\boldsymbol{x} \sim \mu_0}$$
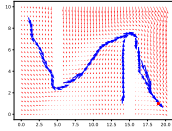
[Cerviño et al., ICML'23]

---

## (Manifold) smoothness

Problem
Learn a smooth (Lipschitz on a manifold) controller that imitates a behavior from limited trajectories
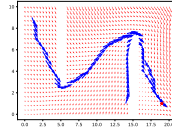
- Labeled data ({Position, Action}) and unlabeled data ({Position})



Dataset      ERM (85% success)      Manifold smoothness (94% success)

[Cerviño et al., ICML'23]

---

## Agenda

Adversarially robust learning

Semi-infinite learning

Probabilistic robustness

---

## Constrained learning challenges

$$\min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^{N} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) \xrightarrow[\text{PACC}]{\text{PAC}} \min_{\boldsymbol{\theta}} \mathbb{E}_{(\boldsymbol{x},y) \sim \mathfrak{D}} \Big[\ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big)\Big]$$
$$\text{s. to } \frac{1}{N} \sum_{n=1}^{N} \Big[\max_{\boldsymbol{\delta} \in \Delta} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n + \boldsymbol{\delta}), y_n\big)\Big] \le c \xrightarrow{\text{PAC}} \text{s. to } \mathbb{E}_{(\boldsymbol{x},y) \sim \mathfrak{D}} \Big[\max_{\boldsymbol{\delta} \in \Delta} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n + \boldsymbol{\delta}), y_n\big)\Big] \le c$$

**Challenges**

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

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

---

## Constrained learning challenges

$$\min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^{N} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big) \xrightarrow[\text{PACC}]{\text{PAC}} \min_{\boldsymbol{\theta}} \mathbb{E}_{(\boldsymbol{x},y) \sim \mathfrak{D}} \Big[\ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n), y_n\big)\Big]$$
$$\text{s. to } \frac{1}{N} \sum_{n=1}^{N} \Big[\max_{\boldsymbol{\delta} \in \Delta} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n + \boldsymbol{\delta}), y_n\big)\Big] \le c \xrightarrow{\text{PAC?}} \text{s. to } \mathbb{E}_{(\boldsymbol{x},y) \sim \mathfrak{D}} \Big[\max_{\boldsymbol{\delta} \in \Delta} \ell\big(f_{\boldsymbol{\theta}}(\boldsymbol{x}_n + \boldsymbol{\delta}), y_n\big)\Big] \le c$$

**Challenges**

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

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

---

## Statistical complexity

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

- Is robust learning harder than non-robust learning? Do we need more samples?
  **A:** YES *and* NO

[Cullina, Bhagoji, Mittal. PAC-learning in the presence of evasion adversaries, NeurIPS'18] ✅ ❌
[Yin, Ramchandran, Bartlett. Rademacher Complexity for Adversarially Robust Generalization, ICML'19] ✅
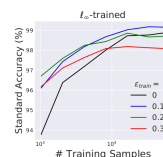[Montasser, Hanneke, Srebro. VC Classes are Adversarially Robustly Learnable, but Only Improperly, COLT'19] ✅ ❌
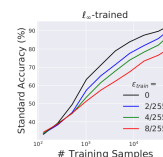[Awasthi, Frank, Mohri. Adversarial Learning Guarantees for Linear Hypotheses and Neural Networks, ICML'20] ✅
[Montasser, Hanneke, Srebro. Adversarially robust learning: A generic minimax optimal learner & characterization, NeurIPS'22] ✅ ❌
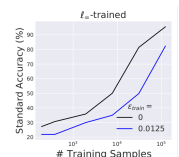
---

## Nominal performance of robust models



(a) MNIST      (b) CIFAR-10      (c) Restricted ImageNet

[Tsipras et al., ICLR'19]

## "Softer" robustness

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

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

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

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

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

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

---

## "Softer" robustness

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

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

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

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

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

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

---

## Towards probabilistic robustness

$$\min_{\theta} \ \frac{1}{N} \sum_{n=1}^{N} \left[ t(x_n, y_n) \right]$$
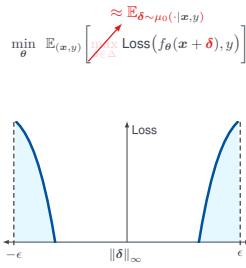
subject to
$$\text{Loss}(f_\theta(x_n + \delta_0), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_1), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_e), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_\pi), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_4), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_{e^2}), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_{\pi^e}), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_{\pi^3}), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_{25}), y_n) \leq t(x_n, y_n)$$

- 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

---

## Towards probabilistic robustness

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

subject to
$$\text{Loss}(f_\theta(x_n + \delta_0), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_1), y_n) \leq t(x_n, y_n)$$
$$\cancel{\text{Loss}(f_\theta(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n)}$$
$$\cancel{\text{Loss}(f_\theta(x_n + \delta_e), y_n) \leq t(x_n, y_n)}$$
$$\text{Loss}(f_\theta(x_n + \delta_\pi), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_4), y_n) \leq t(x_n, y_n)$$
$$\cancel{\text{Loss}(f_\theta(x_n + \delta_{e^2}), y_n) \leq t(x_n, y_n)}$$
$$\text{Loss}(f_\theta(x_n + \delta_{\pi^e}), y_n) \leq t(x_n, y_n)$$
$$\cancel{\text{Loss}(f_\theta(x_n + \delta_{\pi^3}), y_n) \leq t(x_n, y_n)}$$
$$\text{Loss}(f_\theta(x_n + \delta_{25}), y_n) \leq t(x_n, y_n)$$

---

## Towards probabilistic robustness

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

subject to
$$\text{Loss}(f_\theta(x_n + \delta_0), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_1), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_e), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_\pi), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_4), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_{e^2}), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_{\pi^e}), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_{\pi^3}), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_{25}), y_n) \leq t(x_n, y_n)$$

$$\min_{\theta} \ \mathbb{E}_{(x,y)} \left[ \overset{\approx \ \mathbb{E}_{\delta \sim \mu_0(\cdot | x, y)}}{\frac{1}{N} \sum_{\Delta} \text{Loss}(f_\theta(x+\delta), y)} \right]$$
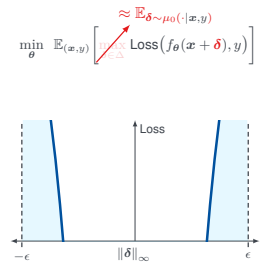


[Robey et al., ICML'22 (spotlight)]

---

## Towards probabilistic robustness

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

subject to
$$\text{Loss}(f_\theta(x_n + \delta_0), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_1), y_n) \leq t(x_n, y_n)$$
$$\cancel{\text{Loss}(f_\theta(x_n + \delta_{\sqrt{2}}), y_n) \leq t(x_n, y_n)}$$
$$\cancel{\text{Loss}(f_\theta(x_n + \delta_e), y_n) \leq t(x_n, y_n)}$$
$$\text{Loss}(f_\theta(x_n + \delta_\pi), y_n) \leq t(x_n, y_n)$$
$$\text{Loss}(f_\theta(x_n + \delta_4), y_n) \leq t(x_n, y_n)$$
$$\cancel{\text{Loss}(f_\theta(x_n + \delta_{e^2}), y_n) \leq t(x_n, y_n)}$$
$$\text{Loss}(f_\theta(x_n + \delta_{\pi^e}), y_n) \leq t(x_n, y_n)$$
$$\cancel{\text{Loss}(f_\theta(x_n + \delta_{\pi^3}), y_n) \leq t(x_n, y_n)}$$
$$\text{Loss}(f_\theta(x_n + \delta_{25}), y_n) \leq t(x_n, y_n)$$

$$\min_{\theta} \ \mathbb{E}_{(x,y)} \left[ \overset{\approx \ \mathbb{E}_{\delta \sim \mu_0(\cdot | x, y)}}{\frac{1}{N} \sum_{\Delta} \text{Loss}(f_\theta(x+\delta), y)} \right]$$
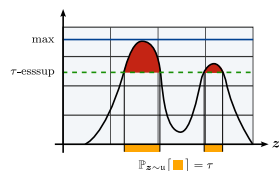


[Robey et al., ICML'22 (spotlight)]

---

## 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)
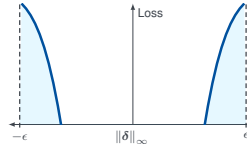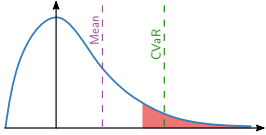  - $\tau = 0$: adversarial robustness (ess sup)



[Robey et al., ICML'22 (spotlight)]

---

## Probabilistic robustness

$$\min_{\theta} \ \mathbb{E}_{(x,y)} \left[ \overset{\mathbb{E}_{\delta \sim \mu_0(\cdot | x, y)}}{\frac{1}{\Delta} \sum_{\Delta} \text{Loss}(f_\theta(x+\delta), y)} \right]$$

$$\min_{\theta} \ \mathbb{E}_{(x,y)} \left[ \overset{\mathbb{E}_{\delta \sim \mu_0(\cdot | x, y)}}{\tau\text{-esssup}_{\Delta} \text{Loss}(f_\theta(x+\delta), y)} \right]$$



[Robey et al., ICML'22 (spotlight)]

## Probabilistic robustness and Risk

- Conditional value at risk:

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

$$= \inf_{\alpha \in \mathbb{R}} \; \alpha + \frac{\mathbb{E}_z\left[[f(z) - \alpha]_+\right]}{1 - \rho}$$

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

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

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

---

## Probabilistically robust learning

1: **for** $n = 1, \ldots, N$:
2: $\quad \alpha_0 = 0$
3: $\quad$ **for** $t = 1, \ldots, T$:
4: $\quad\quad \delta_t \sim \mathsf{Random}(\Delta)$
5: $\quad\quad \alpha \leftarrow \alpha - \frac{\eta}{\tau}\left(\tau - \mathbb{I}\left[\mathsf{Loss}(f_\theta(x_n + \delta_t), y_n) \geq \alpha\right]\right)$ $\quad$ SGD (CVaR)
6: $\quad$ **end**
7: $\quad \theta \leftarrow \theta - \eta\nabla_\theta \underbrace{\left[\mathsf{Loss}(f_\theta(x_n + \delta_T), y_n) - \alpha\right]_+}_{\approx \mathrm{CVaR}_{1-\tau}\left[\mathsf{Loss}(f_\theta(x_n + \delta), y_n)\right]}$ $\quad$ SGD ($\theta$)
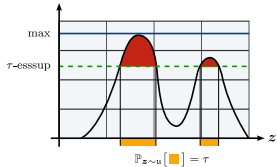8: **end**

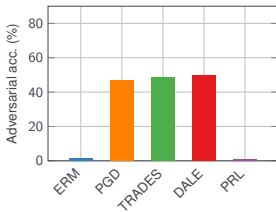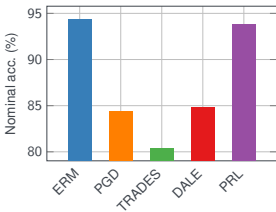[Robey et al., ICML'22 (spotlight)]

---

## Probabilistic robustness

- Probabilistic robustness

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

- $\tau = 1/2$: 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] ✓ ✗
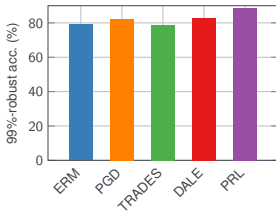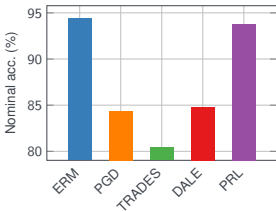- Better performance trade-off
  [Robey et al., ICML'22 (spotlight)] ✓

[Robey et al., ICML'22 (spotlight)]

---

## Probabilistically robust learning



[Robey et al., ICML'22 (spotlight)]

---

## Probabilistically robust learning



[Robey et al., ICML'22 (spotlight)]

---

## Summary

- **Semi-infinite constrained learning is ~~the~~ a tool to enforce worst-case requirements**

- **Semi-infinite constrained learning...**

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

---

## Summary

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

- **Semi-infinite constrained learning...**

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

---

## Summary

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

- **Semi-infinite constrained learning...**
  Learning problem with an infinite number of constraints

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

# Summary

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

- **Semi-infinite constrained learning...**
  Learning problem with an infinite number of constraints

- **...but possible. How?**
  Using a hybrid sampling–optimization algorithm or, in the case of probabilistic robustness,
  a *tight* convex relaxation (CVaR) [Robey et al., ICML'22]

111

**Break**