

CONSTRAINED LEARNING







Why learn under requirements?



Fact: learning has big shortcomings!



Chamon et al.

Constrained Learning

3

How we learn under requirements today





Model structure

Typical for geometrical invariances (CNNs, GNNs, U-net, Dragonnet)

✓ Meets requirements by design

× Hard to design, transfer, or combine



Model structure Typical for geometrical invariances (CNNs, GNNs, U-net, Dragonnet) Meets requirements by design ×

imes Hard to design, transfer, or combine

Performance metric Ubiquitous (GANs, VAEs...)

 $\underset{\boldsymbol{\theta}}{\operatorname{minimize}} \quad f_0(\boldsymbol{\theta}$

$$(\boldsymbol{\theta}) + \sum_{i=1}^m w_i f_i(\boldsymbol{\theta})$$

× Not guaranteed to meet requirements

 $\checkmark\,$ Easy to combine and train

 \checkmark



Develop a theory of constrained learning to

provide tools that enable learning under requirements



- $P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$
- ▶ l_0 is a bounded, Lipschitz continuous, convex function
- \mathcal{F} is a function space (e.g., L_2)
- \mathcal{D} is unknown except for $(\boldsymbol{x}_n, \boldsymbol{y}_n) \sim \mathcal{D}$ (learning)



(P-CSL)

$$P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$$

subject to
$$\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$$

- ▶ l_0, l_i are bounded, Lipschitz continuous, convex functions
- \mathcal{F} is a function space (e.g., L_2)
- \mathcal{D} is unknown except for $(\boldsymbol{x}_n, \boldsymbol{y}_n) \sim \mathcal{D}$ (learning)



(P-CSL)

$$P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(x), y \right) \right]$$

subject to
$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\ell_i \left(\phi(x), y \right) \right] \le c_i$$

- ▶ l_0, l_i are bounded, Lipschitz continuous, convex functions
- \mathcal{F} is a function space (e.g., L_2)
- \mathcal{D} is unknown except for $(\boldsymbol{x}_n, \boldsymbol{y}_n) \sim \mathcal{D}$ (learning)



Definition





(P-CSL)

Definition



- ▶ ℓ_0, ℓ_i are bounded, Lipschitz continuous, convex functions
- Infinite dimensional

• \mathcal{D} is unknown except for $(x_n, y_n) \sim \mathcal{D}$ (learning)



(P-CSL)

Definition



- ▶ ℓ_0, ℓ_i are bounded, Lipschirz continuous, convex functions
- Infinite dimensional

• \mathcal{D} is unknown except for $(\boldsymbol{x}_n, \boldsymbol{y}_n) \sim \mathcal{D}$ (learning)



(P-CSL)



- ℓ_0, ℓ_i are bounded, Lipschitz continuous, convex functions
- Infinite dimensional
- Cannot evaluate expectations







$$P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(x), y \right) \right]$$
subject to $\quad \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\ell_i \left(\phi(x), y \right) \right] \leq c_i$
Main questions
Would (P-CSL) enable learning with requirements?
(P-CSL)



Definition

$$P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$$

subject to
$$\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$$
(P-CSL)

Main questions

Would (P-CSL) enable learning with requirements?

What does it mean to solve (P-CSL)?



Definition

$$P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$$

subject to
$$\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$$
(P-CSL)

Main questions

Would (P-CSL) enable learning with requirements?

What does it mean to solve (P-CSL)?

```
Can we solve (P-CSL) and how?
```



Definition

 $P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$ subject to $\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$ (P-CSL) Main guestions Would (P-CSL) enable learning with requirements? What does it mean to solve (P-CSL)? Can we solve (P-CSL) and how?



Problem

Estimate the probability of an individual making more than US\$ 50,000 based on personal and socio-economical data without discriminating based on gender.



Problem

Estimate the probability of an individual making more than US\$ 50,000 based on personal and socio-economical data without discriminating based on gender.

$\underset{\phi \in \mathcal{F}}{\text{minimize}} \quad \mathbb{E}\left[-y \log\left(\phi(\boldsymbol{x}, \boldsymbol{z})\right)\right]$

▶ *x* collects the features (socio-economical data)

- z is the protected variable (gender)
- $\phi(x, z) = \Pr[\geq US\$ 50.000]$

Fair classification







Problem

Estimate the probability of an individual making more than US\$ 50,000 based on personal and socio-economical data without discriminating based on gender.

$\underset{\phi \in \mathcal{F}}{\text{minimize}} \quad \mathbb{E}\left[-y \log\left(\phi(\boldsymbol{x}, \boldsymbol{z})\right)\right]$

• x collects the features (socio-economical data)

- z is the protected variable (gender)
- $\phi(x, z) = \Pr[\geq US\$ 50.000]$



Problem

Estimate the probability of an individual making more than US\$ 50,000 based on personal and socio-economical data without discriminating based on gender.

 $\underset{\phi \in \mathcal{F}}{\text{minimize}} \quad \mathbb{E}\left[-y \log\left(\phi(\boldsymbol{x}, z)\right)\right]$

subject to $\mathbb{E}[D_{\mathsf{KL}}(\phi(\boldsymbol{x},\mathsf{Male}) \| \phi(\boldsymbol{x},\mathsf{Female}))] \leq c$

x collects the features (socio-economical data)

- z is the protected variable (gender)
- $\phi(x, z) = \Pr[\geq US\$ 50.000]$

Applications



(P-CSL)

 $P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$ subject to $\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$

Applications

Statistical invariances: Fairness, robustness (distributional shift), transfer learning, unbalanced data (within-group accuracy)...

Learning with constraints: semi-supervised learning, VAEs, GANs...



Definition

 $P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$ subject to $\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$ (P-CSL) Questions Would (P-CSL) enable learning with requirements? What does it mean to solve (P-CSL)? Can we solve (P-CSL) and how?



Definition

 $P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$ subject to $\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$ (P-CSL) Questions Would (P-CSL) enable learning with requirements? Yes What does it mean to solve (P-CSL)? Can we solve (P-CSL) and how?



Definition

 $P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$ subject to $\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$ (P-CSL)

Questions

Would (P-CSL) enable learning with requirements? Yes

What does it mean to solve (P-CSL)?

Can we solve (P-CSL) and how?



(P-CSL)

Definition

 $P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$ subject to $\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$

Challenges

- Infinite dimensional
- Cannot evaluate expectations



For $\epsilon_i > 0$, i = 0, ..., m, and $0 \le \delta < 1/2$, a solution ϕ^{\dagger} of (P-CSL) is said to be probably approximately optimal (PAOpt) if



For $\epsilon_i > 0$, i = 0, ..., m, and $0 \le \delta < 1/2$, a solution ϕ^{\dagger} of (P-CSL) is said to be probably approximately optimal (PAOpt) if

1) Probably near-optimal

 $\Pr\left[\left|P^{\star} - \mathbb{E}\left[\ell_{0}\left(\phi^{\dagger}(\boldsymbol{x}), y\right)\right]\right| > \epsilon_{0}\right] \leq \delta$



For $\epsilon_i > 0$, i = 0, ..., m, and $0 \le \delta < 1/2$, a solution ϕ^{\dagger} of (P-CSL) is said to be probably approximately optimal (PAOpt) if

1) Probably near-optimal (\approx PAC learning)

 $\Pr\left[\left|P^{\star} - \mathbb{E}\left[\ell_{0}\left(\phi^{\dagger}(\boldsymbol{x}), y\right)\right]\right| > \epsilon_{0}\right] \leq \delta$



For $\epsilon_i > 0$, i = 0, ..., m, and $0 \le \delta < 1/2$, a solution ϕ^{\dagger} of (P-CSL) is said to be probably approximately optimal (PAOpt) if

1) Probably near-optimal (\approx PAC learning)

$$\Pr\left[\left|P^{\star} - \mathbb{E}\left[\ell_{0}\left(\phi^{\dagger}(\boldsymbol{x}), \boldsymbol{y}\right)\right]\right| > \epsilon_{0}\right] \leq \delta$$

2) Probably approximately feasible

$$\Pr\left[\left|\mathbb{E}\left[\ell_{i}\left(\phi^{\dagger}(\boldsymbol{x}),y\right)\right]\right|>c_{i}+\epsilon_{i}\right]\leq\delta$$



Definition

 $P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$ subject to $\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$ (P-CSL)

Questions

Would (P-CSL) enable learning with requirements? Yes

What does it mean to solve (P-CSL)?

Can we solve (P-CSL) and how?



Definition

$$P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$$

subject to
$$\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$$
 (P-CSL)

Questions

Would (P-CSL) enable learning with requirements? Yes

What does it mean to solve (P-CSL)? PAOpt

Can we solve (P-CSL) and how?



(P-CSL)

Definition

$$P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$$

subject to
$$\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$$

Questions

Would (P-CSL) enable learning with requirements? Yes

What does it mean to solve (P-CSL)? PAOpt

```
Can we solve (P-CSL) and how?
```


$$P^{\star} = \min_{\phi \in \mathcal{F}} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(x), y \right) \right]$$



16

$$P^{\star} = \min_{\phi \in \mathcal{F}} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(x), y \right) \right]$$



$$P^{\star} = \min_{\phi \in \mathcal{F}} \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$$

$$\downarrow$$

$$\hat{P}^{\star}_{\epsilon} = \min_{\theta \in \mathbb{R}^p} \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(f(\theta, \boldsymbol{x}), y \right) \right]$$



 $P^{\star} = \min_{\phi \in \mathcal{F}} \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$ Approximation error Definition For each $\phi \in \mathcal{F}$ there exists $\boldsymbol{\theta} \in \mathbb{R}^p$ such $\hat{P}_{\epsilon}^{\star} = \min_{\boldsymbol{\theta} \in \mathbb{R}^{p}} \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_{0} \left(f(\boldsymbol{\theta}, \boldsymbol{x}), y \right) \right]$ that $\mathbb{E}\left[\left|f(\boldsymbol{\theta}, \boldsymbol{x}) - \phi(\boldsymbol{x})\right|\right] < \nu.$



$$P^{\star} = \min_{\phi \in \mathcal{F}} \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$$

$$\downarrow$$

$$\hat{P}^{\star}_{\epsilon} = \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(f(\boldsymbol{\theta}, \boldsymbol{x}), y \right) \right]$$



$$P^{\star} = \min_{\phi \in \mathcal{F}} \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$$

$$\downarrow$$

$$\hat{P}^{\star}_{\epsilon} = \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(f(\boldsymbol{\theta}, \boldsymbol{x}), y \right) \right]$$

$$\downarrow$$

$$\hat{P}^{\star}_{\epsilon, N} = \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{1}{N} \sum_{n=1}^{N} \ell_0 \left(f(\boldsymbol{\theta}, \boldsymbol{x}_n), y_n \right)$$



$$P^{\star} = \min_{\phi \in \mathcal{F}} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\ell_{0} \left(\phi(x), y \right) \right]$$

$$\hat{P}_{\epsilon}^{\star} = \min_{\theta \in \mathbb{R}^{p}} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\ell_{0} \left(f(\theta, x) \cdot y \right) \right]$$
Theorem (VC learning theory)
Under mild assumptions,

$$P^{\star} \approx \hat{P}_{\epsilon,N}^{\star} \text{ with high probability.}$$



















Theorem

Let f be an ν -parameterization of \mathcal{F} . If $(\theta^{\dagger}, \lambda^{\dagger})$ achieve $\hat{D}_{\epsilon,N}^{\star}$, then $f(\theta^{\dagger}, \cdot)$ is a probably approximately optimal solution of (P-CSL):

 $\left|D_{\epsilon,N}^{\star} - P^{\star}\right| \leq \tilde{O}\left(\nu + \frac{1}{\sqrt{N}}\right)$ with high probability $\mathbb{E}\left[\ell_i(f(\theta^{\dagger}, \boldsymbol{x}), y)\right] \leq c_i + \tilde{O}\left(\frac{1}{\sqrt{N}}\right)$ with high probability



Theorem

Let f be an ν -parameterization of \mathcal{F} . If $(\theta^{\dagger}, \lambda^{\dagger})$ achieve $\hat{D}_{\epsilon,N}^{\star}$, then $f(\theta^{\dagger}, \cdot)$ is a probably approximately optimal solution of (P-CSL):

 $\left|D_{\epsilon,N}^{\star} - P^{\star}\right| \leq \tilde{O}\left(\nu + \frac{1}{\sqrt{N}}\right)$ with high probability $\mathbb{E}\left[\ell_i(f(\theta^{\dagger}, \boldsymbol{x}), y)\right] \leq c_i + \tilde{O}\left(\frac{1}{\sqrt{N}}\right)$ with high probability

Depends on. . .

parametrization quality (ν) requirements difficulty (λ_p^{\star}) sample size (N)



Theorem

Let f be an ν -parameterization of \mathcal{F} . If $(\theta^{\dagger}, \lambda^{\dagger})$ achieve $\hat{D}_{\epsilon,N}^{\star}$, then $f(\theta^{\dagger}, \cdot)$ is a probably approximately optimal solution of (P-CSL):

 $\mathbb{E}\left[\ell_i(f(\boldsymbol{\theta}^{\dagger}, \boldsymbol{x}), y)\right] \leq c_i + 2R \left[\frac{1}{N} + \log\left(\frac{4(2N)^{d_{VC}}}{\delta}\right)\right]$

$$\left| \boldsymbol{D}_{\boldsymbol{\epsilon},\boldsymbol{N}}^{\star} - \boldsymbol{P}^{\star} \right| \leq \left(1 + \left\| \boldsymbol{\lambda}_{p}^{\star} \right\|_{1} \right) L\nu + 2E_{1} \left(\frac{1}{N} \left[1 + \log \left(\frac{4(2N)^{d_{VC}}}{\delta} \right) \right] \quad \text{with prob. } 1 - \delta$$

with prob. $1-\delta$

Depends on. . .

parametrization quality (ν) requirements difficulty (λ_p^{\star}) sample size (N)



Theorem

Let f be an ν -parameterization of \mathcal{F} . If $(\theta^{\dagger}, \lambda^{\dagger})$ achieve $\hat{D}_{\epsilon,N}^{\star}$, then $f(\theta^{\dagger}, \cdot)$ is a probably approximately optimal solution of (P-CSL):

 $\mathbb{E}\left[\ell_i(f(\boldsymbol{\theta}^{\dagger}, \boldsymbol{x}), y)\right] \leq c_i + 2\boldsymbol{\theta} \left[\frac{1}{\lambda} \left(1 + \log\left(\frac{4(2N)^{d_{VC}}}{\delta} \right) \right]$

$$\left| D_{\epsilon,N}^{\star} - P^{\star} \right| \leq \left(1 + \left\| \lambda_p^{\star} \right\|_1 \right) L \nu + 2R_1 \left(\frac{1}{N} \left[1 - \log \left(\frac{4(2N)^{d_{VC}}}{\delta} \right) \right] \text{ with prob. } 1 - \delta$$

with prob. $1-\delta$

Depends on. . .

parametrization quality (u) requirements difficulty (λ_p^{\star}) sample size (



Theorem

Let f be an ν -parameterization of \mathcal{F} . If $(\theta^{\dagger}, \lambda^{\dagger})$ achieve $\hat{D}_{\epsilon,N}^{\star}$, then $f(\theta^{\dagger}, \cdot)$ is a probably approximately optimal solution of (P-CSL):

 $\mathbb{E}\left[\ell_i(f(\boldsymbol{\theta}^{\dagger}, \boldsymbol{x}), y)\right] \leq c_i + 2\beta \left[\frac{1}{N} + \log\left(\frac{4(2N)^{d_{VC}}}{\delta}\right)\right]$

$$\left|D_{\epsilon,N}^{\star} - P^{\star}\right| \leq \left(1 + \left\|\boldsymbol{\lambda}_{p}^{\star}\right\|_{1}\right) L\nu + 2B \left(\frac{1}{N}\left[1 - \log\left(\frac{4(2N)^{d_{vc}}}{\delta}\right)\right] \quad \text{with prob. } 1 - \delta$$

with prob. $1-\delta$

Depends on. . .

parametrization quality (u) requirements difficulty (λ_p^{\star}) sample size



Theorem

Æ

Let f be an ν -parameterization of \mathcal{F} . If $(\theta^{\dagger}, \lambda^{\dagger})$ achieve $\hat{D}_{\epsilon,N}^{\star}$, then $f(\theta^{\dagger}, \cdot)$ is a probably approximately optimal solution of (P-CSL):

$$\begin{split} \left| D_{\epsilon,N}^{\star} - P^{\star} \right| &\leq \left(1 + \left\| \boldsymbol{\lambda}_{p}^{\star} \right\|_{1} \right) L\nu + 2R\sqrt{\frac{1}{N}} \left[1 + \log\left(\frac{4(2N)^{d_{VC}}}{\delta}\right) \right] & \text{with prob. } 1 - \delta \\ \left[\ell_{i}(f(\boldsymbol{\theta}^{\dagger}, \boldsymbol{x}), \boldsymbol{y}) \right] &\leq c_{i} + 2R\sqrt{\frac{1}{N}} \left[1 + \log\left(\frac{4(2N)^{d_{VC}}}{\delta}\right) \right] & \text{with prob. } 1 - \delta \end{split}$$

Depends on. . .

parametrization quality (u) requirements difficulty (λ_p^{\star}) sample size (N)



Theorem

E

Let f be an ν -parameterization of \mathcal{F} . If $(\theta^{\dagger}, \lambda^{\dagger})$ achieve $\hat{D}_{\epsilon,N}^{\star}$, then $f(\theta^{\dagger}, \cdot)$ is a probably approximately optimal solution of (P-CSL):

$$\begin{split} \left| D_{\epsilon,N}^{\star} - P^{\star} \right| &\leq \left(1 + \left\| \boldsymbol{\lambda}_{p}^{\star} \right\|_{1} \right) L\nu + 2R\sqrt{\frac{1}{N}} \left[1 + \log\left(\frac{4(2N)^{d_{vc}}}{\delta}\right) \right] & \text{with prob. } 1 - \delta \\ \left[\ell_{i}(f(\boldsymbol{\theta}^{\dagger}, \boldsymbol{x}), y) \right] &\leq c_{i} + 26 \left[\frac{1}{N} \left[1 + \log\left(\frac{4(2N)^{d_{vc}}}{\delta}\right) \right] & \text{with prob. } 1 - \delta \end{split}$$

Depends on. . .

parametrization quality (u) requirements difficulty (λ_p^{\star}) sample size (N)



Theorem

Let f be an ν -parameterization of \mathcal{F} . If $(\theta^{\dagger}, \lambda^{\dagger})$ achieve $\hat{D}_{\epsilon,N}^{\star}$, then $f(\theta^{\dagger}, \cdot)$ is a probably approximately optimal solution of (P-CSL):

$$\begin{split} \left| D_{\epsilon,N}^{\star} - P^{\star} \right| &\leq \left(1 + \left\| \boldsymbol{\lambda}_{p}^{\star} \right\|_{1} \right) L\nu + 2R\sqrt{\frac{1}{N}} \left[1 + \log\left(\frac{4(2N)^{d_{\mathcal{V}c}}}{\delta}\right) \right] & \text{with prob. } 1 - \delta \\ & \mathsf{E}\left[\ell_{i}(f(\boldsymbol{\theta}^{\dagger}, \boldsymbol{x}), y) \right] \leq c_{i} + 2R\sqrt{\frac{1}{N}} \left[1 + \log\left(\frac{4(2N)^{d_{\mathcal{V}c}}}{\delta}\right) \right] & \text{with prob. } 1 - \delta \end{split}$$

Depends on...

parametrization quality (u) requirements difficulty (λ_p^{\star}) sample size (N)

Agenda



(P-CSL)

Definition

$$P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$$

subject to
$$\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$$

Questions

Would (P-CSL) enable learning with requirements? Yes

What does it mean to solve (P-CSL)? PAOpt

Can we solve (P-CSL) and how?

Agenda



(P-CSL)

Definition

$$P^{\star} = \min_{\phi \in \mathcal{F}} \quad \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_0 \left(\phi(\boldsymbol{x}), y \right) \right]$$

subject to
$$\mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\ell_i \left(\phi(\boldsymbol{x}), y \right) \right] \le c_i$$

Questions

Would (P-CSL) enable learning with requirements? Yes

What does it mean to solve (P-CSL)? PAOpt

Can we solve (P-CSL) and how? Yes



Estimate the probability of an individual making more than US\$ 50,000 based on personal and socio-economical data without discriminating based on gender.

 $\min_{\phi \in \mathcal{F}} \mathbb{E} \left[-y \log \left(\phi(\boldsymbol{x}, z) \right) \right]$

subject to $\mathbb{E}\left[D_{\mathsf{KL}}(\phi(\boldsymbol{x},\mathsf{Male}) \| \phi(\boldsymbol{x},\mathsf{Female}))\right] \leq c$

x collects the features (socio-economical data)

z is the protected variable (gender)

• $\phi(x, z) = \Pr[\geq US\$ 50.000]$



Estimate the probability of an individual making more than US\$ 50,000 based on personal and socio-economical data without discriminating based on gender.

$$\max_{\lambda \ge 0} \min_{\boldsymbol{\theta} \in \mathcal{H}} - \frac{1}{N} \sum_{n=1}^{N} \left(y_n \log \left(f(\boldsymbol{\theta}, \boldsymbol{x}_n, \boldsymbol{z}_n) \right) - \lambda \left[D_{\mathsf{KL}} \left(f(\boldsymbol{\theta}, \boldsymbol{x}_n, \mathsf{Male}) \parallel f(\boldsymbol{\theta}, \boldsymbol{x}_n, \mathsf{Female}) \right) - c \right] \right)$$

x collects the features (socio-economical data)

- z is the protected variable (gender)
- $f(\theta, \boldsymbol{x}, \boldsymbol{z}) = \Pr\left[\geq \mathsf{US}\$ \ 50.000\right]$



Estimate the probability of an individual making more than US\$ 50,000 based on personal and socio-economical data without discriminating based on gender.

$$\max_{\lambda \ge 0} \min_{\boldsymbol{\theta} \in \mathcal{H}} - \frac{1}{N} \sum_{n=1}^{N} \left(y_n \log \left(f(\boldsymbol{\theta}, \boldsymbol{x}_n, \boldsymbol{z}_n) \right) - \lambda \left[D_{\mathsf{KL}} \left(f(\boldsymbol{\theta}, \boldsymbol{x}_n, \mathsf{Male}) \, \| \, f(\boldsymbol{\theta}, \boldsymbol{x}_n, \mathsf{Female}) \right) - c \right] \right)$$

x collects the features (socio-economical data)

- z is the protected variable (gender)
- $f(\theta, \boldsymbol{x}, \boldsymbol{z}) = \Pr\left[\geq \mathsf{US}\$ \ 50.000\right]$



Estimate the probability of an individual making more than US\$ 50,000 based on personal and socio-economical data without discriminating based on gender.

$$\max_{\lambda \ge 0} \min_{\boldsymbol{\theta} \in \mathcal{H}} - \frac{1}{N} \sum_{n=1}^{N} \left(y_n \log \left(f(\boldsymbol{\theta}, \boldsymbol{x}_n, \boldsymbol{z}_n) \right) - \lambda \left[D_{\mathsf{KL}} \left(f(\boldsymbol{\theta}, \boldsymbol{x}_n, \mathsf{Male}) \parallel f(\boldsymbol{\theta}, \boldsymbol{x}_n, \mathsf{Female}) \right) - c \right] \right)$$

x collects the features (socio-economical data)

- z is the protected variable (gender)
- $f(\theta, \boldsymbol{x}, \boldsymbol{z}) = \Pr\left[\geq \mathsf{US}\$ \ 50.000\right]$



Estimate the probability of an individual making more than US\$ 50,000 based on personal and socio-economical data without discriminating based on gender.

$$\max_{\lambda \ge 0} \min_{\boldsymbol{\theta} \in \mathcal{H}} - \frac{1}{N} \sum_{n=1}^{N} \left(y_n \log \left(f(\boldsymbol{\theta}, \boldsymbol{x}_n, \boldsymbol{z}_n) \right) - \lambda \left[D_{\mathsf{KL}} \left(f(\boldsymbol{\theta}, \boldsymbol{x}_n, \mathsf{Male}) \parallel f(\boldsymbol{\theta}, \boldsymbol{x}_n, \mathsf{Female}) \right) - c \right] \right)$$

x collects the features (socio-economical data)

z is the protected variable (gender)

►
$$f(\theta, x, z) = \Pr[\geq \mathsf{US} \ 50.000] \rightarrow \mathsf{neural network}$$

Back to fairness





Back to fairness





Fair classification





Extensions and further applications







-, e'a

Extensions and further applications



- ► Non-convex losses ℓ₀, ℓ_i (under mild assumptions on D) [future work]
- Primal-dual algorithm convergence [future work]



- ► Non-convex losses l₀, l_i (under mild assumptions on D) [future work]
- Primal-dual algorithm convergence [future work]
- Constrained (e.g., safe) reinforcement learning

[PCFR, "Constrained Reinforcement Learning Has Zero Duality Gap." *NeurIPS*, 2019] [PFCR, "Safe Policies for Reinforcement Learning via Primal-Dual Methods." *ArXiv*, 2019]



(1) Does constrained learning enable learning with requirements?

(2) What does it mean to solve a constrained learning problem?

(3) Can we solve (P-CSL) and how?



- (1) **Does constrained learning enable learning with requirements?** Yes, e.g., fair classification.
- (2) What does it mean to solve a constrained learning problem?

(3) Can we solve (P-CSL) and how?



- (1) **Does constrained learning enable learning with requirements?** Yes, e.g., fair classification.
- (2) What does it mean to solve a constrained learning problem? Obtaining a probably approximately optimal solution.
- (3) Can we solve (P-CSL) and how?


- (1) **Does constrained learning enable learning with requirements?** Yes, e.g., fair classification.
- (2) What does it mean to solve a constrained learning problem? Obtaining a probably approximately optimal solution.
- (3) Can we solve (P-CSL) and how?

Yes, using their empirical dual and the approximation depends on the sample size, the model complexity, and the requirements' difficulty.



CONSTRAINED LEARNING

