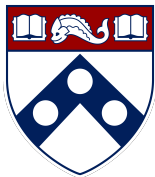


Optimal multi-distribution learning



Yuxin Chen

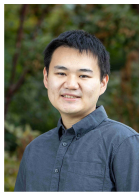
Statistics & Data Science, Wharton, UPenn



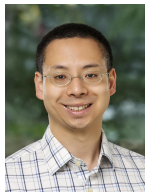
Zihan Zhang
Princeton



Wenhao Zhan
Princeton

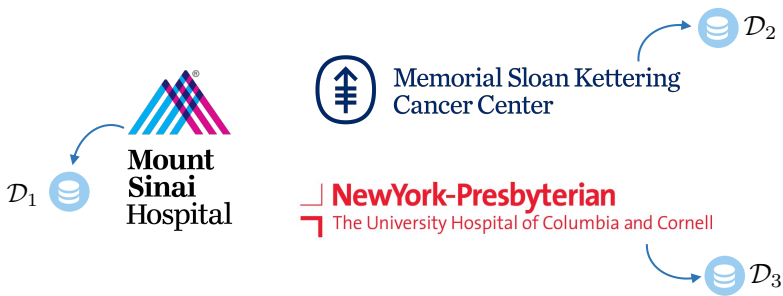


Simon Du
UWashington



Jason Lee
Princeton

“Optimal multi-distribution learning,” Z. Zhang, W. Zhan, Y. Chen, S. Du, J. Lee,
arXiv:2312.05134, 2023



In multi-distribution learning, an agent aims to learn a *shared model* to fit multiple (unknown) data distributions

- diverse data sources (e.g., localities, communities, populations)
- heterogeneous objectives \rightarrow need a balance



- k unknown data distributions $\mathcal{D}_1, \dots, \mathcal{D}_k$ (e.g., localities, communities, populations)
- hypothesis class \mathcal{H} : VC dimension d
- known loss function ℓ (e.g., misclassification error)

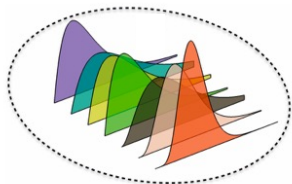


- k unknown data distributions $\mathcal{D}_1, \dots, \mathcal{D}_k$ (e.g., localities, communities, populations)
- hypothesis class \mathcal{H} : VC dimension d
- known loss function ℓ (e.g., misclassification error)

goal: learn an ε -optimal possibly random hypothesis \hat{h} (in **min-max** sense)

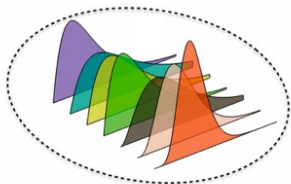
$$\max_{1 \leq i \leq k} \mathbb{E}_{(x,y) \sim \mathcal{D}_i, \hat{h}} [\ell(\hat{h}, (x,y))] \leq \min_{h \in \mathcal{H}} \max_{1 \leq i \leq k} \mathbb{E}_{(x,y) \sim \mathcal{D}_i} [\ell(h, (x,y))] + \varepsilon$$

Mohri et al. '19, Sagawa et al. '19, Blum et al. '17, Buhlmann et al. '15, Guo '23 ...



distributionally robust learning

Mohri et al. '19, Sagawa et al. '19, Blum et al. '17, Buhlmann et al. '15, Guo '23 ...

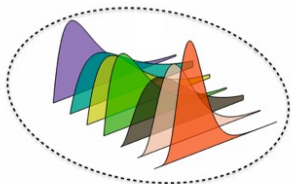


distributionally robust learning



min-max fairness

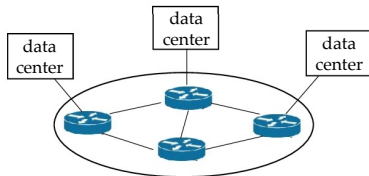
Mohri et al. '19, Sagawa et al. '19, Blum et al. '17, Buhlmann et al. '15, Guo '23 ...



distributionally robust learning



min-max fairness



collaborative learning

Adaptive vs. non-adaptive sampling

- **non-adaptive sampling:** pre-determine sample-size budgets for each distribution beforehand
 - loss of data efficiency

Adaptive vs. non-adaptive sampling

- **non-adaptive sampling**: pre-determine sample-size budgets for each distribution beforehand
 - loss of data efficiency
- **adaptive sampling**: sample on demand during learning process
 - this talk

Adaptive vs. non-adaptive sampling

- **non-adaptive sampling**: pre-determine sample-size budgets for each distribution beforehand
→ loss of data efficiency
- **adaptive sampling**: sample on demand during learning process
→ this talk

learning 1 distribution



$$\frac{d}{\epsilon^2}$$

Adaptive vs. non-adaptive sampling

- **non-adaptive sampling**: pre-determine sample-size budgets for each distribution beforehand
→ loss of data efficiency
- **adaptive sampling**: sample on demand during learning process
→ this talk

learning 1 distribution



$$\frac{d}{\epsilon^2}$$

learning k distributions



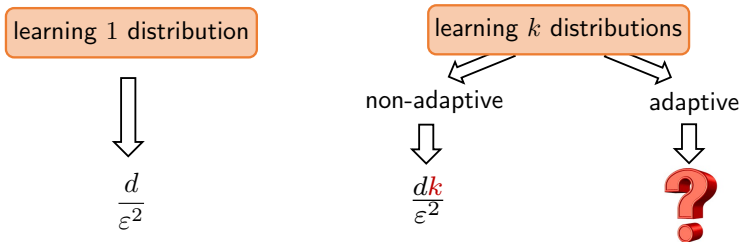
non-adaptive



$$\frac{dk}{\epsilon^2}$$

Adaptive vs. non-adaptive sampling

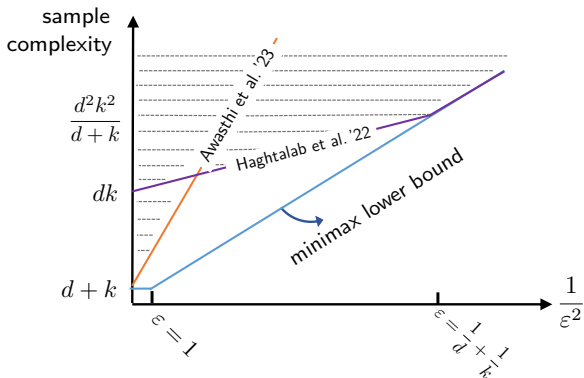
- **non-adaptive sampling**: pre-determine sample-size budgets for each distribution beforehand
→ loss of data efficiency
- **adaptive sampling**: sample on demand during learning process
→ this talk



Prior works: VC classes

paper	sample complexity
Haghtalab et al. '22	$\frac{d+k}{\epsilon^2} + \frac{dk}{\epsilon}$
Awasthi et al. '23	$\frac{d}{\epsilon^4} + \frac{k}{\epsilon^2}$
(lower bound) Haghtalab et al. '22	$\frac{d+k}{\epsilon^2}$

Prior works: VC classes



paper	sample complexity
Haghtalab et al. '22	$\frac{d+k}{\epsilon^2} + \frac{dk}{\epsilon}$
Awasthi et al. '23	$\frac{d}{\epsilon^4} + \frac{k}{\epsilon^2}$
(lower bound) Haghtalab et al. '22	$\frac{d+k}{\epsilon^2}$

Can we close the gap between achievability and lower bound?

**Open Problem: The Sample Complexity of Multi-Distribution
Learning for VC Classes**

Pranjal Awasthi

Google Research, Mountain View, CA, USA

Nika Haghtalab

University of California, Berkeley, CA, USA

Eric Zhao

University of California, Berkeley, CA, USA

PRANJALAWASTHI@GOOGLE.COM

NIKA@BERKELEY.EDU

ERIC.ZH@BERKELEY.EDU

Main results

Theorem 1 (Zhang, Zhan, Chen, Du, Lee '23)

We can design an algorithm that returns randomized hypothesis \hat{h} s.t.

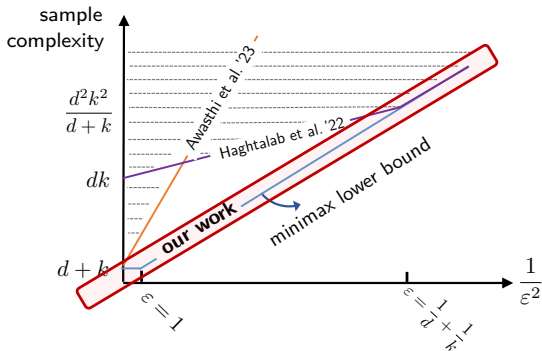
$$\max_{1 \leq i \leq k} \mathbb{E}_{(x,y) \sim \mathcal{D}_i, \hat{h}} \left[\ell(\hat{h}, (x, y)) \right] \leq \min_{h \in \mathcal{H}} \max_{1 \leq i \leq k} \mathbb{E}_{(x,y) \sim \mathcal{D}_i} \left[\ell(h, (x, y)) \right] + \varepsilon,$$

with sample complexity

$$\tilde{O} \left(\frac{d + k}{\varepsilon^2} \right)$$

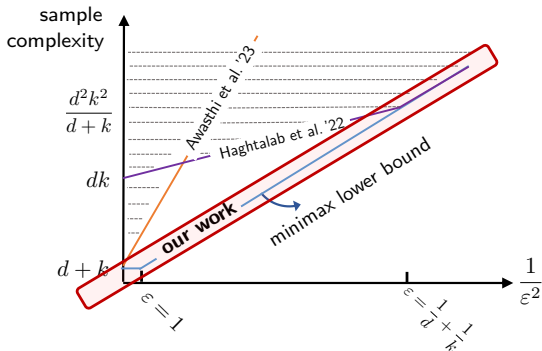
- matches the minimax lower bound (up to log factors)

Main results



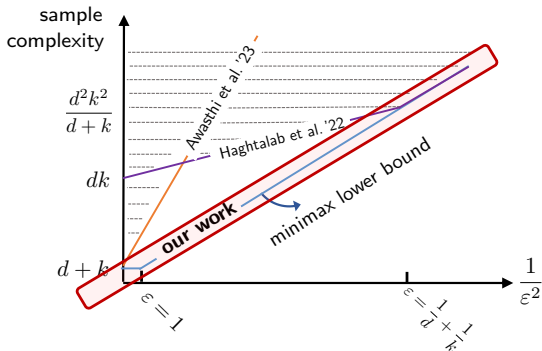
- matches the minimax lower bound (up to log factors)
- solves a COLT open problem (concurrent work: Peng '23)

Main results



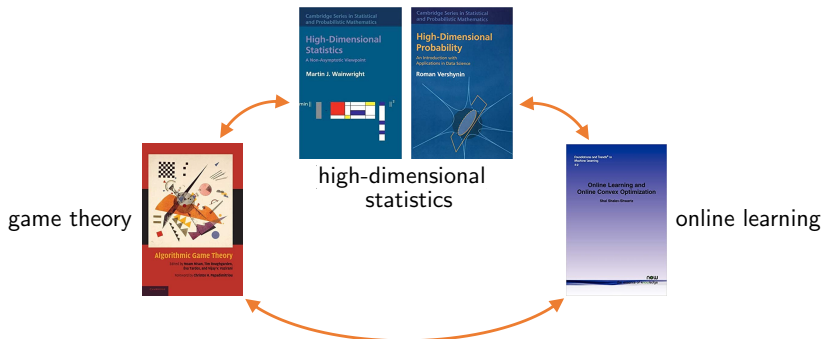
- matches the minimax lower bound (up to log factors)
- solves a COLT open problem (concurrent work: Peng '23)
- can be extended to Rademacher classes

Main results



- matches the minimax lower bound (up to log factors)
- solves a COLT open problem (concurrent work: Peng '23)
- can be extended to Rademacher classes
- algorithm is oracle-efficient (solves another COLT open problem)
only needs to call ERM oracle

Algorithm design



A game-theoretic view

$$\min_{h \in \mathcal{H}} \max_{1 \leq i \leq k} \mathbb{E}_{(x,y) \sim \mathcal{D}_i} [\ell(h, (x, y))]$$



min-player:

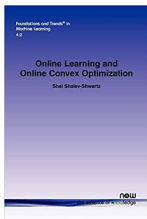
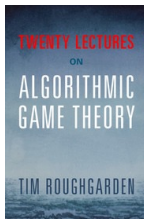
finding most favorable hypothesis



max-player:

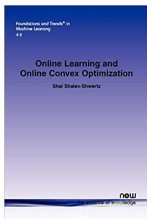
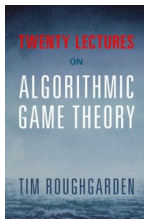
finding least favorable distribution

Preliminaries: learning in games

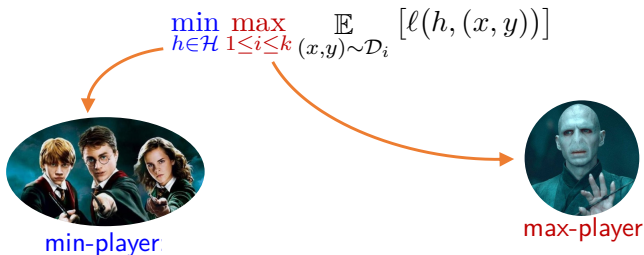


- **no-regret algorithm**: online algorithm w/ $\underbrace{\text{sub-linear regret}}_{\frac{1}{T} \text{Regret}(T) \rightarrow 0}$ over any adversary
 - e.g., **Hedge** algorithm (equivalent to online mirror descent)

Preliminaries: learning in games



- **no-regret algorithm**: online algorithm w/ $\underbrace{\text{sub-linear regret}}_{\frac{1}{T} \text{Regret}(T) \rightarrow 0}$ over any adversary
 - e.g., **Hedge** algorithm (equivalent to online mirror descent)
- **best-response**: play argmin or argmax (not always no-regret)



- min-player/max-player: no-regret/no-regret (Haghtalab et al. '22)

$$\frac{d + k}{\varepsilon^2} + \frac{dk}{\varepsilon} \quad (\text{burn-in due to covering of } \mathcal{H})$$

- min-player/max-player: best-response/no-regret (Awasthi et al. '23)

$$\frac{d}{\varepsilon^4} + \frac{k}{\varepsilon^2} \quad (\text{lack of sample reuse})$$

Our approach: best-response/no-regret

At iteration t :

- min-player computes **empirical best response**

$$h^t \approx \arg \min_{h \in \mathcal{H}} L(h, w^t)$$

$$\circ L(h, w) := \sum_{i=1}^k w_i \mathbb{E}_{(x,y) \sim \mathcal{D}_i} [\ell(h, (x, y))] \text{ (loss w.r.t. weighted dist)}$$

Our approach: best-response/no-regret

At iteration t :

- min-player computes **empirical best response**

$$h^t \approx \arg \min_{h \in \mathcal{H}} L(h, w^t)$$

- $L(h, w) := \sum_{i=1}^k w_i \mathbb{E}_{(x,y) \sim \mathcal{D}_i} [\ell(h, (x, y))]$ (loss w.r.t. weighted dist)

- max-player runs **Hedge** to update $\underbrace{\text{mixed distribution } w^t \in \Delta_k}_{\text{weighted distribution } \sum_i w_i^t \mathcal{D}_i}$

Our approach: best-response/no-regret

At iteration t :

- min-player computes **empirical best response**

$$h^t \approx \arg \min_{h \in \mathcal{H}} L(h, w^t)$$

- $L(h, w) := \sum_{i=1}^k w_i \mathbb{E}_{(x,y) \sim \mathcal{D}_i} [\ell(h, (x, y))]$ (loss w.r.t. weighted dist)

- max-player runs **Hedge** to update mixed distribution $w^t \in \Delta_k$
weighted distribution $\sum_i w_i^t \mathcal{D}_i$

$$w_i^t \propto w_i^{t-1} \exp(\eta \hat{r}_i^t) \quad \text{with } \hat{r}_i^t : \text{ empirical risk for } \mathcal{D}_i$$

Our approach: best-response/no-regret

At iteration t :

- min-player computes **empirical best response**

$$h^t \approx \arg \min_{h \in \mathcal{H}} L(h, w^t)$$

$$\circ L(h, w) := \sum_{i=1}^k w_i \mathbb{E}_{(x,y) \sim \mathcal{D}_i} [\ell(h, (x, y))] \quad (\text{loss w.r.t. weighted dist})$$

- max-player runs **Hedge** to update mixed distribution $w^t \in \Delta_k$
weighted distribution $\sum_i w_i^t \mathcal{D}_i$

$$w_i^t \propto w_i^{t-1} \exp(\eta \hat{r}_i^t) \quad \text{with } \hat{r}_i^t : \text{ empirical risk for } \mathcal{D}_i$$

Output: randomized hypothesis $\hat{h} \sim \text{Uniform}(\{h^t\}_{1 \leq t \leq T})$

Key algorithmic distinction from prior work

adaptive sampling + sample reuse
#samples from \mathcal{D}_i based on $\{w_i^t\}$

Key algorithmic distinction from prior work

adaptive sampling + sample reuse
#samples from \mathcal{D}_i based on $\{w_i^t\}$

Sampling strategy at iteration t :

- **best-response:** have $\frac{d+k}{\epsilon^2} w_i^t$ samples available from \mathcal{D}_i
reuse samples

Key algorithmic distinction from prior work

adaptive sampling + sample reuse
#samples from \mathcal{D}_i based on $\{w_i^t\}$

Sampling strategy at iteration t :

- **best-response:** have $\frac{d+k}{\epsilon^2} \max_{1 \leq \tau \leq t} w_i^\tau$ samples available from \mathcal{D}_i
reuse samples

Key algorithmic distinction from prior work

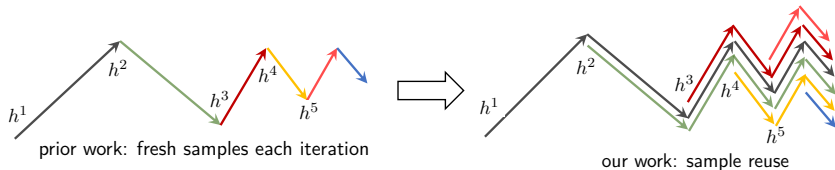
adaptive sampling + sample reuse
#samples from \mathcal{D}_i based on $\{w_i^\tau\}$

Sampling strategy at iteration t :

- **best-response:** have $\frac{d+k}{\epsilon^2} \max_{1 \leq \tau \leq t} w_i^\tau$ samples available from \mathcal{D}_i
reuse samples
- **no-regret:** draw $k \max_{1 \leq \tau \leq t} w_i^\tau$ samples from \mathcal{D}_i
fresh samples

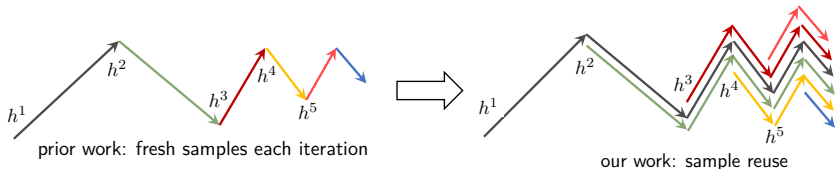
Key technical challenges

1. complicated statistical dependency due to sample reuse



Key technical challenges

1. complicated statistical dependency due to sample reuse

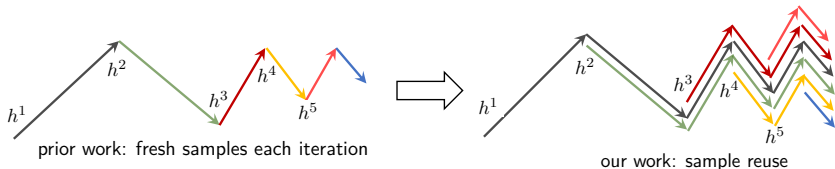


2. need to bound the algorithm trajectory in a fine-grained manner

$$\text{sample complexity} \asymp \frac{d + k}{\varepsilon^2} \underbrace{\sum_{i=1}^k \max_{1 \leq t \leq T} w_i^t}$$

Key technical challenges

1. complicated statistical dependency due to sample reuse



2. need to bound the algorithm trajectory in a fine-grained manner

$$\text{sample complexity} \asymp \frac{d+k}{\varepsilon^2} \underbrace{\sum_{i=1}^k \max_{1 \leq t \leq T} w_i^t}_{\tilde{O}(1)}$$

concentration + doubling trick + combinatorics

Concurrent work: Peng et al. '23

Peng et al. '23 established a sample complexity of

$$\frac{d + k}{\epsilon^2} \left(\frac{k}{\epsilon} \right)^{o(1)}$$

which also solved the COLT open problem

Concurrent work: Peng et al. '23

Peng et al. '23 established a sample complexity of

$$\frac{d + k}{\epsilon^2} \left(\frac{k}{\epsilon} \right)^{o(1)}$$

which also solved the COLT open problem

- optimal up to some sub-polynomial term
- a very different algorithm
 - recursive structure to eliminate non-optimal hypotheses

Necessesity of randomization



Our alg. returns randomized hypothesis ...

Necessesity of randomization



Our alg. returns randomized hypothesis ...

Question: is it possible to find an ϵ -optimal deterministic hypothesis w/ the same sample complexity (**another COLT open problem**)?

Necessesity of randomization



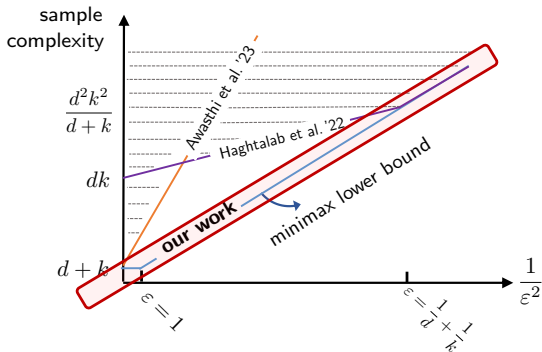
Our alg. returns randomized hypothesis ...

Question: is it possible to find an ε -optimal deterministic hypothesis w/ the same sample complexity (**another COLT open problem**)?

Answer: No!

- finding an ε -optimal deterministic policy needs $\Omega\left(\frac{dk}{\varepsilon^2}\right)$ samples

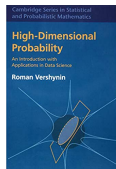
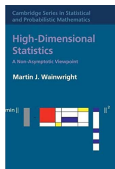
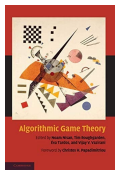
Summary: multi-distribution learning



- settles the sample complexity of MDL under on-demand sampling
- solves 3 COLT open problems posed by Awasthi et al. '23

Concluding remarks

Advancing frontier of statistical learning requires integrated thinking of modern statistics, optimization & game theory



online learning & games

(high-dimensional) statistics

"Optimal multi-distribution learning," Z. Zhang, W. Zhan, Y. Chen, S. Du, J. Lee, arXiv:2312.05134, 2023