

# Efficiently Certifiable Guarantees for Learning with Distribution Shift

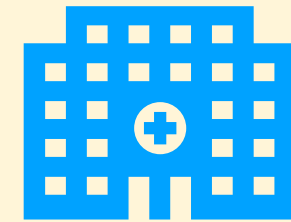
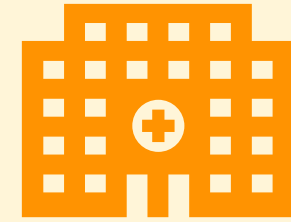


Konstantinos Stavropoulos (UT Austin) based on joint works with: Gautam Chandrasekaran, Surbhi Goel, Adam Klivans, Vasilis Kontonis, Lin Lin Lee, Abhishek Shetty, Arsen Vasilyan

## Learning in the presence of distribution shift

Hospital A: *Labeled data*  
(treatment outcomes)

Hospital B: *Only unlabeled data*  
(patient medical records)



**Goal: Decide treatments for B using A**



**Critical Setting: Guarantees Needed**

**Information-theoretic impossibility**

e.g. uniformity testing in high dimensions requires exponential number of samples



**Attempt #1:** Estimate Statistical Distance between Training and Test distributions

**Attempt #2:** Estimate the Discrepancy (Domain Adaptation)

[BBCP'06] [BCKPW'08] [MMR'09]

Discrepancy distance between distributions  $\mathcal{D}, \mathcal{D}'$ ,  
with respect to class of boolean functions  $\mathcal{C}$ :

$$\text{disc}_{\mathcal{C}}(\mathcal{D}, \mathcal{D}') = \max_{f_1, f_2 \in \mathcal{C}} \left| \Pr_{\mathcal{D}}[f_1(x) \neq f_2(x)] - \Pr_{\mathcal{D}'}[f_1(x) \neq f_2(x)] \right|$$

Out of distribution generalization for classification:

$$\Pr_{\mathcal{D}'}[h(x) \neq f^*(x)] \leq \Pr_{\mathcal{D}}[h(x) \neq f^*(x)] + \text{disc}_{\mathcal{C}}(\mathcal{D}, \mathcal{D}')$$



Small discrepancy between training and test distributions  
 $\Rightarrow$  low test error  
(assuming labels always generated by unknown  $f^* \in \mathcal{C}$ )



Sample complexity of estimating the discrepancy  
is bounded by the VC-dimension of  $\mathcal{C}$



Testing if the discrepancy is small is NP-hard,  
even when  $\mathcal{C}$  is the class of linear classifiers

[CKKS'24]  
[BGS'18]

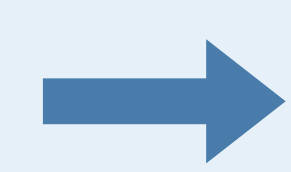
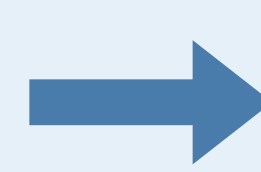
## Our Approach: Testable Learning with Distribution Shift

**Goal: Certify at least one of the following**

1. The test distribution is different from the training (detect distribution shift)
2. The error of the output classifier  $h$  on the test distribution is low

**Input**

$(x \sim \mathcal{D}, f^*(x))$   
 $(x \sim \mathcal{D}')$



**Output**

Either (**REJECT**,  $\perp$ )  
or (**ACCEPT**,  $h$ )

- **Soundness:** Upon acceptance  $\Pr_{x \sim \mathcal{D}'} [h(x) \neq f^*(x)] \leq \epsilon$  (w.h.p.)
- **Completeness:** If  $\mathcal{D} = \mathcal{D}'$ , then accept (w.h.p.)



Certify low test error  
with **no assumptions**  
on test marginal



Results from classical learning theory for  
**most concept classes** can be enhanced  
with **efficient testers** for TDS learning

Concept Class	Training Marginal	TDS Runtime	PAC Runtime
Halfspaces	Standard d-dim Gaussian	$d^{\tilde{O}(\log 1/\epsilon)}$	$\text{poly}(d, 1/\epsilon)$
Intersections of k Halfspaces	Standard d-dim Gaussian	$d^{\tilde{O}(\log 1/\epsilon)}(k/\epsilon)^{O(k^3)}$	$\text{poly}(d)(k/\epsilon)^{O(k^2)}$
Degree-2 PTFs	Standard Gaussian or Uniform	$d^{\tilde{O}(1/\epsilon^9)}$	$\text{poly}(d, 1/\epsilon)$
Circuits of size s, depth t	Uniform on d-dim hypercube	$d^{O(\log(s))^{O(t)} \log(1/\epsilon)}$	$d^{O(\log(s))^{t-1} \log(1/\epsilon)}$

[KS'24a,b] [CKKS'24]

++

## Broader Context: Testable Learning

[Rubinfeld & Vasilyan '23]

**Goal: Get rid of strong assumptions**

**Instead of assumptions:** Provide certifiable guarantees when

1. Direct validation of guarantees
  2. Direct verification of assumptions
- Are intractable or impossible

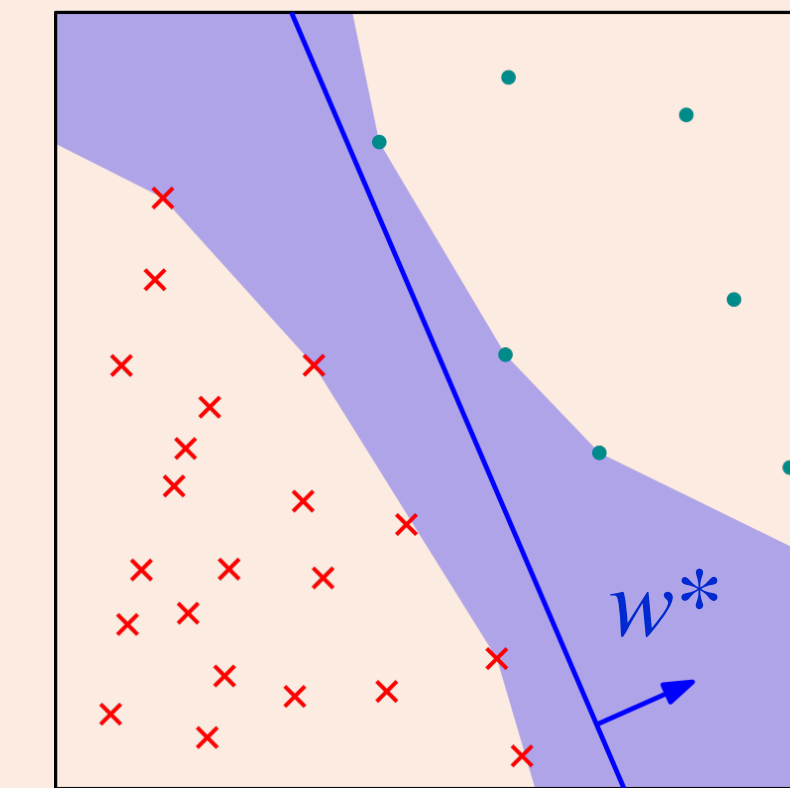
**Testable Agnostic Learning** [RV'23]: No distribution shift, but unknown error benchmark ("optimum error"). *Assumption to be removed:* marginal is well-behaved  
**Testable Noise Assumptions** [GKS'25]: Again, no distribution shift. *Assumption to be removed:* The label noise is structured

## Case Study: Halfspaces

Suppose:  $f^*(x) = \text{sign}(w^* \cdot x - \tau^*)$ ,  $w^* \in \mathbb{S}^{d-1}$ ,  $\tau^* \in \mathbb{R}$

**Simpler result:** There is an  $\epsilon$ -TDS learner with runtime  $\text{poly}(d) \cdot 2^{O(1/\epsilon)}$

**Case I: Low bias**  
( $|\tau^*| \leq 1/\epsilon^{1/2}$ )



**Parameter Recovery**

**Lemma:**  $\mathbb{E}_{x \sim \mathcal{N}}[f^*(x)x] = \frac{\exp(-\tau^{*2}/2)}{\sqrt{2\pi}} w^*$  [DKS'18]

**Find:**  $\hat{w} : \|\hat{w} - w^*\|_2 \leq O(\epsilon/d)$

$\hat{\tau} : \|\hat{\tau} - \tau^*\|_2 \leq O(\epsilon/d)$

using  $\text{poly}(d/\epsilon) \cdot 2^{O(\tau^{*2})}$  samples

**Concern:**  $f^*, \hat{f}$  disagree often under  $\mathcal{D}'$

**Check the following condition:**

$$\mathbb{P}_{x \sim \mathcal{D}'} \left[ \exists w', \tau' : \begin{array}{l} \|\hat{w} - w'\| \leq O(\epsilon/d) \\ \|\hat{\tau} - \tau'\| \leq O(\epsilon/d) \\ \hat{f}(x) \neq f'(x) \end{array} \right] \leq O(\epsilon)$$

**No information about  $w^*$**

**Concern:**  $w^* \cdot x$  often large under  $\mathcal{D}'$

**Concentration Certificates**

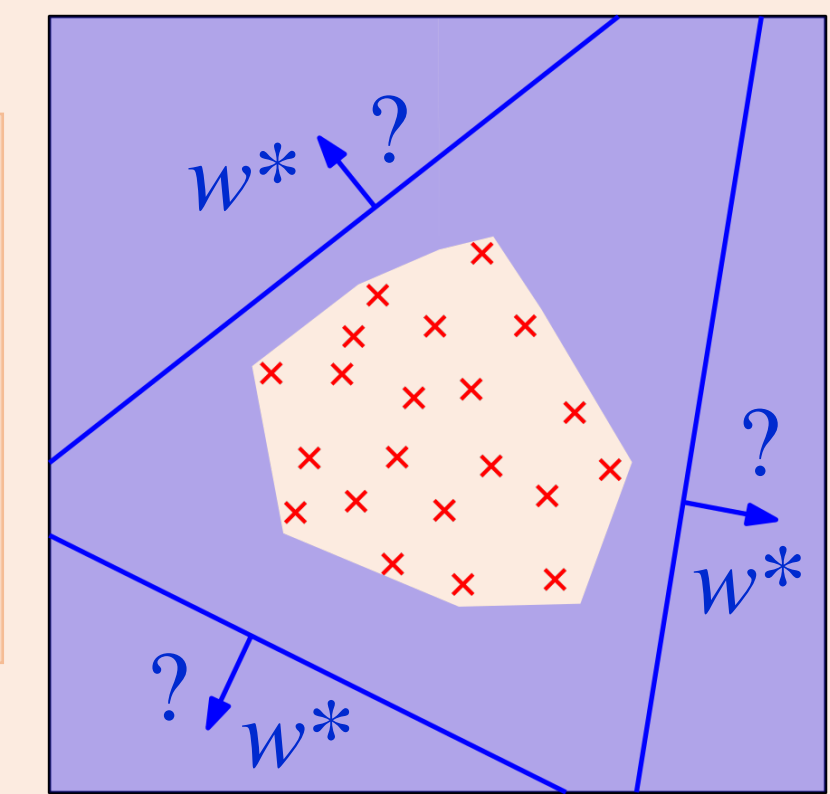
$$\mathbb{P}_{x \sim \mathcal{D}'}[w^* \cdot x > \tau^*] \leq \mathbb{P}_{x \sim \mathcal{D}'}[w^* \cdot x > \epsilon^{-1/2}] \leq \epsilon \mathbb{E}_{x \sim \mathcal{D}'}[(w^* \cdot x)^2]$$

**Check:**  $\mathbb{E}_{\mathcal{D}}[x_i x_j] \approx \mathbb{E}_{\mathcal{N}}[x_i x_j], \forall i, j$

So:  $\mathbb{E}_{x \sim \mathcal{D}'}[(w^* \cdot x)^2] \approx \mathbb{E}_{x \sim \mathcal{N}}[(w^* \cdot x)^2] = 1$

**For tight results:** Check degree  $\log(1/\epsilon)$

**Case II: High bias ( $\tau^* > 1/\epsilon^{1/2}$ )**



**Universal TDS Learners**  
**Handle Benign Shifts**

**Accept** whenever  $\mathcal{D}'$  such that:

1.  $\mathbb{E}_{x \sim \mathcal{D}'}[(v \cdot x)^4] \leq C, \forall v \in \mathbb{S}^{d-1}$

2.  $\mathbb{P}_{x \sim \mathcal{D}'}[|v \cdot x| \leq r] \leq Cr, \forall v \in \mathbb{S}^{d-1}$

Instead of moment matching:

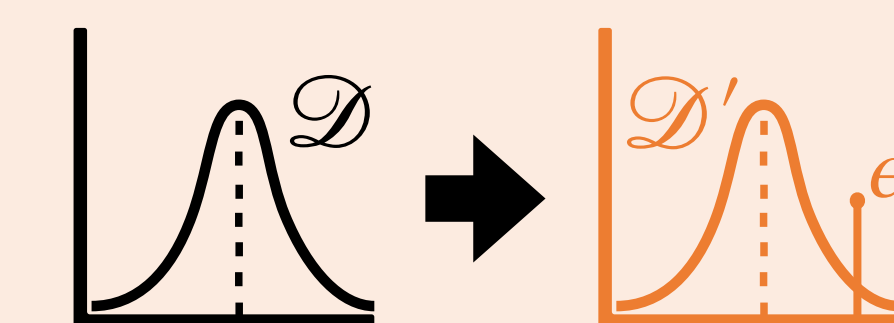
**Check:**  $\sup_{v \in \mathbb{S}^{d-1}} \mathbb{E}_{\mathcal{D}'}[(v \cdot x)^2] \leq C$   
using eigenvalue decomposition

**For improved runtime:** Certify subgaussianity via SoS [DHPT'24]

**Tolerant TDS Learners** [GSSV'24]

**Handle Moderate Amounts of Shift**

**Accept** whenever  $d_{TV}(\mathcal{D}, \mathcal{D}') \leq \epsilon$



**While:**  $\sup_{v \in \mathbb{S}^{d-1}} \mathbb{E}_{\mathcal{D}'}[(v \cdot x)^2] > 10$

**Find**  $r$  s.t.  $\mathbb{P}_{x \sim \mathcal{D}'}[(v_{\max} \cdot x)^2 > r]$

is at least  $2 \mathbb{P}_{x \sim \mathcal{N}}[(v_{\max} \cdot x)^2 > r]$

**Condition**  $\mathcal{D}'$  on  $(v_{\max} \cdot x)^2 \leq r$

**Check:** Mass of rejected region  $O(\epsilon)$