



## Introduction

- ML systems may be biased towards particular groups
- Existing approaches mainly **evaluate** fairness
- Important & challenging to rigorously **certify** fairness, which is our focus

### Main Contributions:

- We formulate** certified fairness problem of an end-to-end ML model
- We propose an effective fairness certification framework** that **for the first time** solves this certified fairness problem by subpopulation decomposition
- We evaluate** our framework on **6** real-world datasets to show its tightness and scalability

## Core Methodology:

### Subpopulation Decomposition

Decompose according to sensitive attribute  $X_s$  and label  $Y$

$$\mathcal{P} = \sum_{s=1}^S \sum_{y=1}^C \Pr[X_s = s, Y = y] \cdot \mathcal{P}_{s,y}$$

$$\mathcal{Q} = \sum_{s=1}^S \sum_{y=1}^C \Pr_{\mathcal{Q}}[X_s = s, Y = y] \cdot \mathcal{Q}_{s,y}$$

## Problem Formulation

Given model  $h_{\theta}(\cdot)$ , compute an upper bound of its expected loss on a fair test distribution  $\mathcal{Q}$ , i.e.,

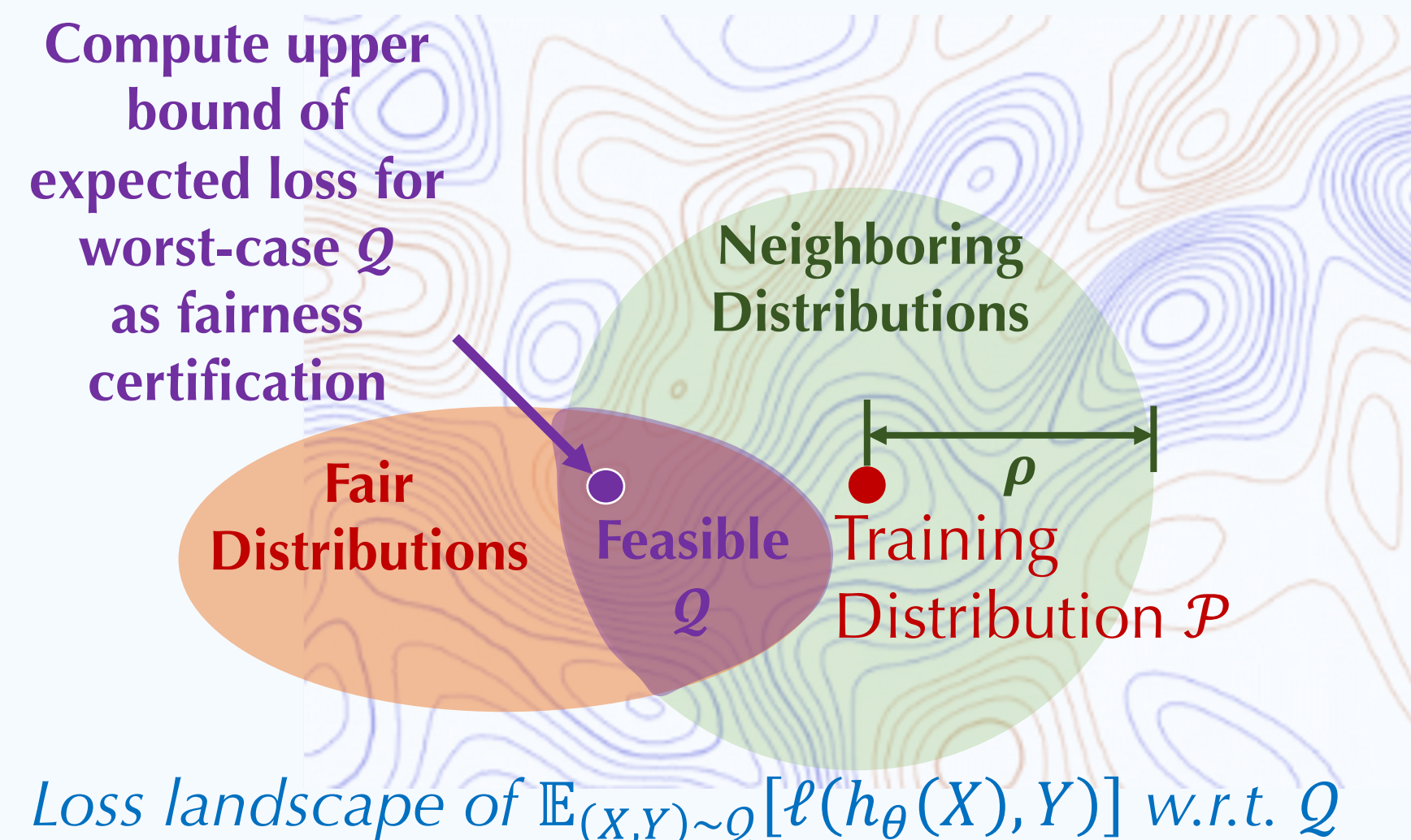
**upper bounding**  $\max_{\mathcal{Q}} \mathbb{E}_{(X,Y) \sim \mathcal{Q}}[\ell(h_{\theta}(X), Y)]$

s. t.  $\text{dist}(\mathcal{P}, \mathcal{Q}) \leq \rho, \quad \mathcal{Q}$  is a fair distribution

Test distribution  $\mathcal{Q}$  not too far from training distribution

Measure performance on distribution with fair base rate

where  $\rho$  User-specified distance threshold  
 $\text{dist}(\mathcal{P}, \mathcal{Q})$  Hellinger distance between distributions



## Theoretical Observations

<b>Distance Constraint</b>	Decomposed to constraints on $\mathcal{P}_{s,y}, \mathcal{Q}_{s,y}$
<b>Fair Distribution Constraint</b>	Equal to constraints on $\Pr_{\mathcal{P}}[X_s = s, Y = y], \Pr_{\mathcal{Q}}[X_s = s, Y = y]$

### Fair Distribution Constraint

- Consider discrete sensitive attribute  $X_s$  and label  $Y$
- Define fair distribution to be distribution with fair base rate:

$$\Pr_{(X,Y) \sim \mathcal{Q}}[Y = y | X_s = s_a] = \Pr_{(X,Y) \sim \mathcal{Q}}[Y = y | X_s = s_b], \forall y, s_a, s_b$$

- Sensitive attribute  $X_s$  has no effect on label  $Y$  at population level

- Such fair distribution admits unconstrained parameterization:

$$\Pr_{(X,Y) \sim \mathcal{Q}}[Y = y | X_s = s] = k_s r_y \quad (k_s, r_y \in [0,1])$$

## Certification Procedure (Informal, Theorem 3)

**Input:** subpopulation statistics & subpopulation level constraints

- Query subpopulation statistics:

$$\Pr_{\mathcal{P}}[X_s = s, Y = y], \mathbb{E}_{(X,Y) \sim \mathcal{P}_{s,y}}[\ell(h_{\theta}(X), Y)]$$

- Divide  $k_s, r_y \in [0,1]$  into grids

- In each grid:

- Known quantities:**

$$\Pr_{\mathcal{P}}[X_s = s, Y = y], \Pr_{\mathcal{Q}}[X_s = s, Y = y], \mathbb{E}_{(X,Y) \sim \mathcal{P}_{s,y}}[\ell(h_{\theta}(X), Y)]$$

- Variables to optimize:**  $\text{dist}(\mathcal{P}_{s,y}, \mathcal{Q}_{s,y})^2$  (subject to distance constraints)

- Key variable to upper bound:**  $\mathbb{E}_{(X,Y) \sim \mathcal{Q}_{s,y}}[\ell(h_{\theta}(X), Y)]$

- Plug in Gramian bound [Weber et al, ICML 2022] to get upper bound
- Optimize the upper bound with low-dimensional convex optimization
- Bypass non-convexity with variable transforms

- Maximization over all grids  $\Rightarrow$  **Output:** Certification of fairness!

Remarks:

- For **sensitive shifting** setting (no distribution shift within each subpopulation, only portions among subpopulations shifted), we have **simpler** fairness certification procedure with **tighter** guarantees
- Framework **amenable to finite sampling error**: with high-confidence intervals of statistics, we provide high-confidence probabilistic certification.
- Framework **support any population loss function**, e.g., can bound group risk discrepancy
- Our fairness notion implies demographic parity (**DP**) and equalized odds (**EO**)

## Experimental Evaluation

### Conclusions

**Tightness:** distance between gray points and black curve

- **Usually tight, especially in sensitive shifting setting**

**Soundness:** gray points always below black curve

- **Always sound**

More results & ablation studies in our paper!

x-axis: distance threshold  $\rho$   
y-axis: expected loss

For sensitive shifting setting:

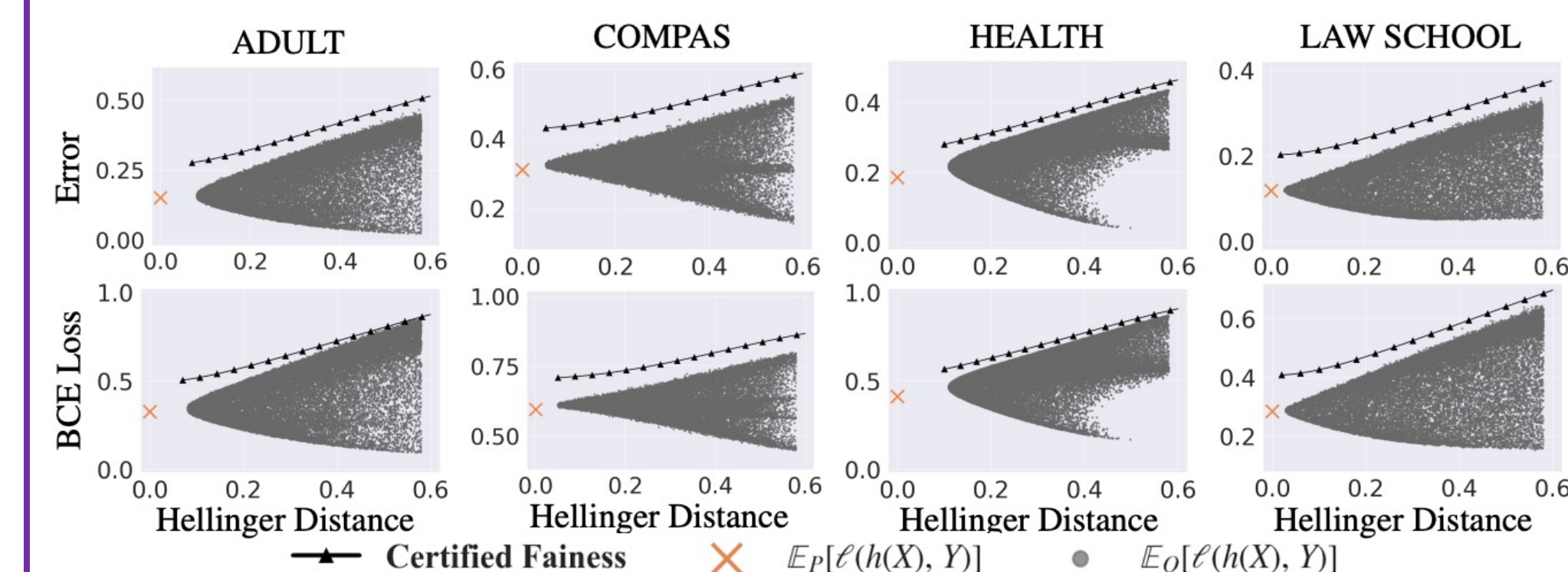


Figure 1: Certified fairness with sensitive shifting. Grey points are results on generated distributions ( $\mathcal{Q}$ ) and the black line is our fairness certificate based on Thm. 2. We observe that our fairness certificate is usually tight.

For general shifting setting:

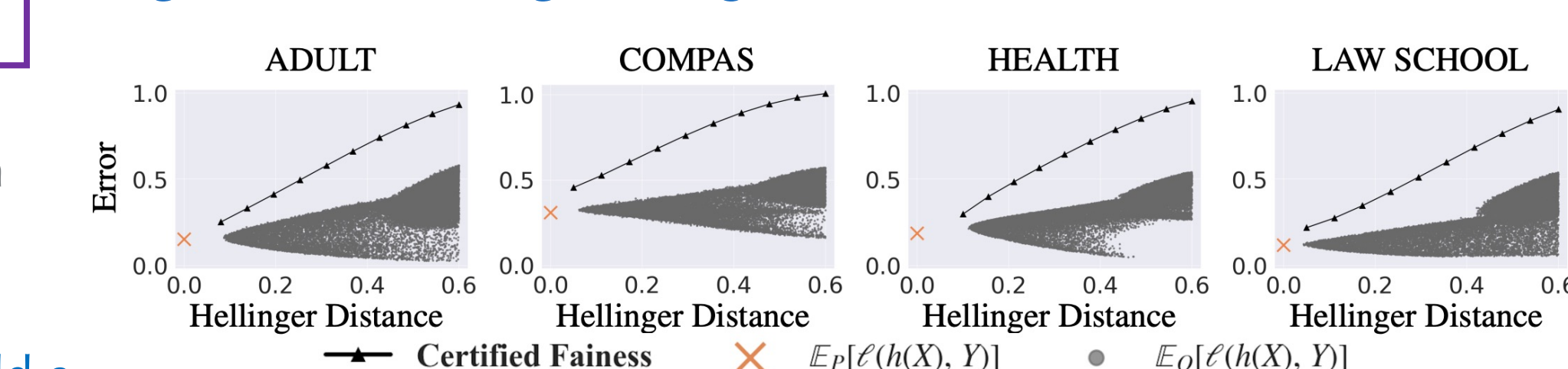


Figure 2: Certified fairness with general shifting. Grey points are results on generated distributions ( $\mathcal{Q}$ ) and the black line is our fairness certificate based on Thm. 3. We observe that our fairness certificate is non-trivial.

### Distance constraint

$$\text{dist}(\mathcal{P}, \mathcal{Q}) \leq \rho \Leftrightarrow$$

$$1 - \rho^2 - \sum_{s=1}^S \sum_{y=1}^C \sqrt{\Pr_{\mathcal{P}}[X_s = s, Y = y] \Pr_{\mathcal{Q}}[X_s = s, Y = y]} (1 - \text{dist}(\mathcal{P}_{s,y}, \mathcal{Q}_{s,y})^2) \leq 0$$