

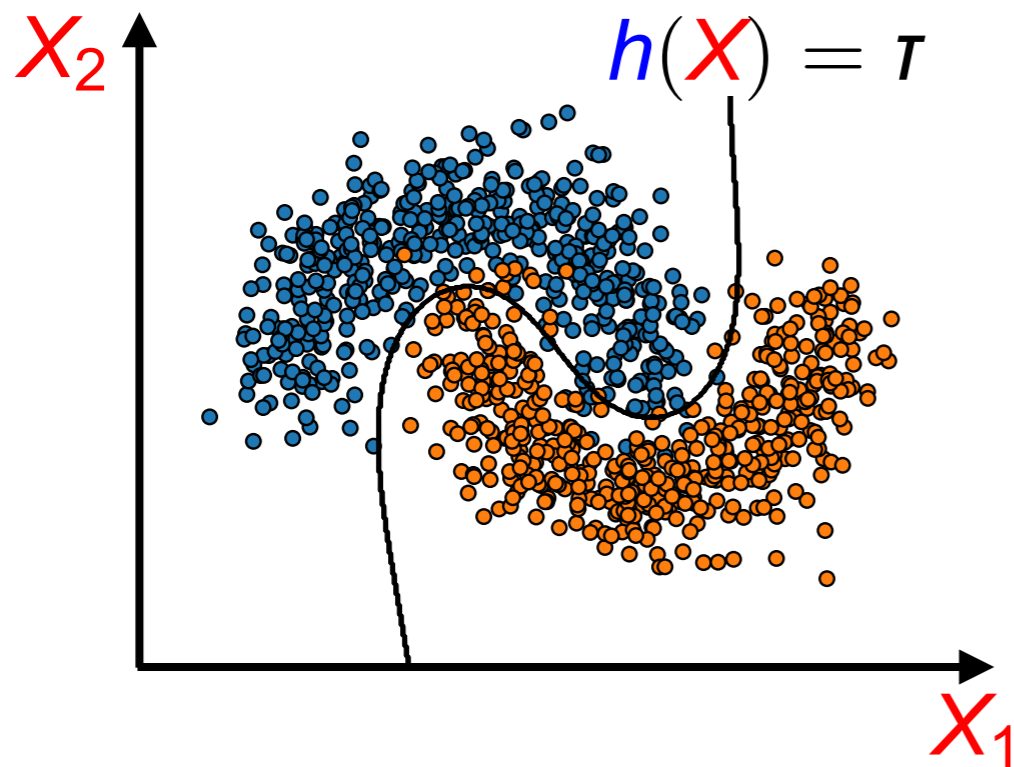
# Metrizing Fairness

Yves Rychener, Bahar Taşkesen, Daniel Kuhn

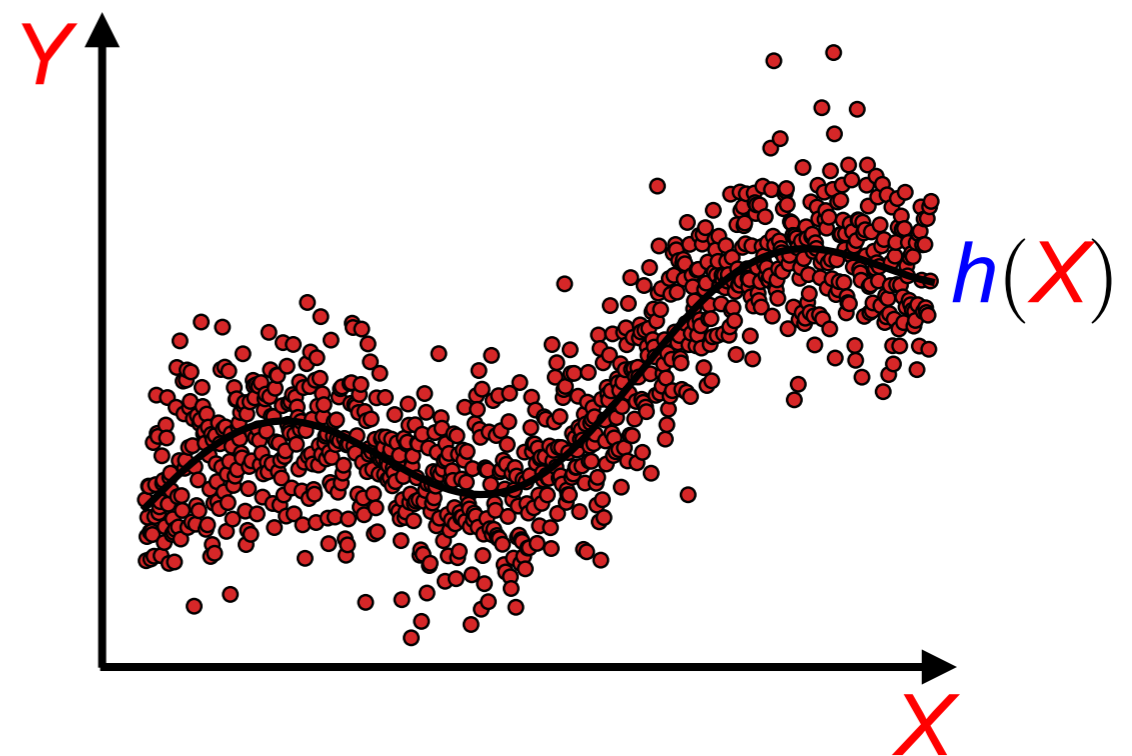
École Polytechnique Fédérale de Lausanne

# Statistical Learning

Learning problem:  $\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)]$



Classification

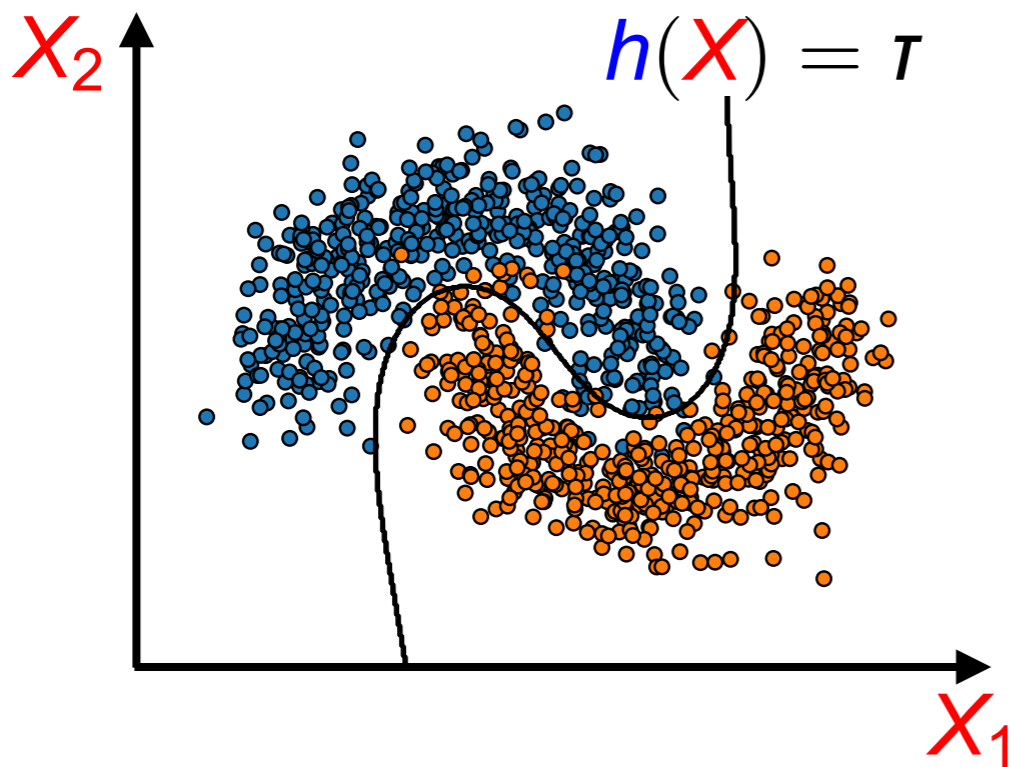


Regression

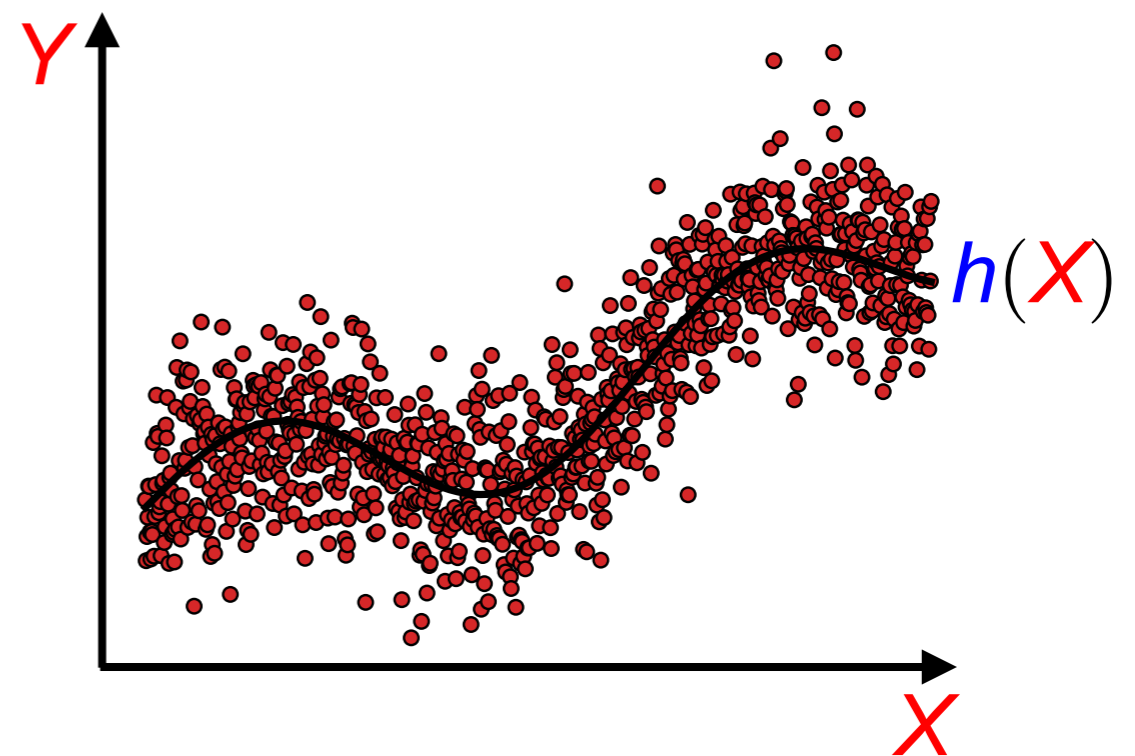
# Statistical Learning

Learning problem:  $\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)]$

$X \in \mathbb{R}^d$  input



Classification



Regression

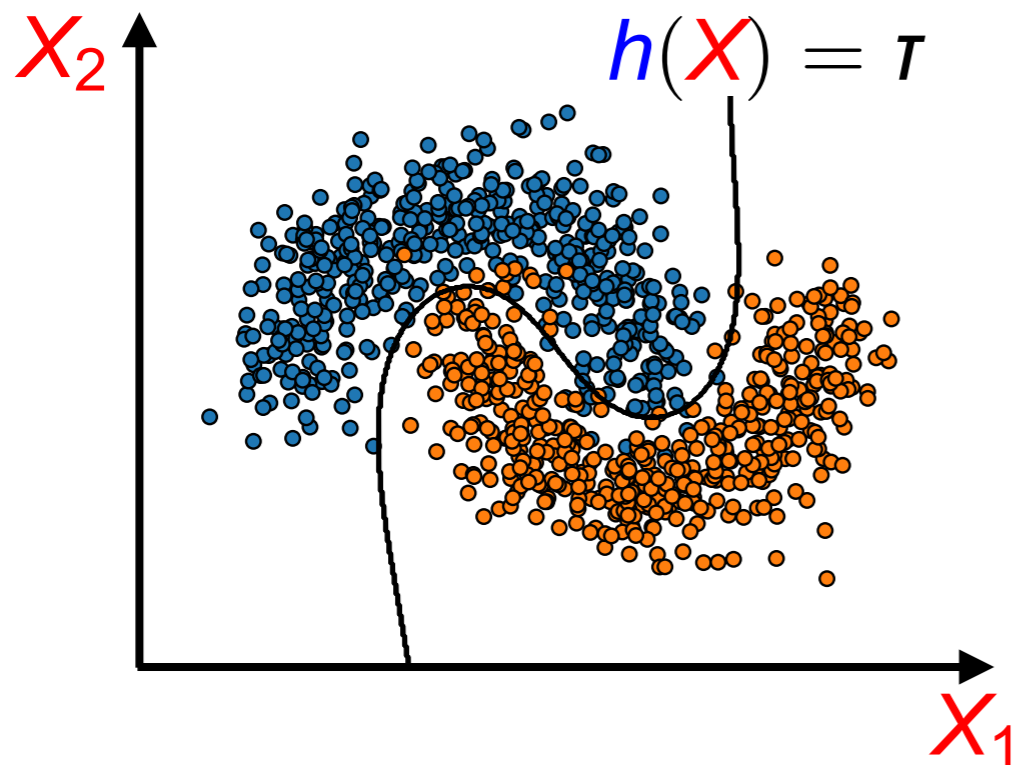
# Statistical Learning

Learning problem:

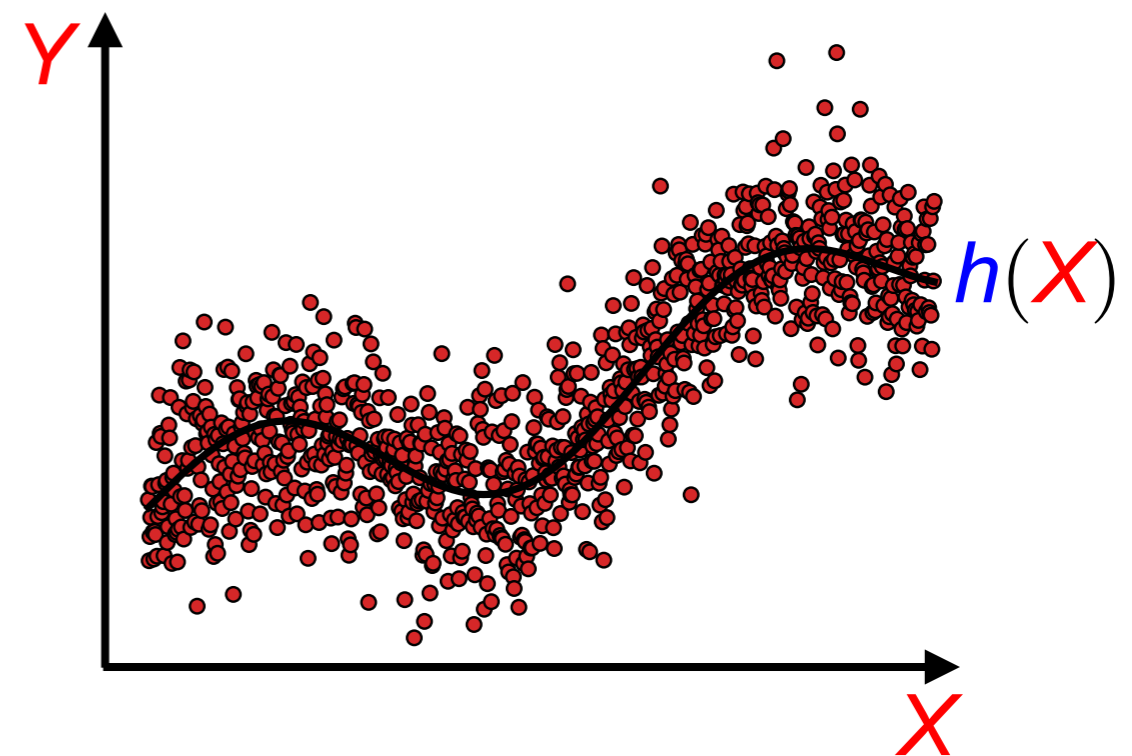
$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)]$$

$Y \in \mathbb{R}$  output

$X \in \mathbb{R}^d$  input



Classification



Regression



# Statistical Learning

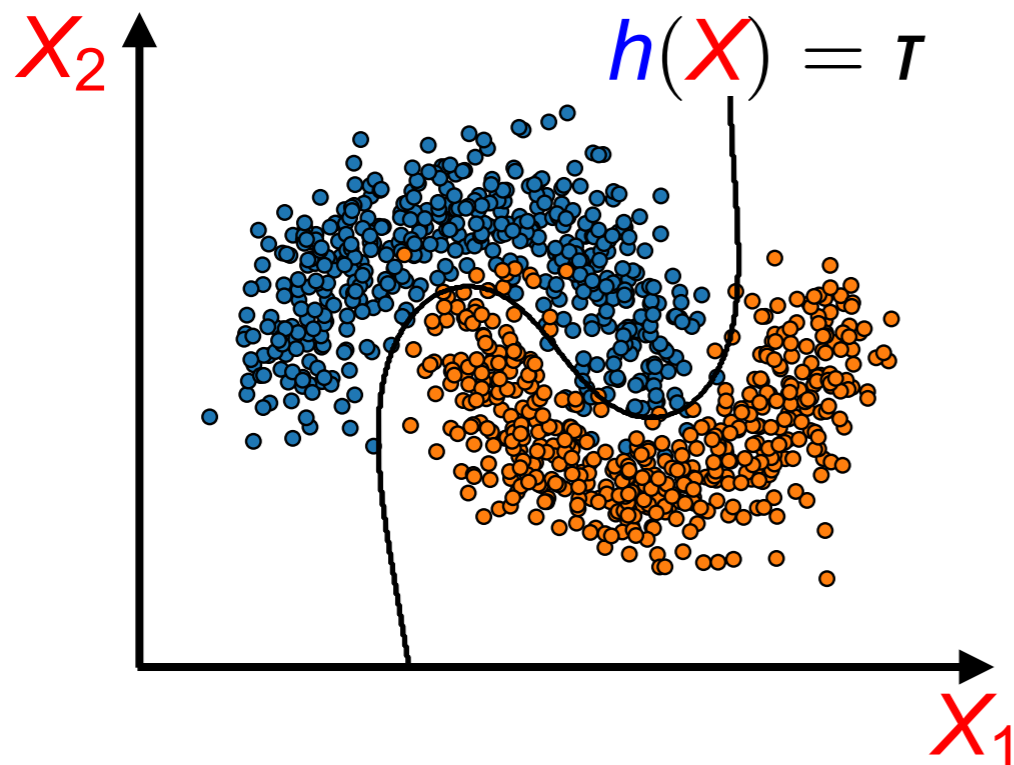
Learning problem:

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)]$$

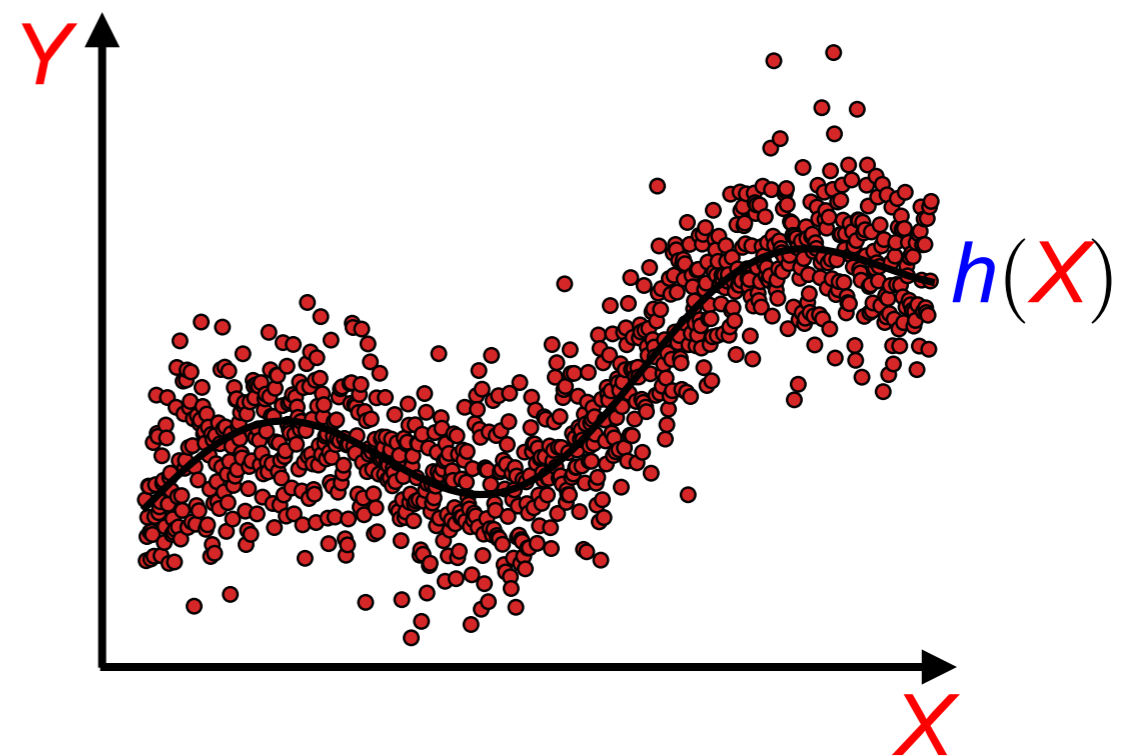
hypothesis  
(predictor)

$X \in \mathbb{R}^d$  input

$Y \in \mathbb{R}$  output



Classification



Regression

# Statistical Learning

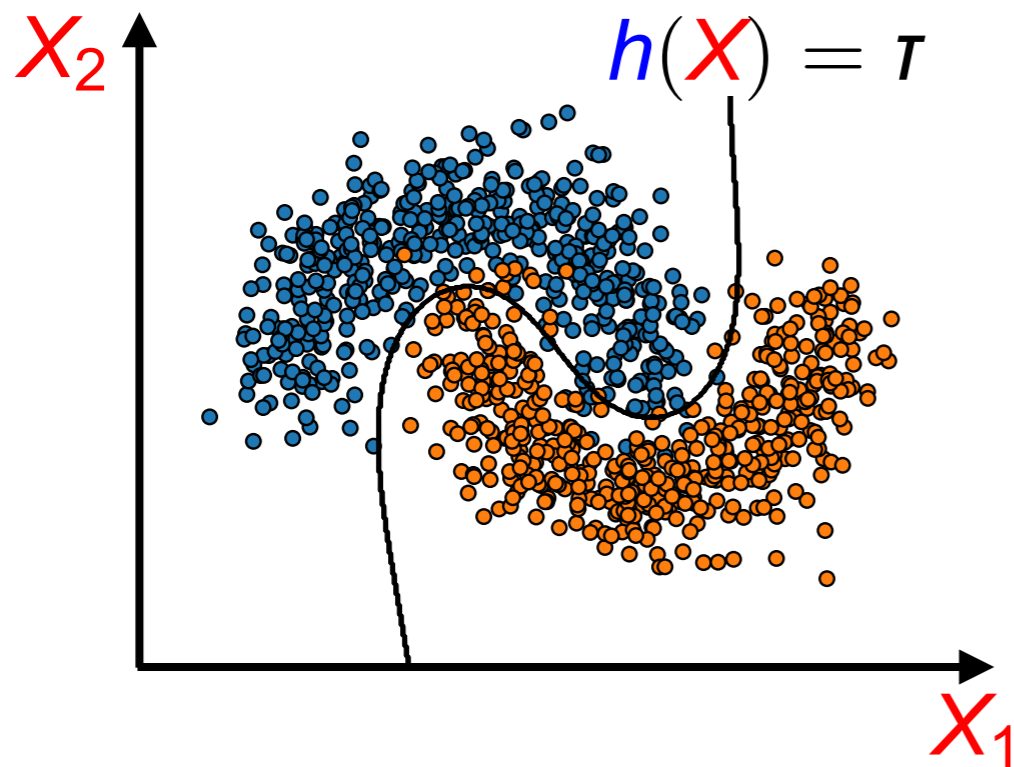
**Learning problem:**  $\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)]$

penalizes dissimilarity of  $h(X)$  and  $Y$

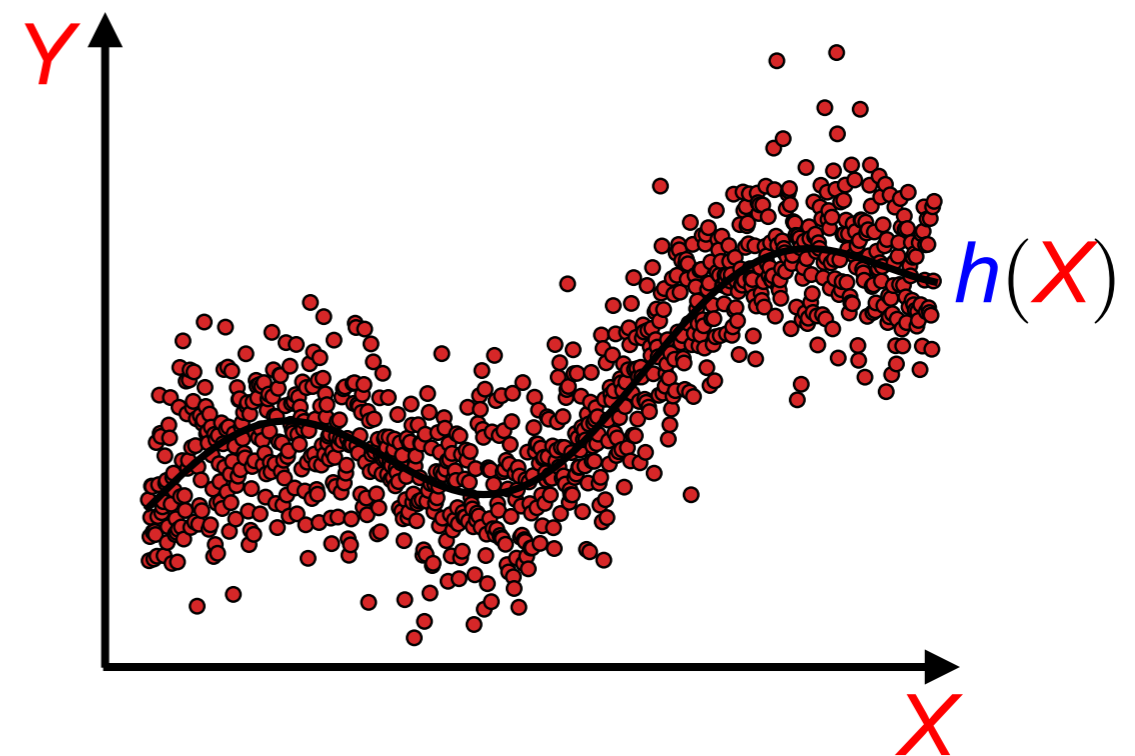
hypothesis (predictor)

$Y \in \mathbb{R}$  output

$X \in \mathbb{R}^d$  input



Classification



Regression

# Weapons of Math Destruction

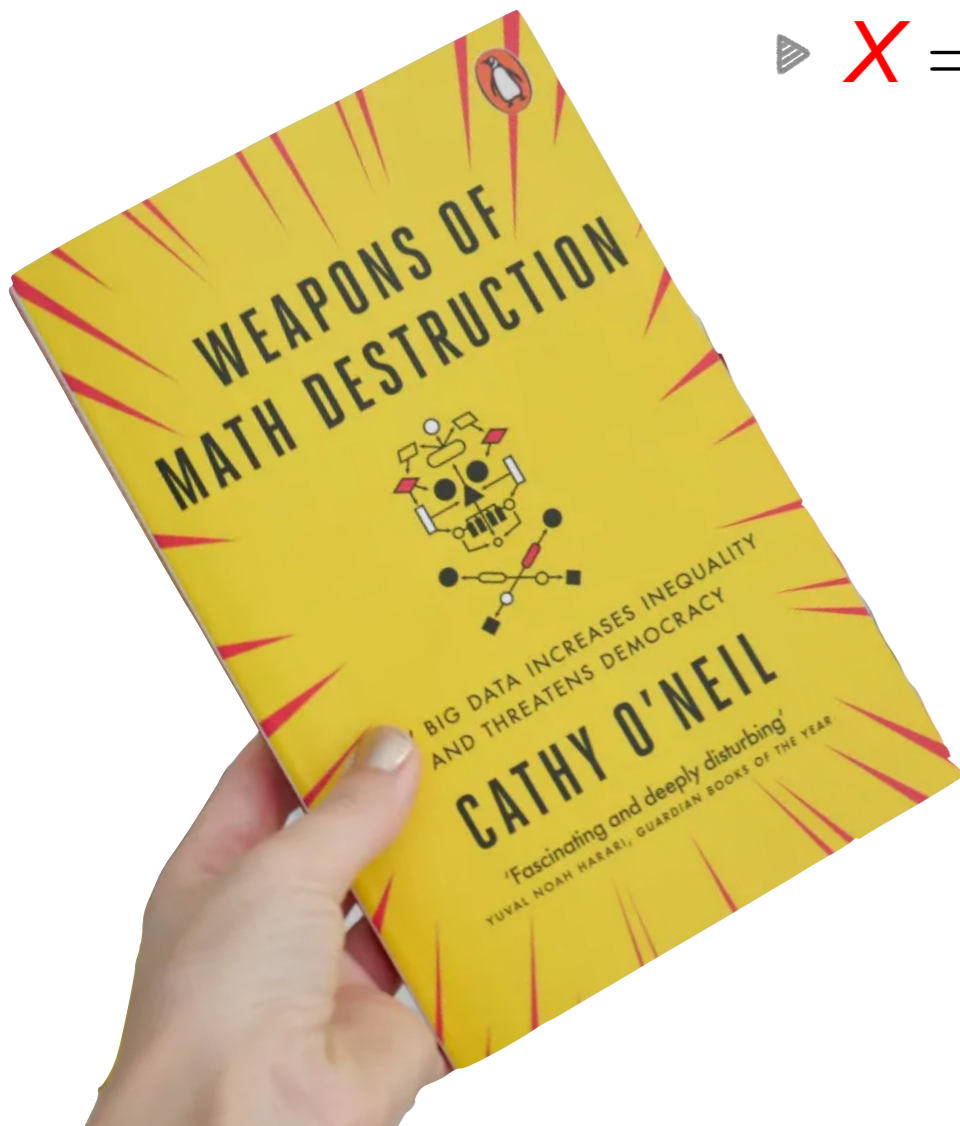
**Learning problem:**  $\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)]$

▶  $X$  = web browsing history ,  $Y$  = consumer behavior

▶  $X$  = credit history,  $Y$  = creditworthiness

▶  $X$  = crime history,  $Y$  = recidivism

▶  $X$  = résumé ,  $Y$  = skills



# Weapons of Math Destruction

## Amazon ditched AI recruiting tool that favored men for technical jobs

Specialists had been building computer programs since 2014 to review résumés in an effort to automate the search process



Amazon's automated hiring tool was found to be inadequate after penalizing the résumés of female candidates. Photograph: Brian Snyder/Reuters

The data on which the AI hiring algorithm was trained created a preference for male candidates.<sup>1)</sup>

## RESEARCH ARTICLE

### ECONOMICS

## Dissecting racial bias in an algorithm used to manage the health of populations

Ziad Obermeyer<sup>1,2\*</sup>, Brian Powers<sup>3</sup>, Christine Vogeli<sup>4</sup>, Sendhil Mullainathan<sup>5\*†</sup>

Health systems rely on commercial prediction algorithms to identify and help patients with complex health needs. We show that a widely used algorithm, typical of this industry-wide approach and affecting millions of patients, exhibits significant racial bias: At a given risk score, Black patients are considerably sicker than White patients, as evidenced by signs of uncontrolled illnesses. Remedying this disparity would increase the percentage of Black patients receiving additional help from 17.7 to 46.5%. The bias arises because the algorithm predicts health care costs rather than illness, but unequal access to care means that we spend less money caring for Black patients than for White patients. Thus, despite health care cost appearing to be an effective proxy for health by some measures of predictive accuracy, large racial biases arise. We suggest that the choice of convenient, seemingly effective proxies for ground truth can be an important source of algorithmic bias in many contexts.

Industry-wide approach affecting millions of patients exhibits significant racial bias.<sup>2)</sup>

1) Dastin, *Reuters*, 2018.

2) Obermeyer et al., *Science*, 2019.



# Sensitive Attributes

**X** contains a sensitive attribute  $A \in \{0, 1\}$  such as:



- ▶ Gender
- ▶ Ethnicity
- ▶ Religion
- ▶ Age
- ▶ Marital Status
- ▶ Citizenship

# Fairness Through Unawareness

**Idea:** Remove  $A$  from  $X$

# Fairness Through Unawareness

**Idea:** Remove  $A$  from  $X$

**Problem:**<sup>1)</sup> Can use other features  $X_1, X_2, X_3, \dots$  to predict  $A$

---

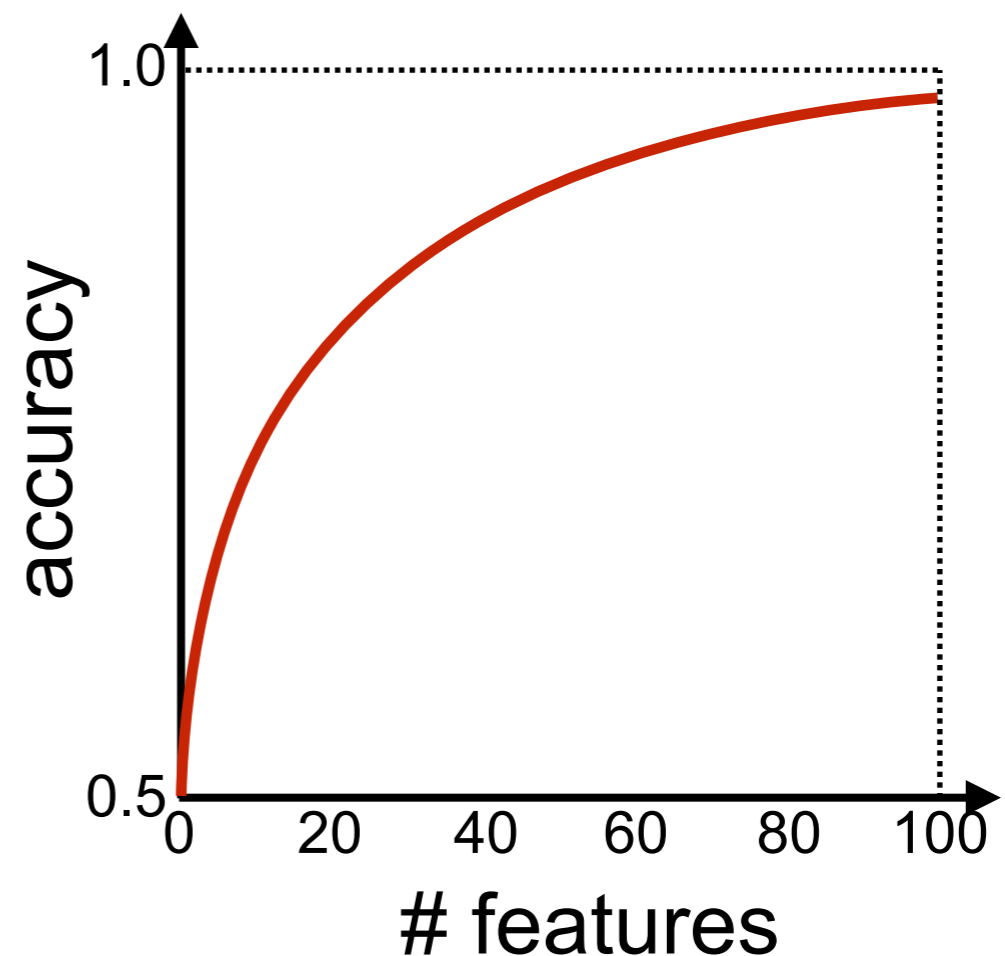
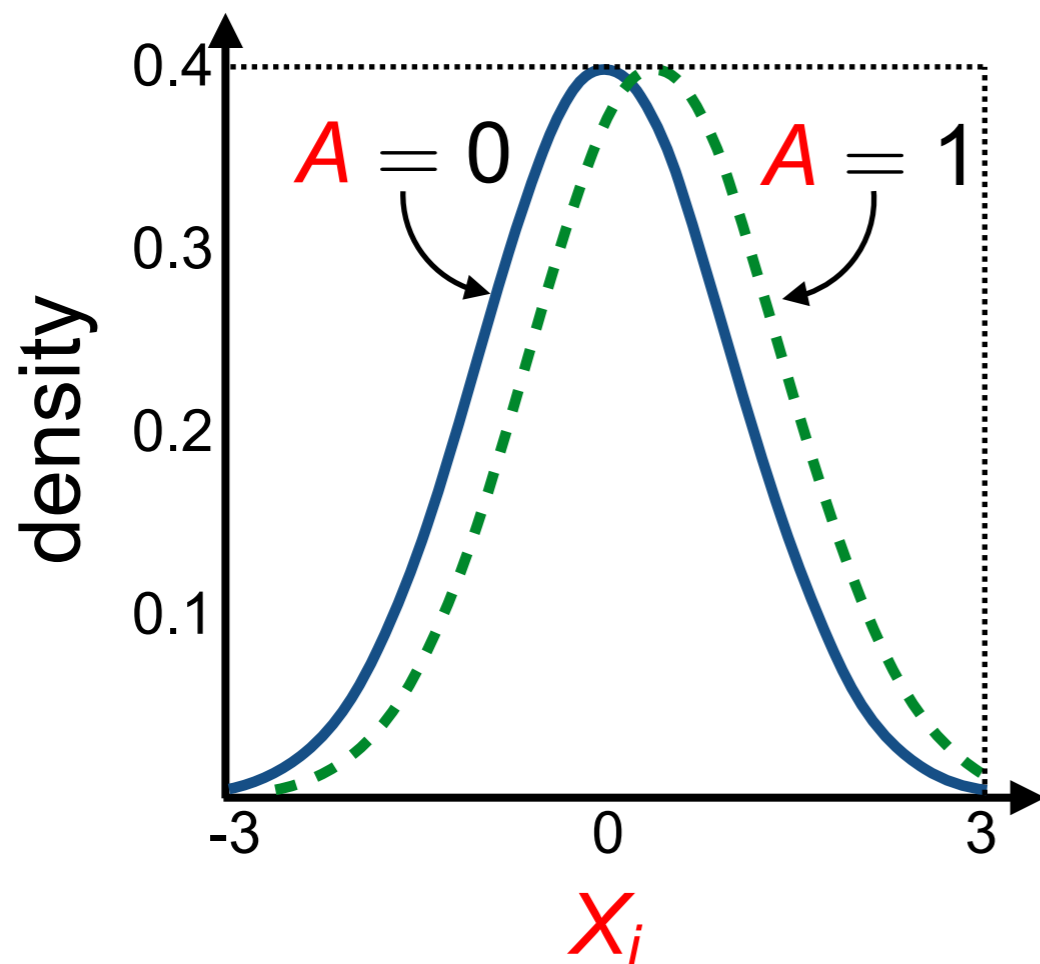
<sup>1)</sup> Barocas, Hardt & Narayanan, *fairmlbook.org*, 2019.



# Fairness Through Unawareness

**Idea:** Remove  $A$  from  $X$

**Problem:**<sup>1)</sup> Can use other features  $X_1, X_2, X_3, \dots$  to predict  $A$



<sup>1)</sup> Barocas, Hardt & Narayanan, *fairmlbook.org*, 2019.

**Statistical parity:**<sup>1)</sup>

$$h(X) | A = 0 \sim h(X) | A = 1$$

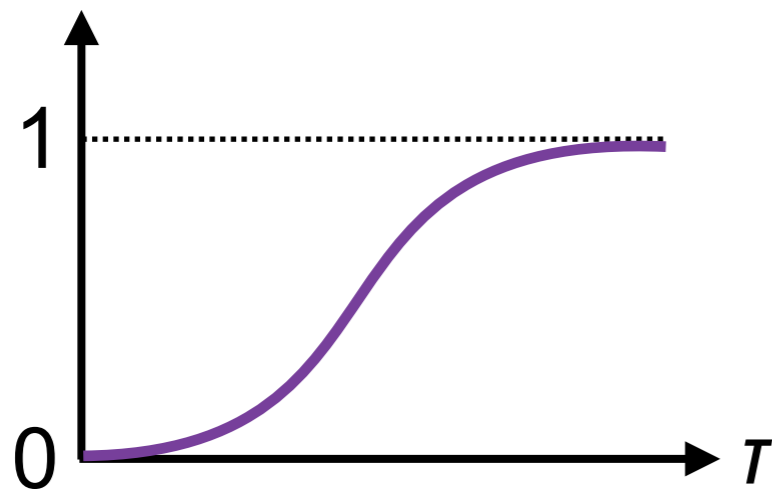
---

<sup>1)</sup> Calders et al., *ICDM*, 2013.

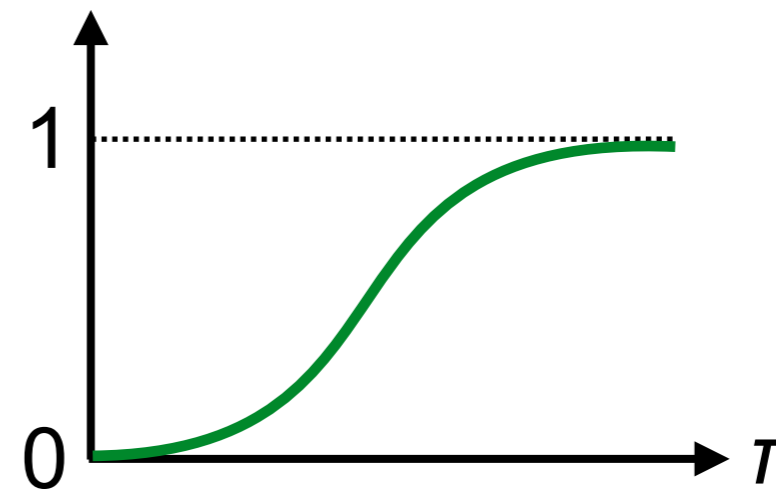
## Statistical parity:<sup>1)</sup>

$$\mathbb{P}(h(X) \leq \tau | A = 0) = \mathbb{P}(h(X) \leq \tau | A = 1) \quad \forall \tau \in \mathbb{R}$$

## CDF of predictor:



$$\mathbb{P}(h(X) \leq \tau | A = 0)$$



$$\mathbb{P}(h(X) \leq \tau | A = 1)$$

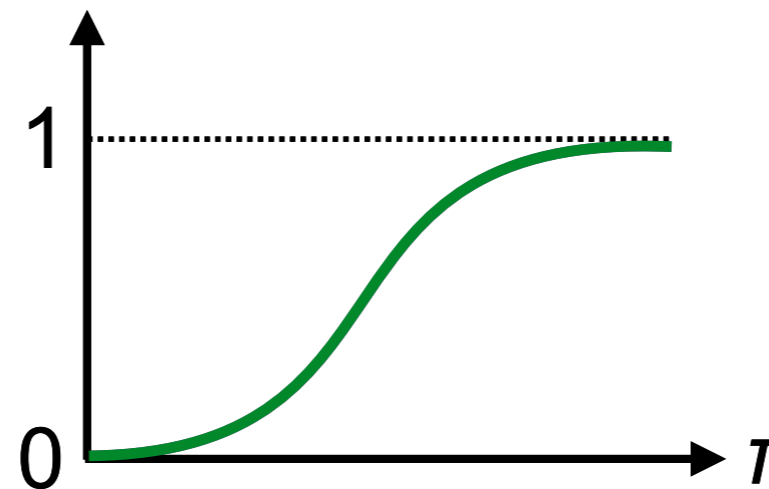
---

<sup>1)</sup> Calders et al., *ICDM*, 2013.

## Statistical parity:<sup>1)</sup>

$$\mathbb{P}(h(X) \leq \tau | A = 0) = \mathbb{P}(h(X) \leq \tau | A = 1) \quad \forall \tau \in \mathbb{R}$$

## CDF of predictor:



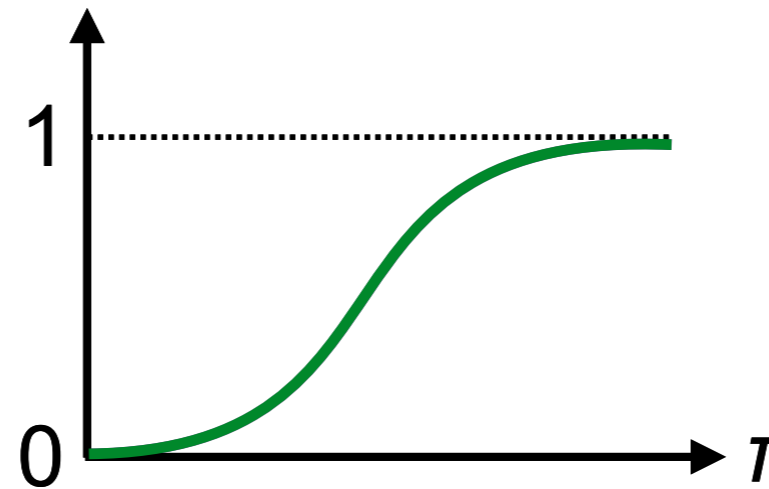
---

<sup>1)</sup> Calders et al., *ICDM*, 2013.

**Statistical parity:**<sup>1)</sup>

$$h(X) \perp A$$

**CDF of predictor:**



---

<sup>1)</sup> Calders et al., *ICDM*, 2013.

# Other Group Fairness Definitions

**Statistical parity:**<sup>1)</sup>

$$h(X) | A = 0 \sim h(X) | A = 1$$

**Equalized odds:**<sup>2)</sup>

$$h(X) | Y = y, A = 0 \sim h(X) | Y = y, A = 1 \quad \forall y \in \mathbb{R}$$

**Risk parity:**<sup>3)</sup>

$$L(h(X), Y) | A = 0 \sim L(h(X), Y) | A = 1$$

---

1) Dwork et al., *ITCS*, 2012.

2) Hardt et al., *NeurIPS*, 2016.

3) Donini et al., *NeurIPS*, 2018.

# Conceptual Analysis of Statistical Parity



# Fair Statistical Learning

**Fair learning problem:**

$$\begin{array}{ll} \min_{h \in \mathcal{H}} & \mathbb{E}[L(h(X), Y)] \\ \text{s.t.} & h(X) \perp A \end{array}$$

# Fair Statistical Learning

**Fair learning problem:**

$$\begin{array}{ll} \min_{h \in \mathcal{H}} & \mathbb{E}[L(h(X), Y)] \\ \text{s.t.} & h(X) \perp A \end{array}$$

SP fairness



# Fair Statistical Learning

**Fair learning problem:**

$$\begin{array}{ll} \min_{h \in \mathcal{H}} & \mathbb{E}[L(h(X), Y)] \\ \text{s.t.} & \boxed{h(X) \perp A} \end{array}$$



SP fairness

$$\iff \mathbb{P}(h(X) \leq \tau | A = 0) = \mathbb{P}(h(X) \leq \tau | A = 1) \quad \forall \tau \in \mathbb{R}$$

$\implies$  reminiscent of chance constraint

$\implies$  intractable

# Idealized Models

## Assumptions:

- ▶  $\mathbb{P}$  is known (sample size =  $\infty$ )
- ▶  $\mathcal{H} = \mathcal{L}(\mathbb{R}^d, \mathbb{R})$  (all measurable hypotheses)
- ▶  $h^*$  is essentially unique

# Automatic Salary Determination

**Goal:** Predict the skill levels of job candidates.

$X_1$  = GPA (normalized to  $[0, 1]$ )  
 $X_2$  = age group (0: age  $> 40$ , 1: age  $\leq 40$ )  
 $Y$  = skill level (normalized to  $[0, 1]$ )  
 $S$  = work experience (normalized to  $[0, 1]$ , unobserved)

# Automatic Salary Determination

**Goal:** Predict the skill levels of job candidates.

$X_1$  = GPA (normalized to  $[0, 1]$ )  
 $X_2$  = age group (0: age  $> 40$ , 1: age  $\leq 40$ )  
 $Y$  = skill level (normalized to  $[0, 1]$ )  
 $S$  = work experience (normalized to  $[0, 1]$ , unobserved)

$$X_1, S \sim \mathcal{U}([0, 1]), \quad X_2 \sim \mathcal{U}(\{0, 1\}), \quad Y = X_1 \cdot X_2 + S \cdot (1 - X_2)$$

# Automatic Salary Determination

**Goal:** Predict the skill levels of job candidates.

$X_1$  = GPA (normalized to  $[0, 1]$ )  
 $X_2$  = age group (0: age  $> 40$ , 1: age  $\leq 40$ )  
 $Y$  = skill level (normalized to  $[0, 1]$ )  
 $S$  = work experience (normalized to  $[0, 1]$ , unobserved)

$$X_1, S \sim \mathcal{U}([0, 1]), \quad X_2 \sim \mathcal{U}(\{0, 1\}), \quad Y = X_1 \cdot X_2 + S \cdot (1 - X_2)$$

$A = X_2$  = sensitive attribute



# Automatic Salary Determination

Learning problem with square loss  $L(\hat{y}, y) = (\hat{y} - y)^2$  :

**Original learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[(h(X) - Y)^2]$$

**Fair learning problem:**

$$\begin{aligned} \min_{h \in \mathcal{H}} \quad & \mathbb{E}[(h(X) - Y)^2] \\ \text{s.t.} \quad & h(X) \perp A \end{aligned}$$

# Automatic Salary Determination

Learning problem with square loss  $L(\hat{y}, y) = (\hat{y} - y)^2$  :

**Original learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[(h(X) - Y)^2]$$

**Fair learning problem:**

$$\begin{aligned} \min_{h \in \mathcal{H}} \quad & \mathbb{E}[(h(X) - Y)^2] \\ \text{s.t.} \quad & h(X) \perp A \end{aligned}$$

$$\implies h^*(X) = \begin{cases} \frac{1}{2} & \text{if } X_2 = 0 \\ X_1 & \text{if } X_2 = 1 \end{cases}$$

# Automatic Salary Determination

Learning problem with square loss  $L(\hat{y}, y) = (\hat{y} - y)^2$  :

**Original learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[(h(X) - Y)^2]$$

$$\implies h^*(X) = \begin{cases} \frac{1}{2} & \text{if } X_2 = 0 \\ X_1 & \text{if } X_2 = 1 \end{cases}$$

**Fair learning problem:**

$$\begin{aligned} \min_{h \in \mathcal{H}} \quad & \mathbb{E}[(h(X) - Y)^2] \\ \text{s.t.} \quad & h(X) \perp A \end{aligned}$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(X_1 - \frac{1}{2})$$

# Automatic Salary Determination

Learning problem with square loss  $L(\hat{y}, y) = (\hat{y} - y)^2$  :

**Original learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[(h(X) - Y)^2]$$

$$\implies h^*(X) = \begin{cases} \frac{1}{2} & \text{if } X_2 = 0 \\ X_1 & \text{if } X_2 = 1 \end{cases}$$

**Fair learning problem:**

$$\begin{aligned} \min_{h \in \mathcal{H}} \quad & \mathbb{E}[(h(X) - Y)^2] \\ \text{s.t.} \quad & h(X) \perp A \end{aligned}$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(X_1 - \frac{1}{2})$$

$$\implies \mathbb{E}[(h^*(X) - Y)^2] = \frac{1}{24} < \frac{91}{96} = \mathbb{E}[(h_{\text{SP}}^*(X) - Y)^2]$$

# Automatic Salary Determination

Learning problem with square loss  $L(\hat{y}, y) = (\hat{y} - y)^2$  :

**Original learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[(h(X) - Y)^2]$$

**Fair learning problem:**

$$\begin{aligned} \min_{h \in \mathcal{H}} \quad & \mathbb{E}[(h(X) - Y)^2] \\ \text{s.t.} \quad & h(X) \perp A \end{aligned}$$

$$\implies h^*(X) = \begin{cases} \frac{1}{2} & \text{if } X_2 = 0 \\ X_1 & \text{if } X_2 = 1 \end{cases}$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(X_1 - \frac{1}{2})$$

**SP deteriorates predictive power!**

# Automatic Salary Determination

**Original learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[(h(X) - Y)^2]$$

$$\implies h^*(X) = \begin{cases} \frac{1}{2} & \text{if } X_2 = 0 \\ X_1 & \text{if } X_2 = 1 \end{cases}$$

$A = X_2 = 1$  (junior candidate) :

$$\implies h^*(X) = Y$$

**Fair learning problem:**

$$\begin{aligned} \min_{h \in \mathcal{H}} \quad & \mathbb{E}[(h(X) - Y)^2] \\ \text{s.t.} \quad & h(X) \perp A \end{aligned}$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(X_1 - \frac{1}{2})$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(Y - \frac{1}{2})$$

# Automatic Salary Determination

**Original learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[(h(X) - Y)^2]$$

$$\implies h^*(X) = \begin{cases} \frac{1}{2} & \text{if } X_2 = 0 \\ X_1 & \text{if } X_2 = 1 \end{cases}$$

**Fair learning problem:**

$$\begin{aligned} \min_{h \in \mathcal{H}} \quad & \mathbb{E}[(h(X) - Y)^2] \\ \text{s.t.} \quad & h(X) \perp A \end{aligned}$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(X_1 - \frac{1}{2})$$

$A = X_2 = 1$  (junior candidate) :

$$\implies h^*(X) = Y$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(Y - \frac{1}{2})$$

salary grows with skill level



# Automatic Salary Determination

**Original learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[(h(X) - Y)^2]$$

$$\implies h^*(X) = \begin{cases} \frac{1}{2} & \text{if } X_2 = 0 \\ X_1 & \text{if } X_2 = 1 \end{cases}$$

$A = X_2 = 0$  (senior candidate) :

$$\implies h^*(X) = \frac{1}{2}$$

**Fair learning problem:**

$$\begin{aligned} \min_{h \in \mathcal{H}} \quad & \mathbb{E}[(h(X) - Y)^2] \\ \text{s.t.} \quad & h(X) \perp A \end{aligned}$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(X_1 - \frac{1}{2})$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(X_1 - \frac{1}{2})$$


# Automatic Salary Determination

**Original learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[(h(X) - Y)^2]$$

$$\implies h^*(X) = \begin{cases} \frac{1}{2} & \text{if } X_2 = 0 \\ X_1 & \text{if } X_2 = 1 \end{cases}$$

$A = X_2 = 0$  (senior candidate) :


$$\implies h^*(X) = \frac{1}{2}$$


uniform salary

**Fair learning problem:**

$$\begin{aligned} \min_{h \in \mathcal{H}} \quad & \mathbb{E}[(h(X) - Y)^2] \\ \text{s.t.} \quad & h(X) \perp A \end{aligned}$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(X_1 - \frac{1}{2})$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(X_1 - \frac{1}{2})$$


random salary

# Automatic Salary Determination

**Original learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[(h(X) - Y)^2]$$

**Fair learning problem:**

$$\begin{aligned} \min_{h \in \mathcal{H}} \quad & \mathbb{E}[(h(X) - Y)^2] \\ \text{s.t.} \quad & h(X) \perp A \end{aligned}$$

$$\implies h^*(X) = \begin{cases} \frac{1}{2} & \text{if } X_2 = 0 \\ X_1 & \text{if } X_2 = 1 \end{cases}$$

$$\implies h_{\text{SP}}^*(X) = \frac{1}{2} + \frac{1}{2}(X_1 - \frac{1}{2})$$

$A = X_2 = 0$  (senior candidate) :

Treatment of senior candidates seems less "fair" when SP is enforced!

random salary

# Optimality Implies Statistical Parity

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)]$$

**Theorem.**  $\mathbb{P}_{Y|X} \perp A \implies h^*(X) \perp A$

# Optimality Implies Statistical Parity

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)]$$

**Theorem.**  $\mathbb{P}_{Y|X} \perp A \implies h^*(X) \perp A$

$\implies$  SP is a necessary optimality condition!

# Training with Biased Data

**True learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y_0)]$$

**Biased learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y_\delta)]$$

# Training with Biased Data

**True learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y_0)]$$



true target  
(no data available)

**Biased learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y_\delta)]$$



biased target  
(data available)

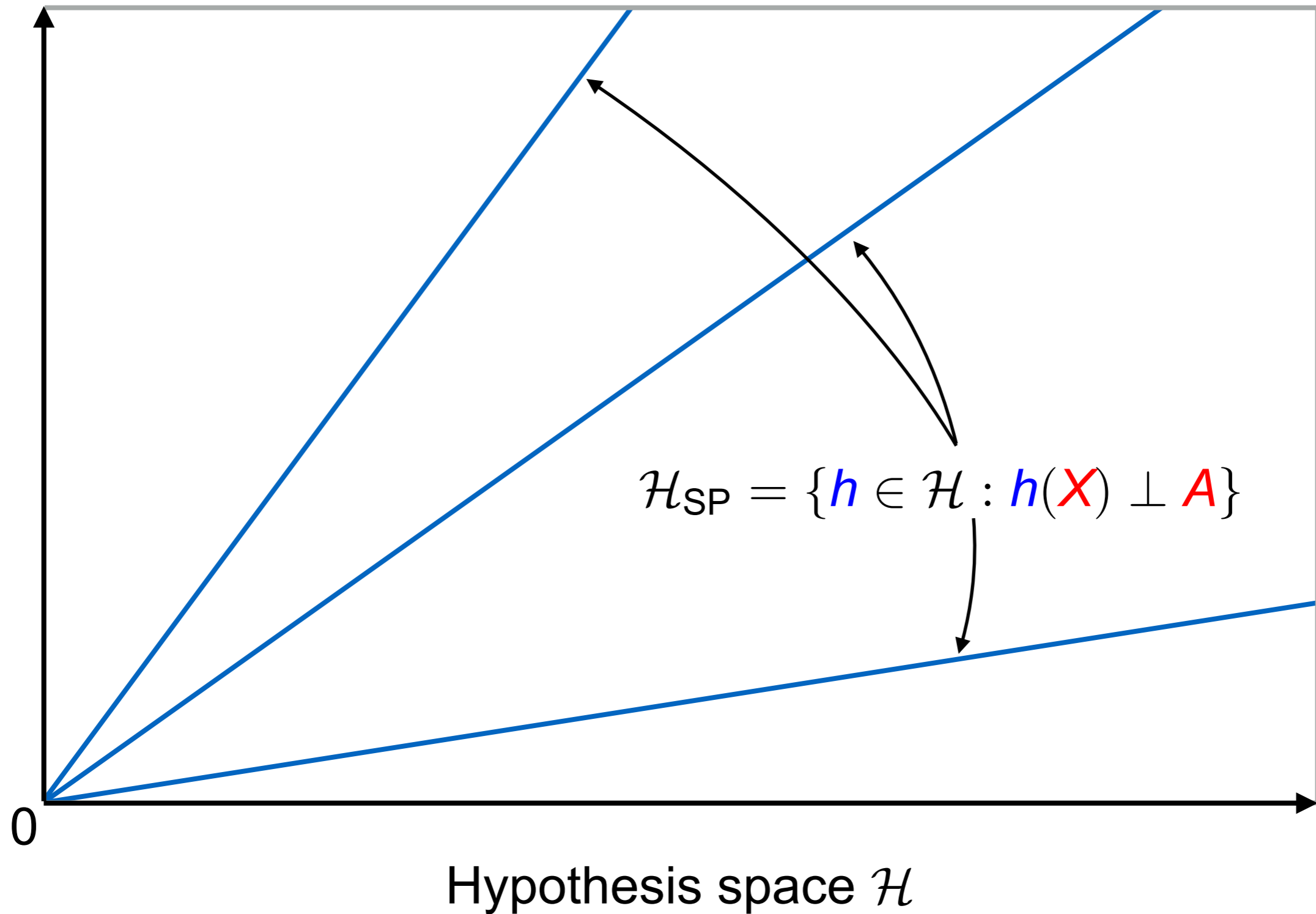
# The Geometry of Statistical Parity



Hypothesis space  $\mathcal{H} = \mathcal{L}(\mathbb{R}^d, \mathbb{R})$

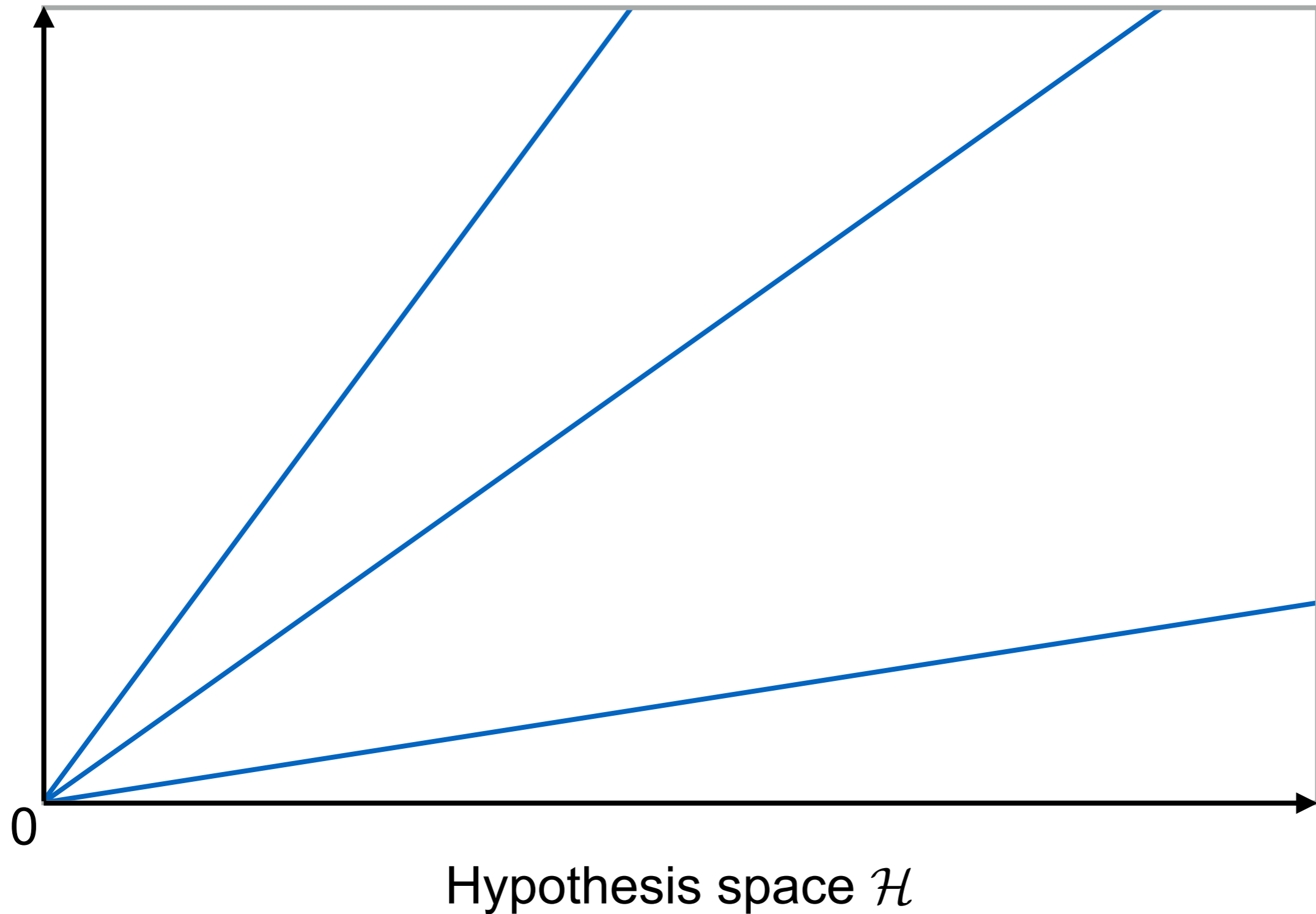


# The Geometry of Statistical Parity



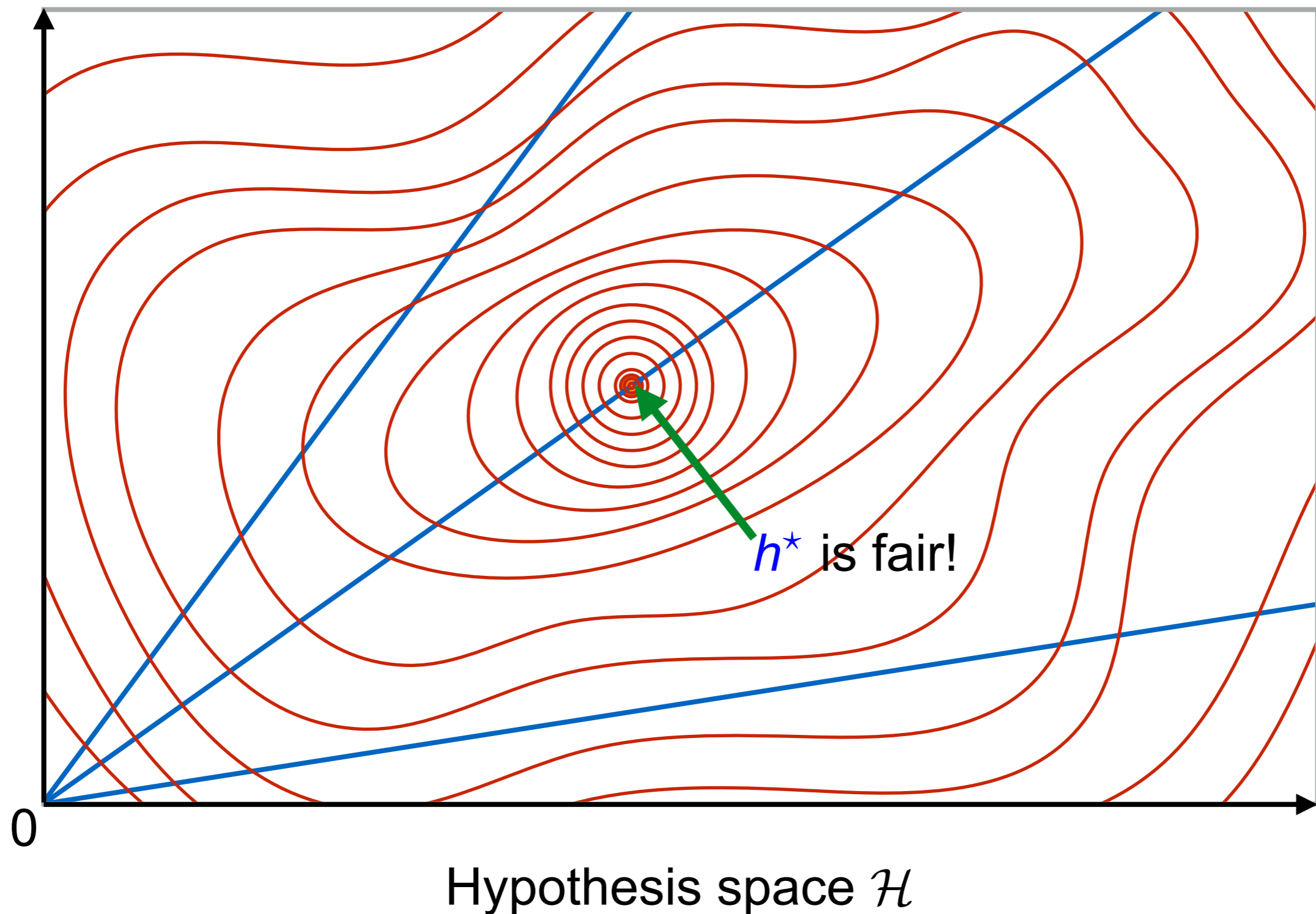
# The Geometry of Statistical Parity

**Assumption:**  $\mathbb{P}_{Y_0|X} \perp A$



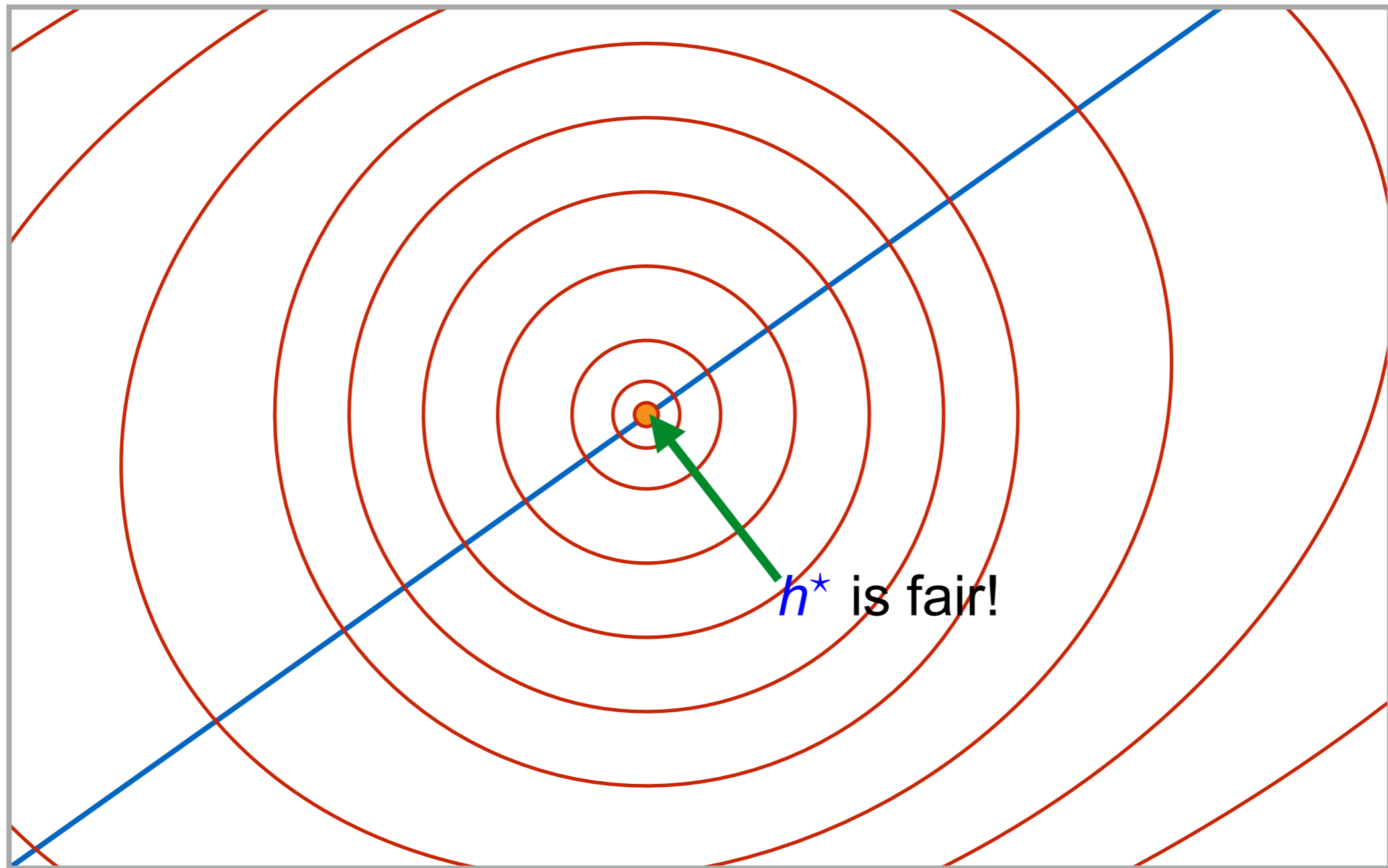
# The Geometry of Statistical Parity

$$\text{Prediction loss} = f(h) = \mathbb{E}[L(h(X, Y_0))]$$



# The Geometry of Statistical Parity

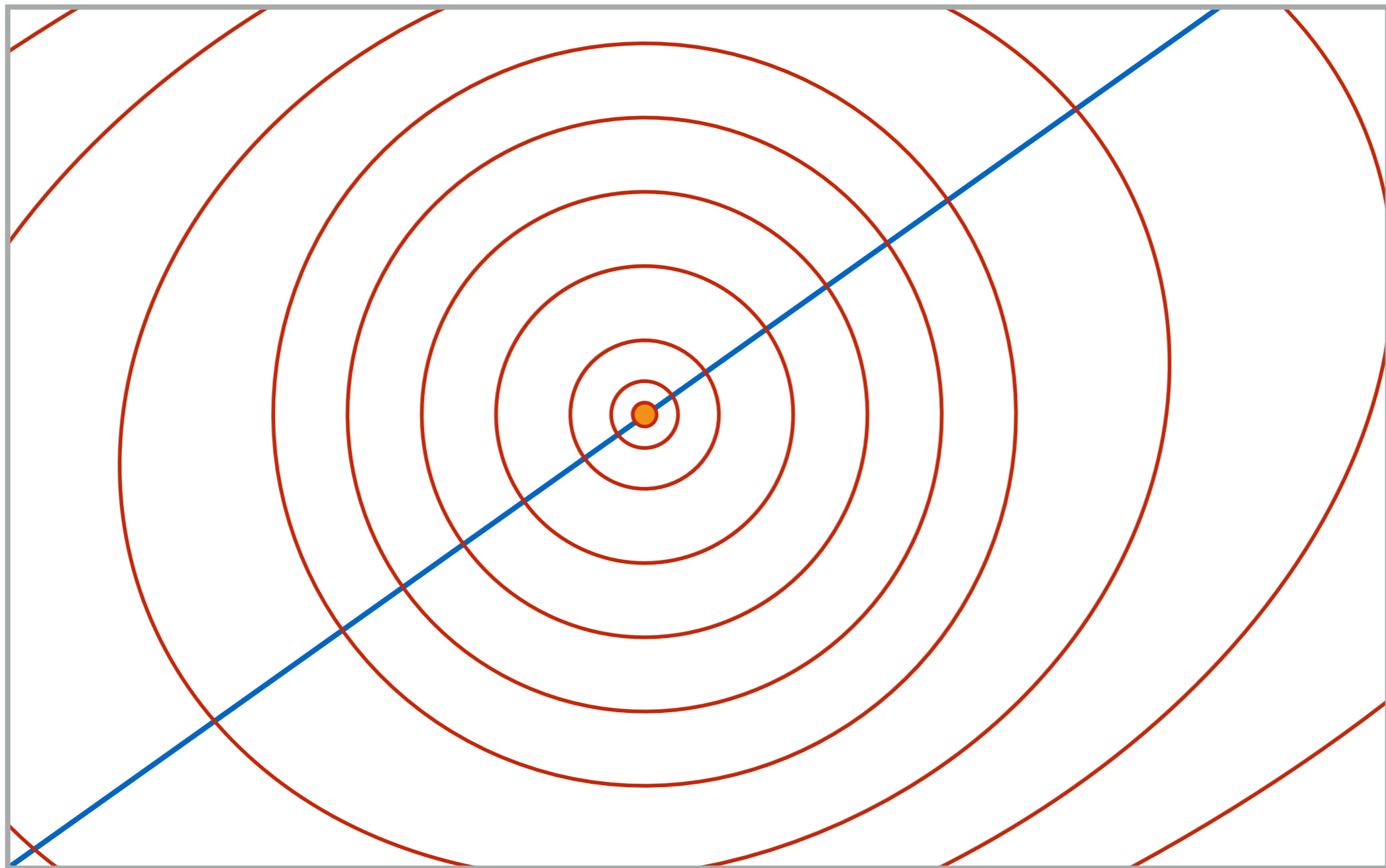
$$\text{Prediction loss} = f(h) = \mathbb{E}[L(h(X, Y_0))]$$



Hypothesis space  $\mathcal{H}$

# The Geometry of Statistical Parity

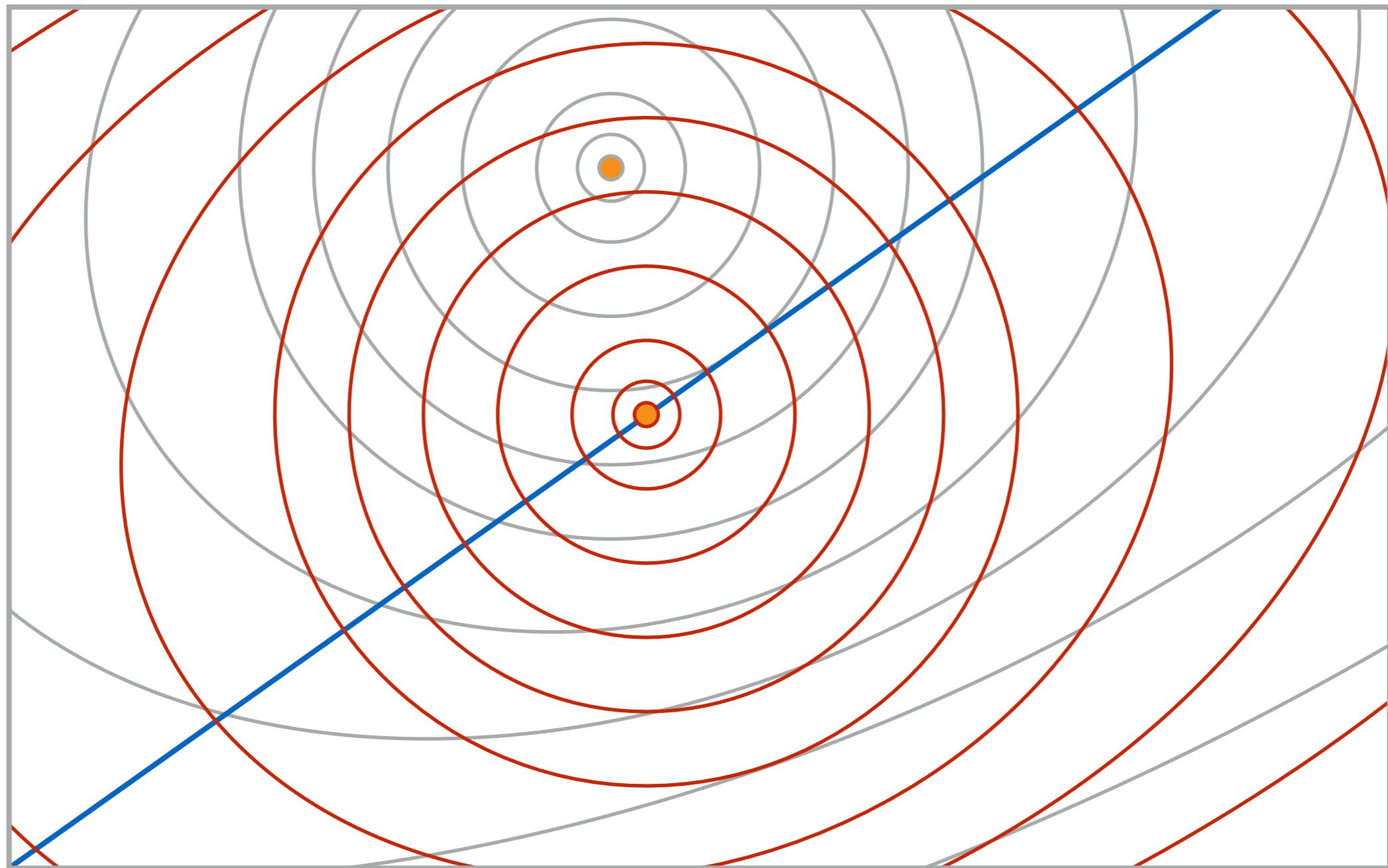
Biased prediction loss =  $f_\delta(h) = \mathbb{E}[L(h(X, Y_\delta))]$



Hypothesis space  $\mathcal{H}$

# The Geometry of Statistical Parity

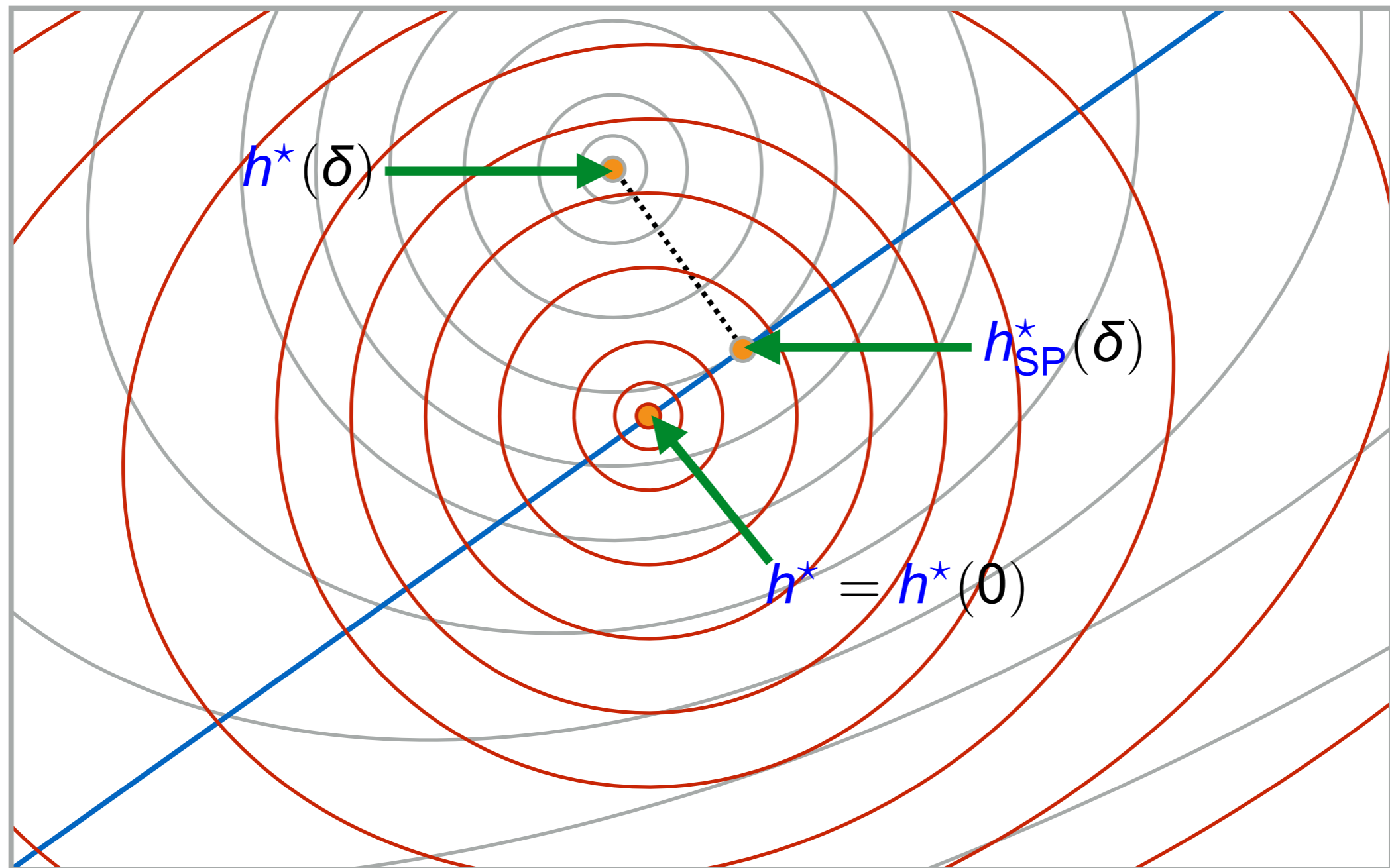
— Contours of  $= f_0(h)$ , — Contours of  $= f_\delta(h)$



Hypothesis space  $\mathcal{H}$

# The Geometry of Statistical Parity

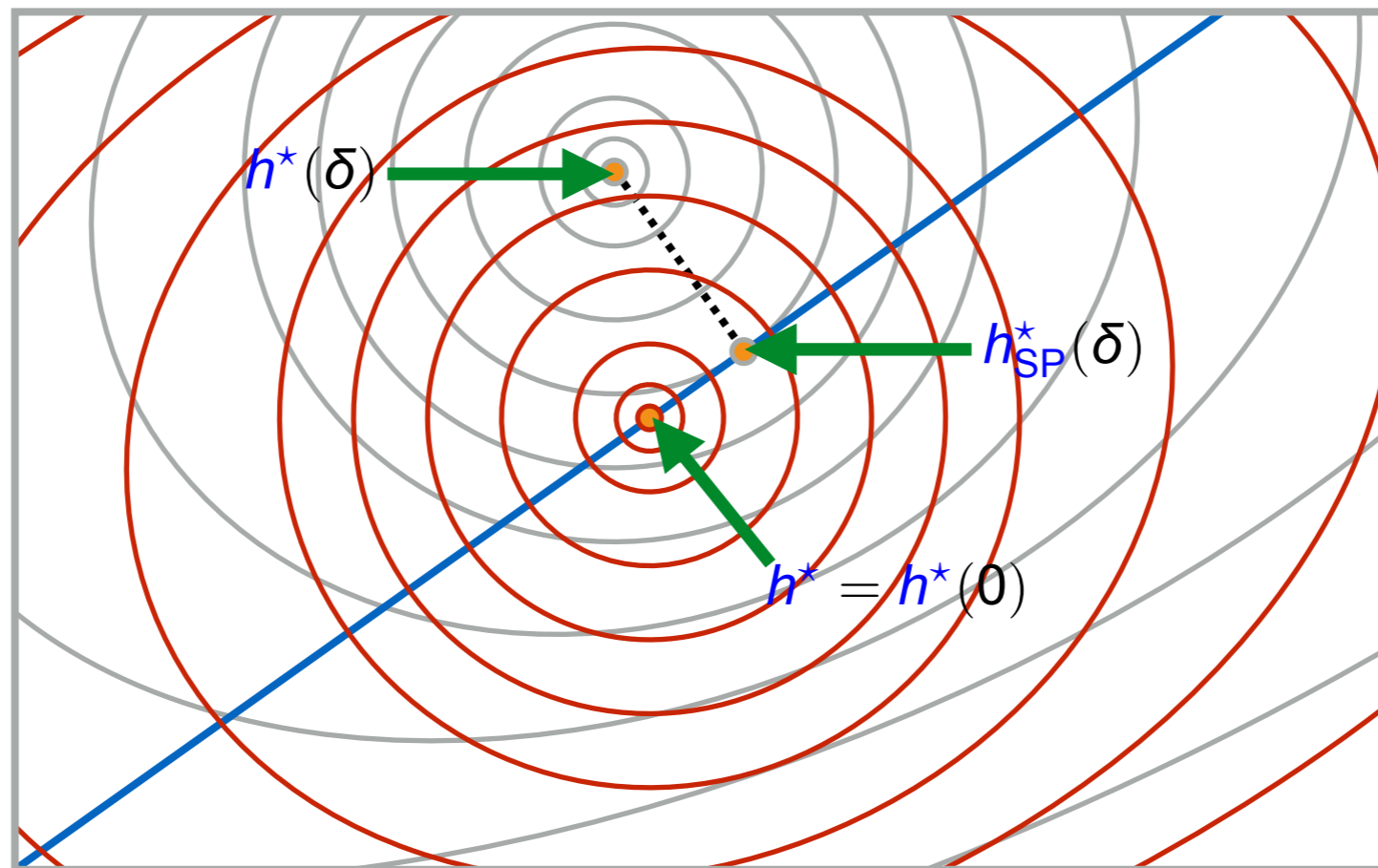
— Contours of  $= f_0(h)$ , — Contours of  $= f_\delta(h)$



Hypothesis space  $\mathcal{H}$

# The Geometry of Statistical Parity

**Theorem:** If  $\mathbb{P}_{Y_0|X} \perp A$  and  $\delta$  is small, then  $h_{SP}^*(\delta)$  is preferable to  $h^*(\delta)$  w.r.t. the true objective  $f_0(h) = \mathbb{E}[L(h(X), Y_0)]$ .





# The Geometry of Statistical Parity

**Theorem:** If  $\mathbb{P}_{Y_0|X} \perp A$  and  $\delta$  is small, then  $h_{SP}^*(\delta)$  is preferable to  $h^*(\delta)$  w.r.t. the true objective  $f_0(h) = \mathbb{E}[L(h(X, Y_0))]$ .

**Win-win situation:**  
Statistical parity improves both  
fairness and predictive power!

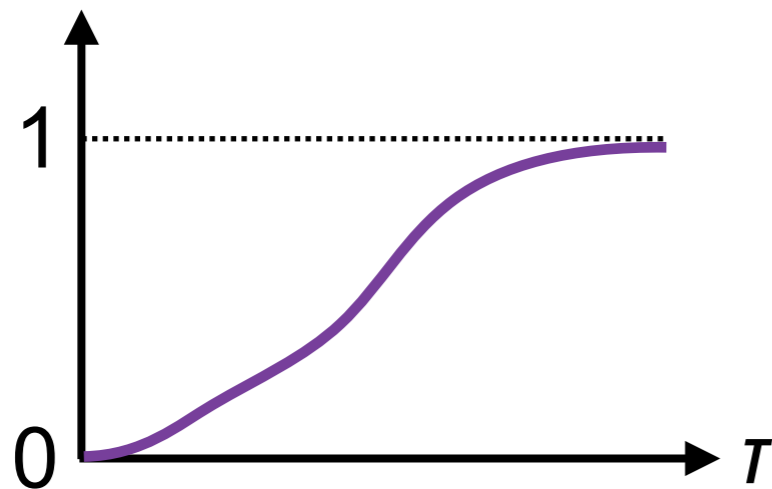
# Unfairness Measures and Integral Probability Metrics

# Relaxing Statistical Parity

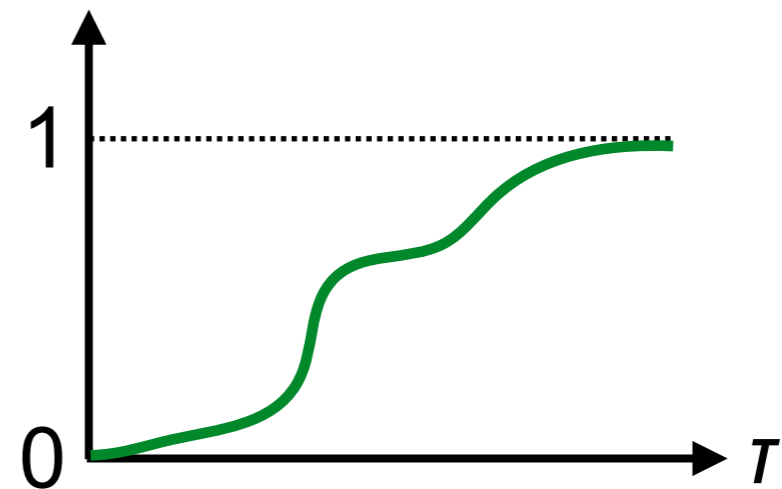
**Statistical parity at level  $\varepsilon$ :**

$$|\mathbb{P}(h(X) \leq \tau | A = 0) - \mathbb{P}(h(X) \leq \tau | A = 1)| \leq \varepsilon \quad \forall \tau \in \mathbb{R}$$

**CDF of predictor:**



$$\mathbb{P}(h(X) \leq \tau | A = 0)$$



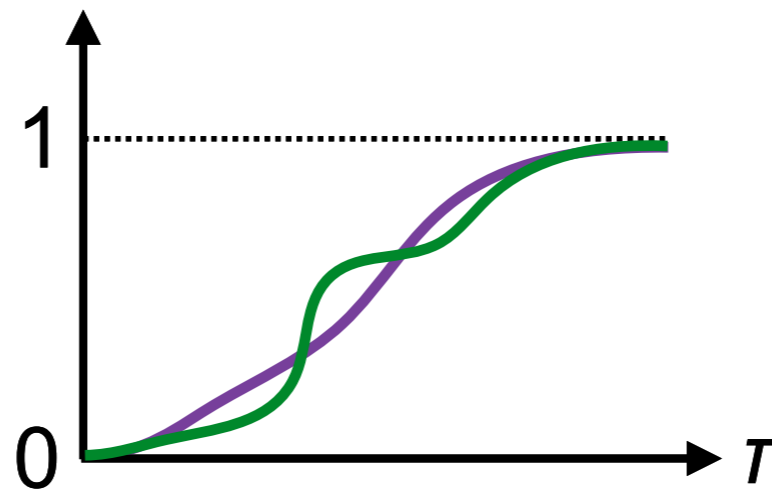
$$\mathbb{P}(h(X) \leq \tau | A = 1)$$

# Relaxing Statistical Parity

**Statistical parity at level  $\varepsilon$ :**

$$|\mathbb{P}(h(X) \leq \tau | A = 0) - \mathbb{P}(h(X) \leq \tau | A = 1)| \leq \varepsilon \quad \forall \tau \in \mathbb{R}$$

**CDF of predictor:**

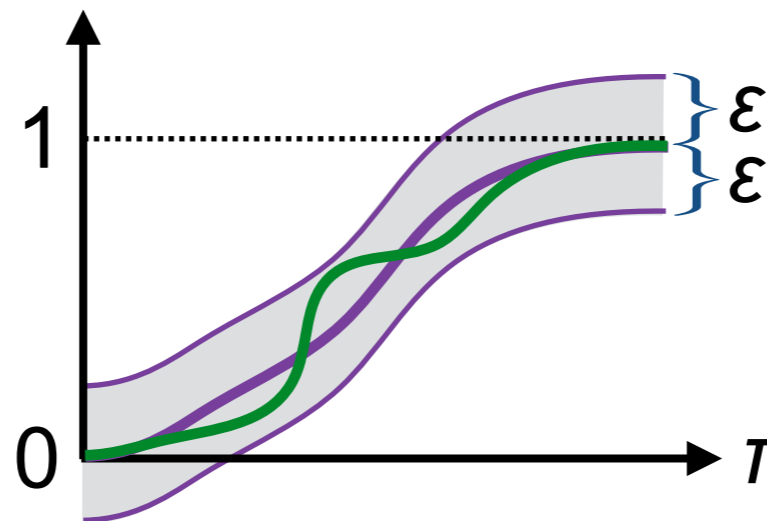


# Relaxing Statistical Parity

**Statistical parity at level  $\varepsilon$ :**

$$|\mathbb{P}(h(X) \leq \tau | A = 0) - \mathbb{P}(h(X) \leq \tau | A = 1)| \leq \varepsilon \quad \forall \tau \in \mathbb{R}$$

**CDF of predictor:**



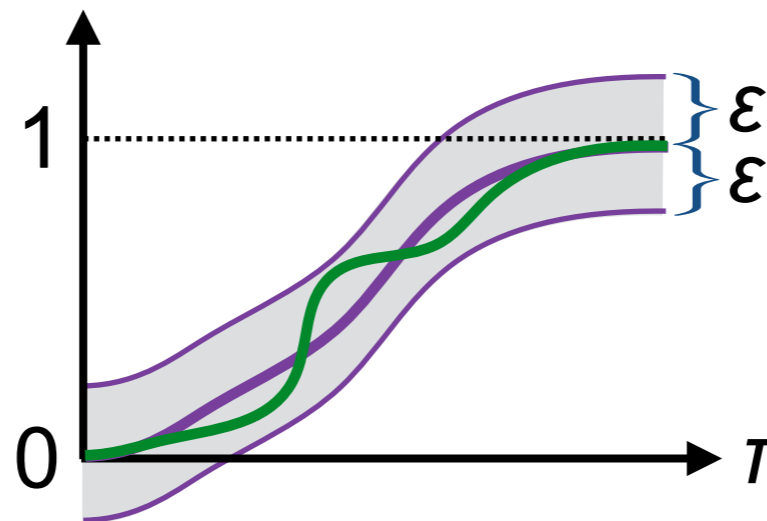
# Relaxing Statistical Parity

**Statistical parity at level  $\varepsilon$ :**

$$\mathcal{D} \left( \mathbb{P}_{h(X)|A=0}, \mathbb{P}_{h(X)|A=1} \right) \leq \varepsilon$$

↑  
Kolmogorov distance

**CDF of predictor:**




# Integral Probability Metrics (IPMs)

$$\mathcal{D}_\Psi(Q_1, Q_2) = \sup_{\psi \in \Psi} \int_{\mathbb{R}} \psi(y) Q_1(\mathrm{d}y) - \int_{\mathbb{R}} \psi(y) Q_2(\mathrm{d}y)$$

# Integral Probability Metrics (IPMs)

$$\mathcal{D}_\Psi(\mathbb{Q}_1, \mathbb{Q}_2) = \sup_{\psi \in \Psi} \int_{\mathbb{R}} \psi(y) \mathbb{Q}_1(\mathrm{d}y) - \int_{\mathbb{R}} \psi(y) \mathbb{Q}_2(\mathrm{d}y)$$

  
generator



# Integral Probability Metrics (IPMs)

$$\mathcal{D}_\Psi(\mathbb{Q}_1, \mathbb{Q}_2) = \sup_{\psi \in \Psi} \int_{\mathbb{R}} \psi(y) \mathbb{Q}_1(dy) - \int_{\mathbb{R}} \psi(y) \mathbb{Q}_2(dy)$$

## Examples:

IPM	$\Psi$
Kolmogorov distance	$\{\psi : \exists \tau \in \mathbb{R} \text{ with } \psi(y) = \pm 1_{y \leq \tau}\}$
Wasserstein distance	$\{\psi : \text{Lip}(\psi) \leq 1\}$
$\mathcal{L}^p$ -distance ( $\frac{1}{p} + \frac{1}{q} = 1$ )	$\{\psi : \ \psi'\ _{\mathcal{L}^q} \leq 1\}$
Kernel distance	$\{\psi : \ \psi\ _{\mathbb{H}_K} \leq 1\}$
Total variation distance	$\{\psi : \ \psi\ _{\mathcal{L}^\infty} \leq 1\}$

# Relaxing Statistical Parity


**Statistical parity at level  $\varepsilon$ :**

$$\mathcal{D}_\Psi \left( \mathbb{P}_{h(X)|A=0}, \mathbb{P}_{h(X)|A=1} \right) \leq \varepsilon$$

# Relaxing Statistical Parity

**Statistical parity at level  $\varepsilon$ :**

$$\mathcal{D}_\Psi \left( \mathbb{P}_{h(X)|A=0}, \mathbb{P}_{h(X)|A=1} \right) \leq \varepsilon$$

any IPM

# Fair Statistical Learning

**Fair learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)] + \rho(\mathcal{D}_\Psi(\mathbb{P}_{h(X)|A=0}, \mathbb{P}_{h(X)|A=1}))$$

# Fair Statistical Learning

**Fair learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)] + \underbrace{\rho(\mathcal{D}_\Psi(\mathbb{P}_{h(X)|A=0}, \mathbb{P}_{h(X)|A=1}))}_{\text{unfairness penalty}}$$

# Numerical Solution of Fair Learning Problems

# Fair Statistical Learning

**Fair learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)] + \rho(\mathcal{D}_\Psi(\mathbb{P}_{h(X)|A=0}, \mathbb{P}_{h(X)|A=1}))$$

# Fair Statistical Learning

**Fair learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)] + \rho(\mathcal{D}_\Psi(\mathbb{P}_{h(X)|A=0}, \mathbb{P}_{h(X)|A=1}))$$



$$\mathcal{H} = \{h_\theta : \theta \in \Theta\}$$



# Fair Statistical Learning

**Fair learning problem:**

$$\min_{h \in \mathcal{H}} \mathbb{E}[L(h(X), Y)] + \rho(\mathcal{D}_\Psi(\mathbb{P}_{h(X)|A=0}, \mathbb{P}_{h(X)|A=1}))$$



$$\mathcal{H} = \{h_\theta : \theta \in \Theta\}$$

- ▶ all linear hypotheses
- ▶ all neural networks with a fixed architecture

# Fair Statistical Learning

**Fair learning problem:**

$$\min_{\theta \in \Theta} \mathbb{E}[L(h_{\theta}(X), Y)] + \rho(\mathcal{D}_{\Psi}(\mathbb{P}_{h_{\theta}(X)|A=0}, \mathbb{P}_{h_{\theta}(X)|A=1}))$$

# Fair Statistical Learning

**Fair learning problem:**

$$\min_{\theta \in \Theta} \mathbb{E}[L(h_{\theta}(X), Y)] + \rho(\mathcal{D}_{\Psi}(\mathbb{P}_{h_{\theta}(X)|A=0}, \mathbb{P}_{h_{\theta}(X)|A=1}))$$

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples

# Fair Statistical Learning

**Fair learning problem:**

$$\min_{\theta \in \Theta} \mathbb{E}[L(h_{\theta}(X), Y)] + \rho(\mathcal{D}_{\Psi}(\mathbb{P}_{h_{\theta}(X)|A=0}, \mathbb{P}_{h_{\theta}(X)|A=1}))$$

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples

**SGD:**  $\theta_{k+1} = \theta_k - \gamma \cdot g_k(\theta_k)$

# Fair Statistical Learning

**Fair learning problem:**

$$\min_{\theta \in \Theta} \mathbb{E}[L(h_{\theta}(X), Y)] + \rho(\mathcal{D}_{\Psi}(\mathbb{P}_{h_{\theta}(X)|A=0}, \mathbb{P}_{h_{\theta}(X)|A=1}))$$

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples

**SGD:**  $\theta_{k+1} = \theta_k - \gamma \cdot g_k(\theta_k)$



unbiased stochastic gradient  
constructed from batch of  $N$  samples

# Empirical Risk Minimization

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples

  $\frac{1}{N} \sum_{i=1}^N L(h_{\theta}(\hat{X}_i), \hat{Y}_i)$  unbiased estimator for  $\mathbb{E}[L(h(X), Y)]$

# Empirical Risk Minimization

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples

  $\frac{1}{N} \sum_{i=1}^N L(h_{\theta}(\hat{X}_i), \hat{Y}_i)$  unbiased estimator for  $\mathbb{E}[L(h(X), Y)]$

 difficult to find unbiased estimator for unfairness penalty

# Towards an Unbiased Estimator

**Fair learning problem:**

$$\min_{\theta \in \Theta} \mathbb{E}[L(h_{\theta}(X), Y)] + \rho(\underbrace{\mathcal{D}_{\Psi}(\mathbb{P}_{h_{\theta}(X)|A=0})}_{Q_0}, \underbrace{\mathbb{P}_{h_{\theta}(X)|A=1})}_{Q_1})$$



# Unfairness Penalty: Squared Kernel Distance

$$\rho(\mathcal{D}_\Psi(Q_0, Q_1))$$

# Unfairness Penalty: Squared Kernel Distance

$$\mathcal{D}_\Psi(Q_0, Q_1)^2$$

# Unfairness Penalty: Squared Kernel Distance

$$\mathcal{D}_\Psi(\mathbb{Q}_0, \mathbb{Q}_1)^2 = \left( \sup_{\|\psi\|_{\mathbb{H}} \leq 1} \int_{\mathbb{R}} \psi(\mathbf{y}) \mathbb{Q}_0(d\mathbf{y}) - \int_{\mathbb{R}} \psi(\mathbf{y}) \mathbb{Q}_1(d\mathbf{y}) \right)^2$$

# Unfairness Penalty: Squared Kernel Distance

**Theorem:**<sup>1)</sup>  $\mathcal{D}_\Psi(Q_0, Q_1)^2 = \int_{\mathbb{R} \times \mathbb{R}} K(y, y') Q_0(dy) Q_0(dy')$   
 $+ \int_{\mathbb{R} \times \mathbb{R}} K(y, y') Q_1(dy) Q_1(dy')$   
 $- 2 \int_{\mathbb{R} \times \mathbb{R}} K(y, y') Q_0(dy) Q_1(dy')$

admits unbiased  
estimator!

<sup>1)</sup> Sriperumbudur et al., *Electron. J. Stat.*, 2012.

# Unfairness Penalty: Squared Kernel Distance

**Theorem:**<sup>1)</sup>  $\mathcal{D}_\Psi(\mathbb{Q}_0, \mathbb{Q}_1)^2 = \int_{\mathbb{R} \times \mathbb{R}} K(\mathbf{y}, \mathbf{y}') \mathbb{Q}_0(d\mathbf{y}) \mathbb{Q}_0(d\mathbf{y}') + \int_{\mathbb{R} \times \mathbb{R}} K(\mathbf{y}, \mathbf{y}') \mathbb{Q}_1(d\mathbf{y}) \mathbb{Q}_1(d\mathbf{y}') - 2 \int_{\mathbb{R} \times \mathbb{R}} K(\mathbf{y}, \mathbf{y}') \mathbb{Q}_0(d\mathbf{y}) \mathbb{Q}_1(d\mathbf{y}')$

**Data:**  $\hat{Y}_i^0, \dots, \hat{Y}_{N_0}^0 \sim \mathbb{Q}_0$  i.i.d.,  $\hat{Y}_i^1, \dots, \hat{Y}_{N_1}^1 \sim \mathbb{Q}_1$  i.i.d.

# Unfairness Penalty: Squared Kernel Distance

**Theorem:**  $^1) \mathcal{D}_\Psi(\mathbb{Q}_0, \mathbb{Q}_1)^2 = \int_{\mathbb{R} \times \mathbb{R}} K(y, y') \mathbb{Q}_0(dy) \mathbb{Q}_0(dy') + \int_{\mathbb{R} \times \mathbb{R}} K(y, y') \mathbb{Q}_1(dy) \mathbb{Q}_1(dy') - 2 \int_{\mathbb{R} \times \mathbb{R}} K(y, y') \mathbb{Q}_0(dy) \mathbb{Q}_1(dy')$

**Data:**  $\hat{Y}_i^0, \dots, \hat{Y}_{N^0}^0 \sim \mathbb{Q}_0$  i.i.d.,  $\hat{Y}_i^1, \dots, \hat{Y}_{N^1}^1 \sim \mathbb{Q}_1$  i.i.d.

**Unbiased estimator for  $\mathcal{D}_\Psi(\mathbb{Q}_0, \mathbb{Q}_1)^2$ :**

$$\sum_{a \in \{0,1\}} \frac{1}{N^a(N^a - 1)} \sum_{\substack{i,j=1 \\ i \neq j}}^{N^a} K(\hat{Y}_i^a, \hat{Y}_j^a) - 2 \frac{1}{N^0 N^1} \sum_{i=1}^{N^0} \sum_{j=1}^{N^1} K(\hat{Y}_i^0, \hat{Y}_j^1)$$

# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

**Batch:** At least 2 samples from each class and  $\bar{N}$  samples in total



# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

**Batch:** At least 2 samples from each class and  $\bar{N}$  samples in total

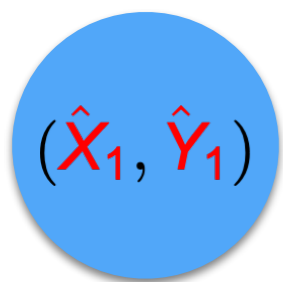
**Example:**  $\bar{N} = 5$

# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

**Batch:** At least 2 samples from each class and  $\bar{N}$  samples in total

**Example:**  $\bar{N} = 5$

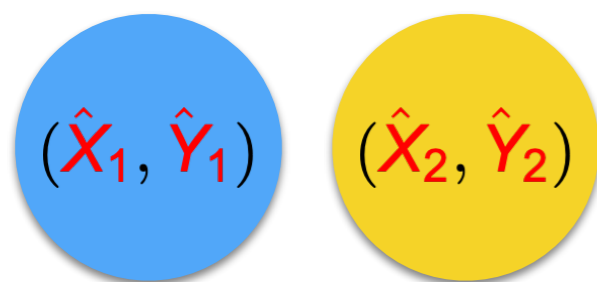


# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

**Batch:** At least 2 samples from each class and  $\bar{N}$  samples in total

**Example:**  $\bar{N} = 5$

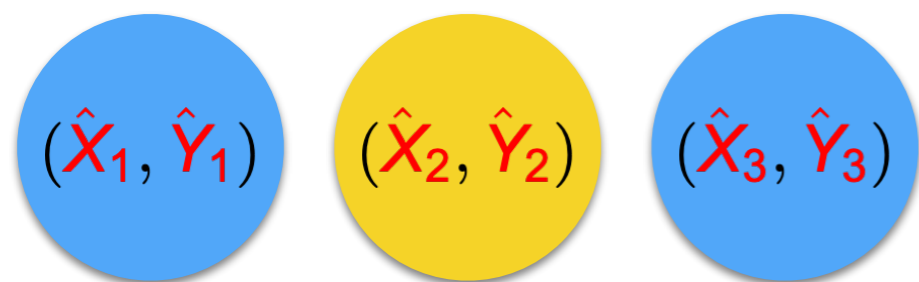


# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

**Batch:** At least 2 samples from each class and  $\bar{N}$  samples in total

**Example:**  $\bar{N} = 5$

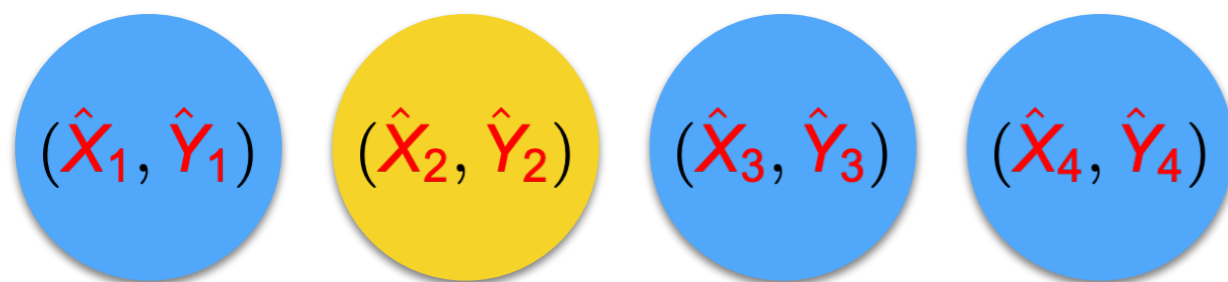


# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

**Batch:** At least 2 samples from each class and  $\bar{N}$  samples in total

**Example:**  $\bar{N} = 5$

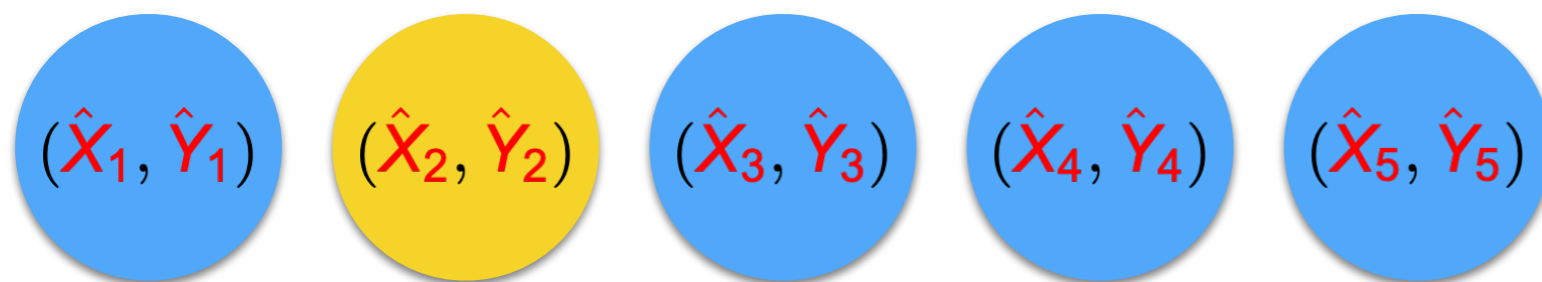


# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

**Batch:** At least 2 samples from each class and  $\bar{N}$  samples in total

**Example:**  $\bar{N} = 5$



first 5 samples contain only  
one sample from class 1

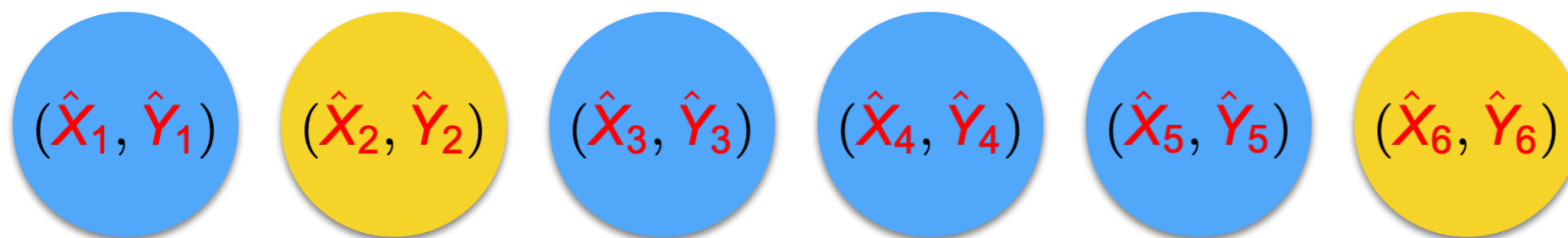
⇒ estimator for unfairness penalty undefined!

# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

**Batch:** At least 2 samples from each class and  $\bar{N}$  samples in total

**Example:**  $\bar{N} = 5$



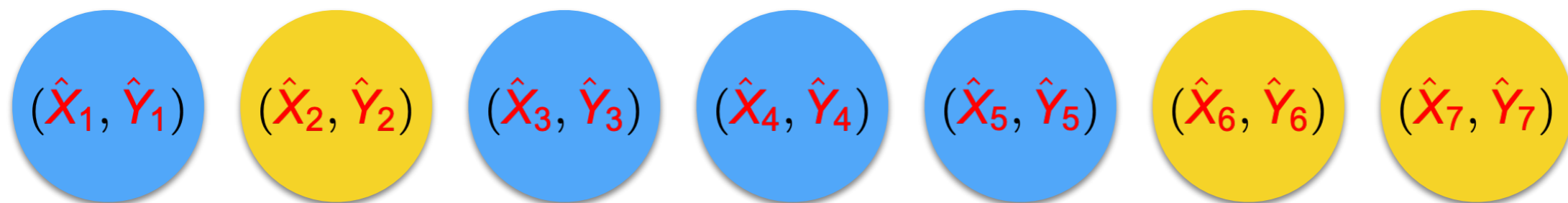
complete batch with 6 samples  
(4 from class 0, 2 from class 1)

# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

**Batch:** At least 2 samples from each class and  $\bar{N}$  samples in total

**Example:**  $\bar{N} = 5$



complete batch with 6 samples  
(4 from class 0, 2 from class 1)

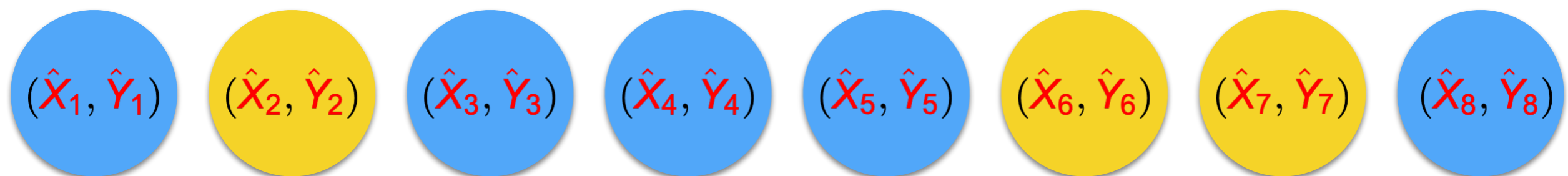


# Random Batches

**Data:**  $(\hat{X}_i, \hat{Y}_i, \hat{A}_i)$ ,  $i \in \mathbb{N}$ , i.i.d. samples (● :  $A_i = 0$ , ● :  $A_i = 1$ )

**Batch:** At least 2 samples from each class and  $\bar{N}$  samples in total

**Example:**  $\bar{N} = 5$



complete batch with 6 samples  
(4 from class 0, 2 from class 1)

next batch

# Unbiased Estimators: Unfairness Penalty

- Notation:**
- ▶  $\mathcal{I}_b \subseteq \mathbb{N}$   $b$ -th batch
  - ▶  $\mathcal{I}_b^a$  = class  $a$  samples in  $\mathcal{I}_b$ ,  $a \in \{0, 1\}$

# Unbiased Estimators: Unfairness Penalty

- Notation:**
- ▶  $\mathcal{I}_b \subseteq \mathbb{N}$   $b$ -th batch
  - ▶  $\mathcal{I}_b^a$  = class  $a$  samples in  $\mathcal{I}_b$ ,  $a \in \{0, 1\}$

**Lemma:** The following estimator of the unfairness penalty is unbiased for every batch  $b$ .

$$\hat{U}_b(\theta) = \begin{cases} \sum_{a \in \{0,1\}} \frac{1}{|\mathcal{I}_b^a|(|\mathcal{I}_b^a| - 1)} \sum_{i \neq j \in \mathcal{I}_b^a} K(h_\theta(\hat{X}_i), h_\theta(\hat{X}_j)) \\ -2 \frac{1}{|\mathcal{I}_b^0| \cdot |\mathcal{I}_b^1|} \sum_{i \in \mathcal{I}_b^0} \sum_{j \in \mathcal{I}_b^1} K(h_\theta(\hat{X}_i), h_\theta(\hat{X}_j)) \end{cases}$$

# Unbiased Estimators: Unfairness Penalty

- Notation:**
- ▶  $\mathcal{I}_b \subseteq \mathbb{N}$   $b$ -th batch
  - ▶  $\mathcal{I}_b^a$  = class  $a$  samples in  $\mathcal{I}_b$ ,  $a \in \{0, 1\}$

**Lemma:** The following estimator of the unfairness penalty is unbiased for every batch  $b$ .

$$\hat{U}_b(\theta) = \begin{cases} \sum_{a \in \{0,1\}} \frac{1}{|\mathcal{I}_b^a|(|\mathcal{I}_b^a| - 1)} \sum_{i \neq j \in \mathcal{I}_b^a} K(h_\theta(\hat{X}_i), h_\theta(\hat{X}_j)) \\ -2 \frac{1}{|\mathcal{I}_b^0| \cdot |\mathcal{I}_b^1|} \sum_{i \in \mathcal{I}_b^0} \sum_{j \in \mathcal{I}_b^1} K(h_\theta(\hat{X}_i), h_\theta(\hat{X}_j)) \end{cases}$$

**Note:** All index sets are random!

# Unbiased Estimators: Unfairness Penalty

- Notation:**
- ▶  $\mathcal{I}_b \subseteq \mathbb{N}$   $b$ -th batch
  - ▶  $\mathcal{I}_b^a$  = class  $a$  samples in  $\mathcal{I}_b$ ,  $a \in \{0, 1\}$

**Lemma:** The following estimator of the unfairness penalty is unbiased for every batch  $b$ .

$$\hat{U}_b(\theta) = \begin{cases} \sum_{a \in \{0,1\}} \frac{1}{|\mathcal{I}_b^a|(|\mathcal{I}_b^a| - 1)} \sum_{i \neq j \in \mathcal{I}_b^a} K(h_\theta(\hat{X}_i), h_\theta(\hat{X}_j)) \\ -2 \frac{1}{|\mathcal{I}_b^0| \cdot |\mathcal{I}_b^1|} \sum_{i \in \mathcal{I}_b^0} \sum_{j \in \mathcal{I}_b^1} K(h_\theta(\hat{X}_i), h_\theta(\hat{X}_j)) \end{cases}$$

$\implies \nabla_\theta \hat{U}_b(\theta)$  is an unbiased stochastic gradient

# Unbiased Estimators: Prediction Loss

**Empirical prediction loss:**

$$\frac{1}{|\mathcal{I}_b|} \sum_{i \in \mathcal{I}_b} L(h_\theta(\hat{X}_i), \hat{Y}_i)$$

# Unbiased Estimators: Prediction Loss

**Empirical prediction loss:**

$$\frac{1}{|\mathcal{I}_b|} \sum_{i \in \mathcal{I}_b} L(h_\theta(\hat{X}_i), \hat{Y}_i)$$

biased because  $\mathcal{I}_b$  is random!

# Unbiased Estimators: Prediction Loss

**Empirical prediction loss:**

$$\frac{1}{|\mathcal{I}_b|} \sum_{a \in \{0,1\}} \sum_{i \in \mathcal{I}_b^a} L(h_\theta(\hat{X}_i), \hat{Y}_i)$$



# Unbiased Estimators: Prediction Loss

**Empirical prediction loss:**

$$\frac{1}{|\mathcal{I}_b|} \sum_{a \in \{0,1\}} \sum_{i \in \mathcal{I}_b^a} \underbrace{\Delta(|\mathcal{I}_b|, |\mathcal{I}_b^a|)}_{\text{bias correction term}} \cdot L(h_\theta(\hat{X}_i), \hat{Y}_i)$$

bias correction term

**Definition:** For  $N \in \{\bar{N}, \bar{N} + 1, \dots\}$  and  $n \in \{2, \dots, N - 2\}$ , set

$$\Delta(N, n) = 1_{N=\bar{N}} + \frac{N}{2(N-1)} 1_{(N>\bar{N}) \wedge (n=2)} + \frac{N}{N-1} 1_{(N>\bar{N}) \wedge (n=N-2)}$$

# Unbiased Estimators: Prediction Loss

**Empirical prediction loss:**

$$\frac{1}{|\mathcal{I}_b|} \sum_{a \in \{0,1\}} \sum_{i \in \mathcal{I}_b^a} \underbrace{\Delta(|\mathcal{I}_b|, |\mathcal{I}_b^a|)}_{\text{bias correction term}} \cdot L(h_\theta(\hat{X}_i), \hat{Y}_i)$$

bias correction term

**Definition:** For  $N \in \{\bar{N}, \bar{N} + 1, \dots\}$  and  $n \in \{2, \dots, N - 2\}$ , set

$$\Delta(N, n) = 1_{N=\bar{N}} + \frac{N}{2(N-1)} 1_{(N>\bar{N}) \wedge (n=2)} + \frac{N}{N-1} 1_{(N>\bar{N}) \wedge (n=N-2)}$$

# Unbiased Estimators: Prediction Loss

**Lemma:** The following estimator of the prediction loss is unbiased for every batch  $b$ .

$$\hat{R}_b(\theta) = \frac{1}{|\mathcal{I}_b|} \sum_{a \in \{0,1\}} \sum_{i \in \mathcal{I}_b^a} \Delta(|\mathcal{I}_b|, |\mathcal{I}_b^a|) \cdot L(h_\theta(\hat{X}_i), \hat{Y}_i)$$

# Unbiased Estimators: Prediction Loss

**Lemma:** The following estimator of the prediction loss is unbiased for every batch  $b$ .

$$\hat{R}_b(\theta) = \frac{1}{|\mathcal{I}_b|} \sum_{a \in \{0,1\}} \sum_{i \in \mathcal{I}_b^a} \Delta(|\mathcal{I}_b|, |\mathcal{I}_b^a|) \cdot L(h_\theta(\hat{X}_i), \hat{Y}_i)$$

$\implies \nabla_\theta \hat{R}_b(\theta)$  is an unbiased stochastic gradient

# SGD Convergence

**Fair learning problem:**

$$\min_{\theta \in \Theta} \mathbb{E}[L(h_{\theta}(X), Y)] + \rho(\mathcal{D}_{\Psi}(\mathbb{P}_{h_{\theta}(X)|A=0}, \mathbb{P}_{h_{\theta}(X)|A=1})) \quad (*)$$

# SGD Convergence

**Fair learning problem:**

$$\min_{\theta \in \Theta} \mathbb{E}[L(h_{\theta}(X), Y)] + \rho(\mathcal{D}_{\Psi}(\mathbb{P}_{h_{\theta}(X)|A=0}, \mathbb{P}_{h_{\theta}(X)|A=1})) \quad (*)$$

**Theorem:** If  $\mathcal{D}_{\Psi}$  is a kernel distance and  $\rho(z) = \lambda z^2$  with  $\lambda \geq 0$ , then  $\nabla_{\theta} \hat{R}_b(\theta) + \lambda \cdot \nabla_{\theta} \hat{U}_b(\theta)$  is an unbiased gradient estimator for  $(*)$ .

# SGD Convergence

**Fair learning problem:**

$$\min_{\theta \in \Theta} \mathbb{E}[L(h_{\theta}(X), Y)] + \rho(\mathcal{D}_{\Psi}(\mathbb{P}_{h_{\theta}(X)|A=0}, \mathbb{P}_{h_{\theta}(X)|A=1})) \quad (*)$$

**Theorem:** If  $\mathcal{D}_{\Psi}$  is a kernel distance and  $\rho(z) = \lambda z^2$  with  $\lambda \geq 0$ , then  $\nabla_{\theta} \hat{R}_b(\theta) + \lambda \cdot \nabla_{\theta} \hat{U}_b(\theta)$  is an unbiased gradient estimator for  $(*)$ .

$\implies$  SGD converges (in expectation) to a stationary point of  $(*)$

# Numerical Experiments



# Regression

## Synthetic data:

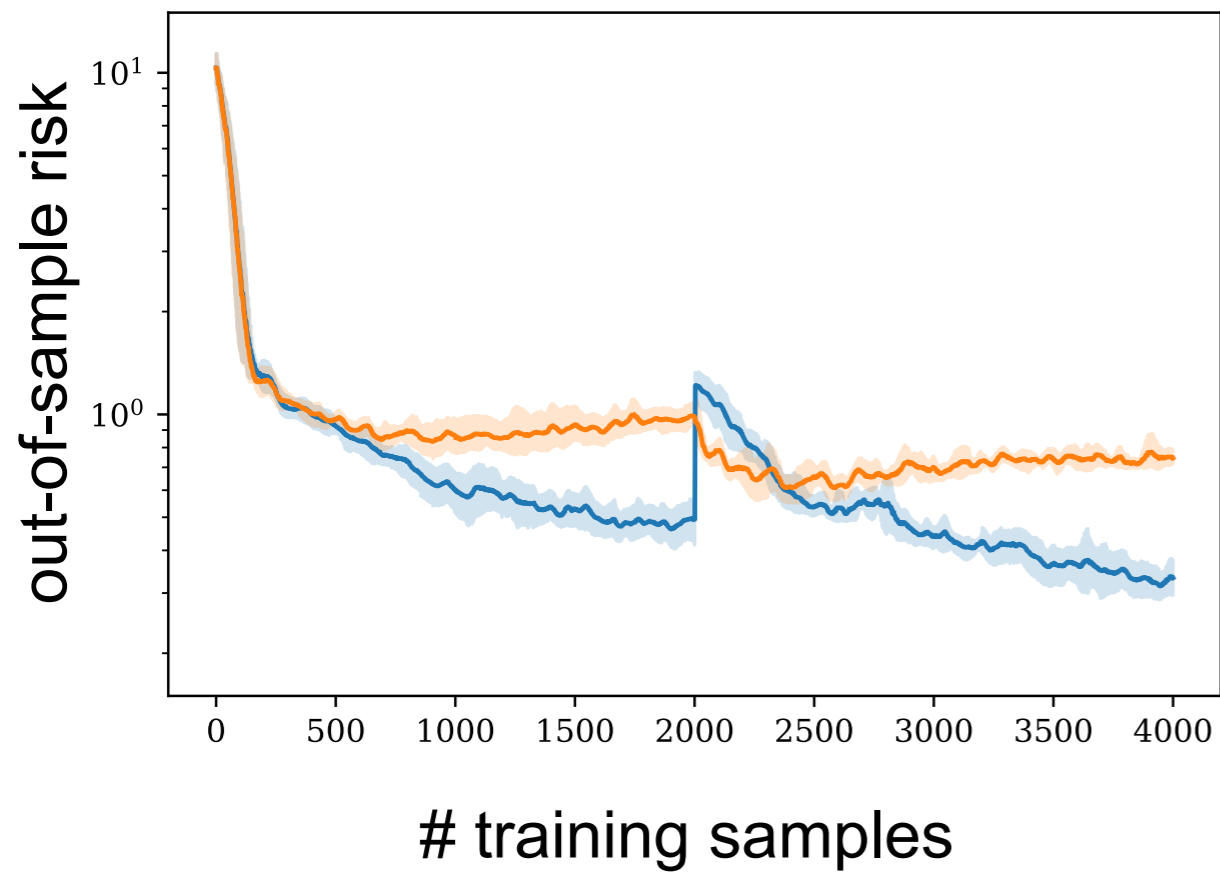
- ▶ Input:  $X \sim \mathcal{U}([0, 1]^9 \times \{0, 1\})$
- ▶ Sensitive attribute:  $A = X_{10}$
- ▶ Output:  $Y = \max\{s_1^\top X, \dots, s_5^\top X\}$

## Regression model:

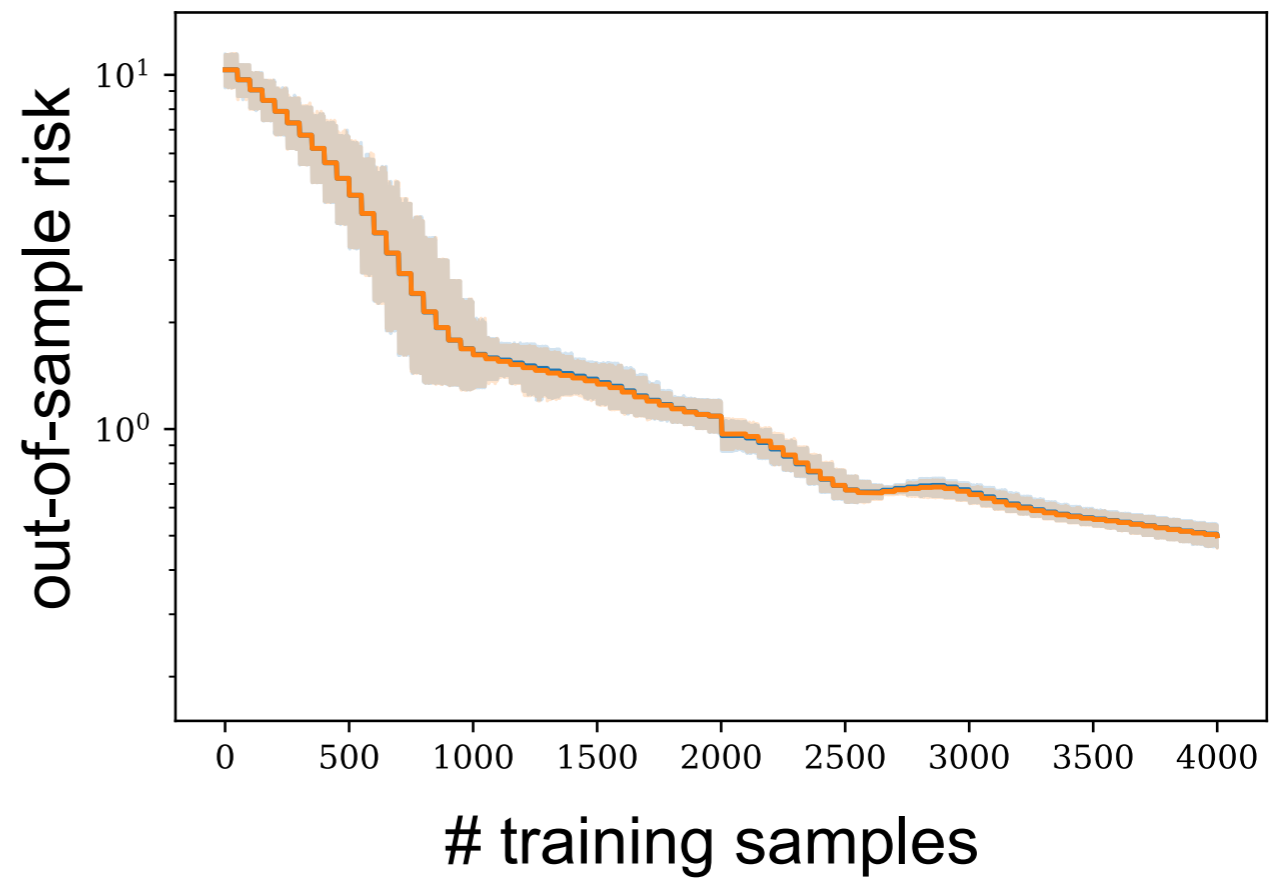
- ▶ Square loss:  $L(\hat{y}, y) = (\hat{y} - y)^2$
- ▶ Predictor: 3-layer NN with 20 hidden nodes

# Regression

Batch size  $\bar{N} = 4$

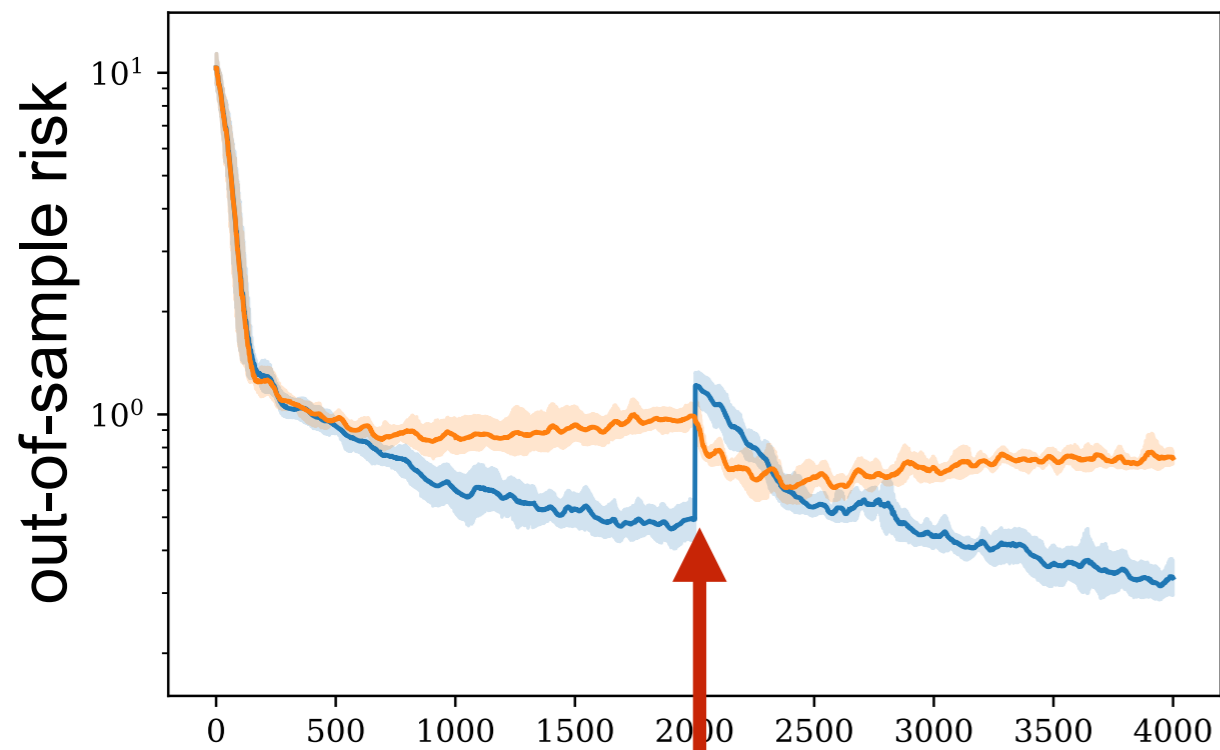


Batch size  $\bar{N} = 50$



# Regression

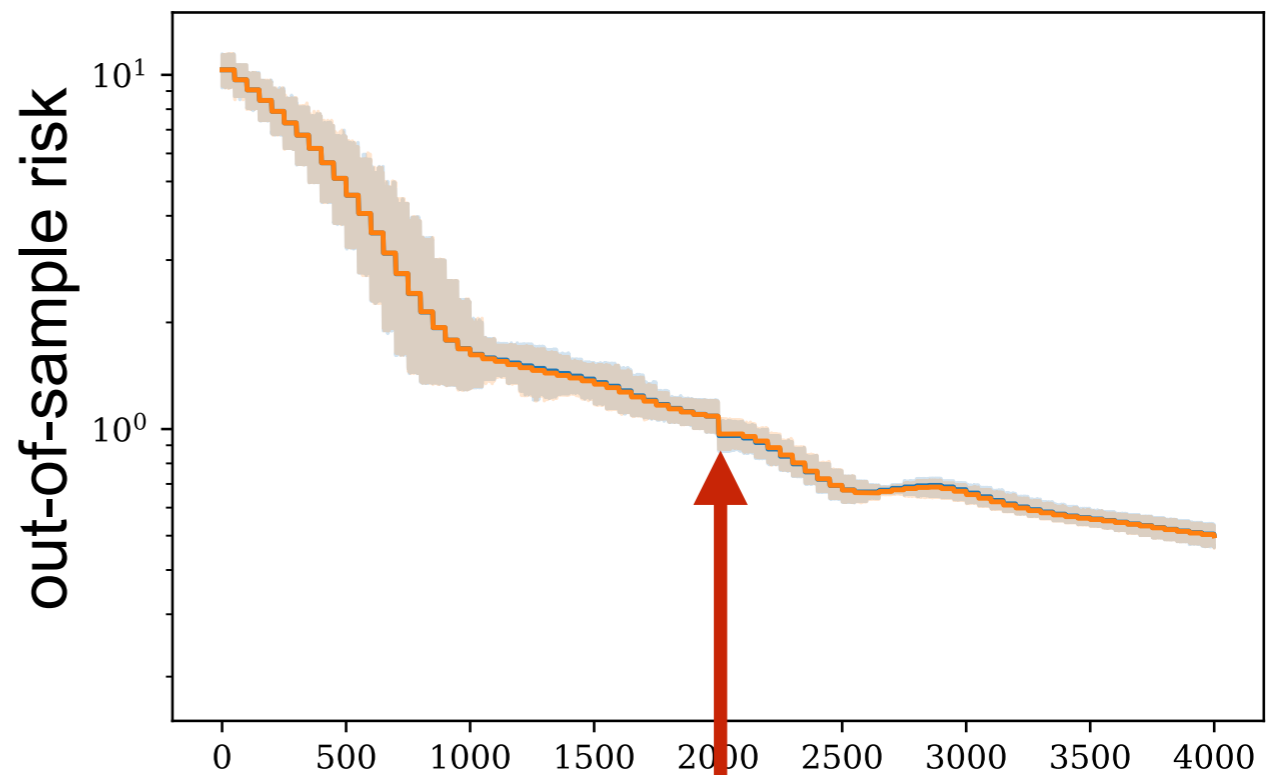
Batch size  $\bar{N} = 4$



# training samples

regime shift:  
new output

Batch size  $\bar{N} = 50$

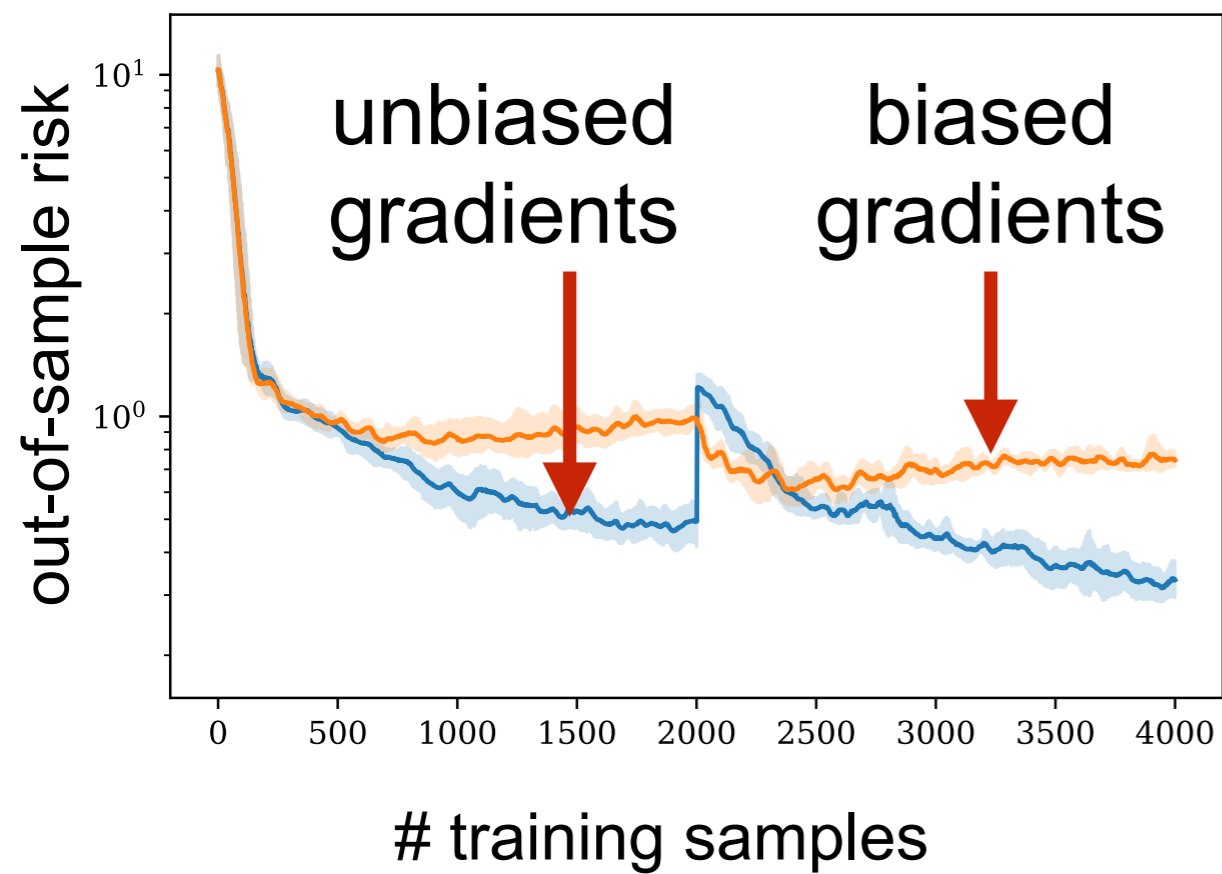


# training samples

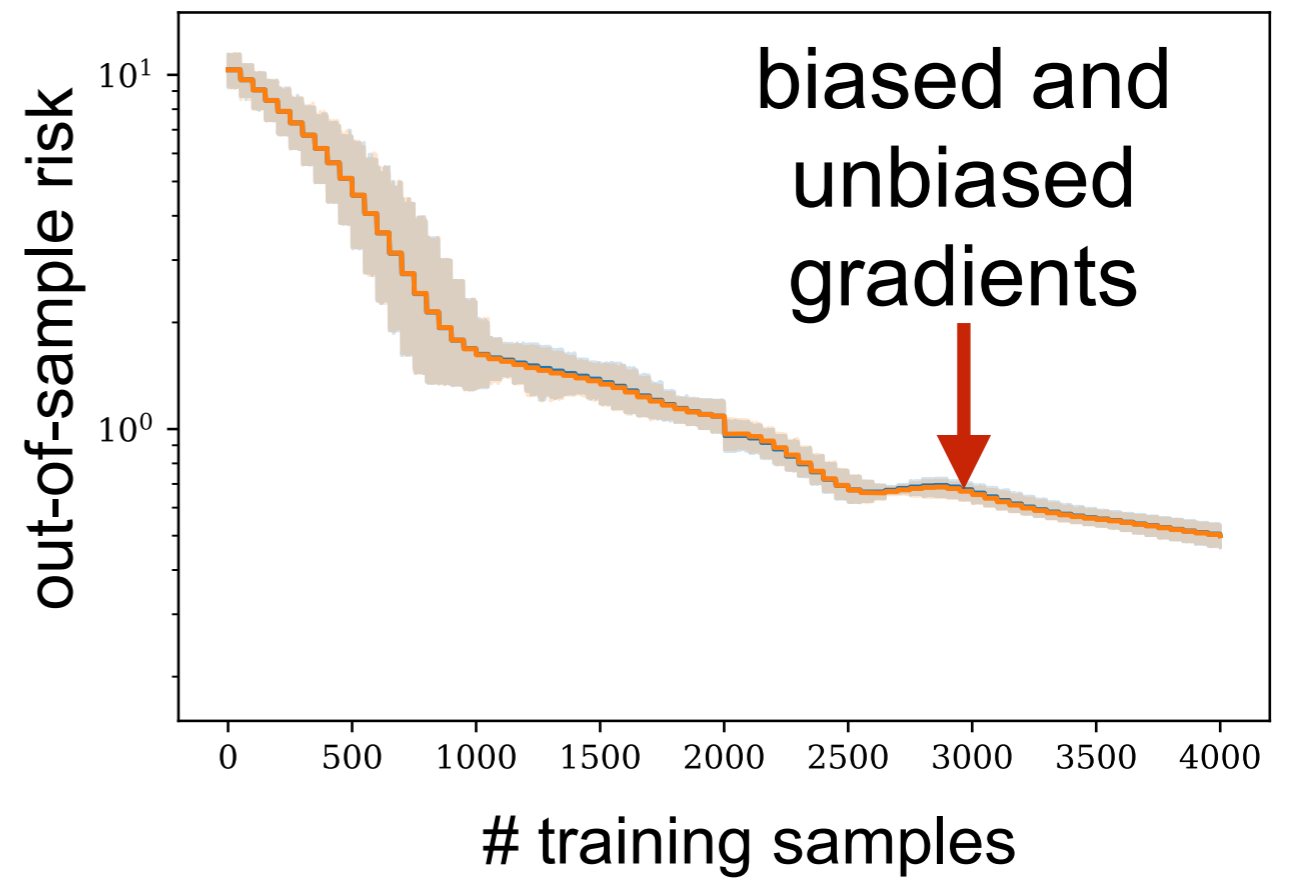
regime shift:  
new output

# Regression

Batch size  $\bar{N} = 4$



Batch size  $\bar{N} = 50$



# Merits of Unbiased Gradient Estimators

		batch size	
		small	large
gradients	biased		
	unbiased		

# Merits of Unbiased Gradient Estimators

		batch size	
		small	large
gradients	biased	fast convergence to bad solution	
	unbiased		

# Merits of Unbiased Gradient Estimators

		batch size	
		small	large
gradients	biased	fast convergence to bad solution	slow convergence to good solution
	unbiased		

# Merits of Unbiased Gradient Estimators

		batch size	
		small	large
gradients	biased	fast convergence to bad solution	slow convergence to good solution
	unbiased		slow convergence to good solution



# Merits of Unbiased Gradient Estimators

		batch size	
		small	large
gradients	biased	fast convergence to bad solution	slow convergence to good solution
	unbiased	fast convergence to good solution	slow convergence to good solution

## Drug dataset:<sup>1)</sup>

- ▶ Input: personality type, level of education, age etc.
- ▶ Sensitive attribute: race
- ▶ Output: “used” vs. “never used” for heroin

## Classification model:

- ▶ Cross-entropy loss:  $L(\hat{y}, y) = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$
- ▶ Predictor: 3-layer NN with 16 hidden nodes

---

<sup>1)</sup> <https://archive.ics.uci.edu/ml/datasets/Drug+consumption+%28quantified%29>

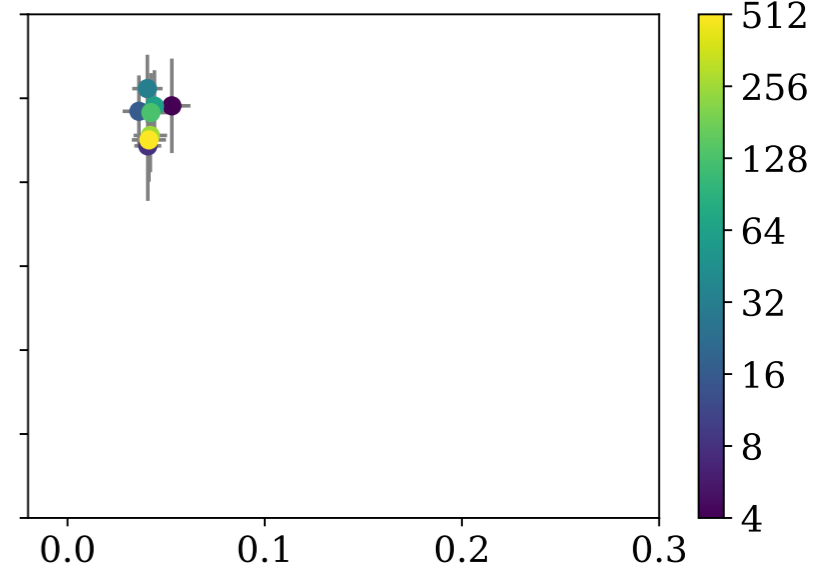
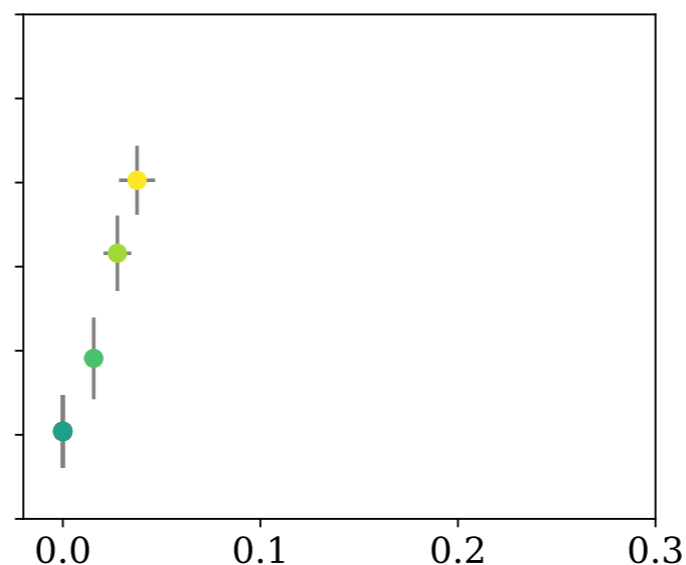
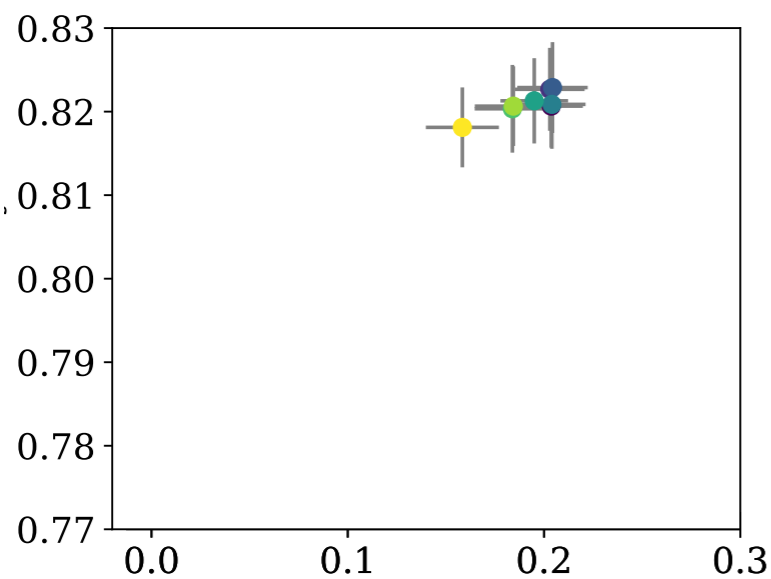
# Classification

no unfairness penalty

biased gradients

unbiased gradients

accuracy ↑



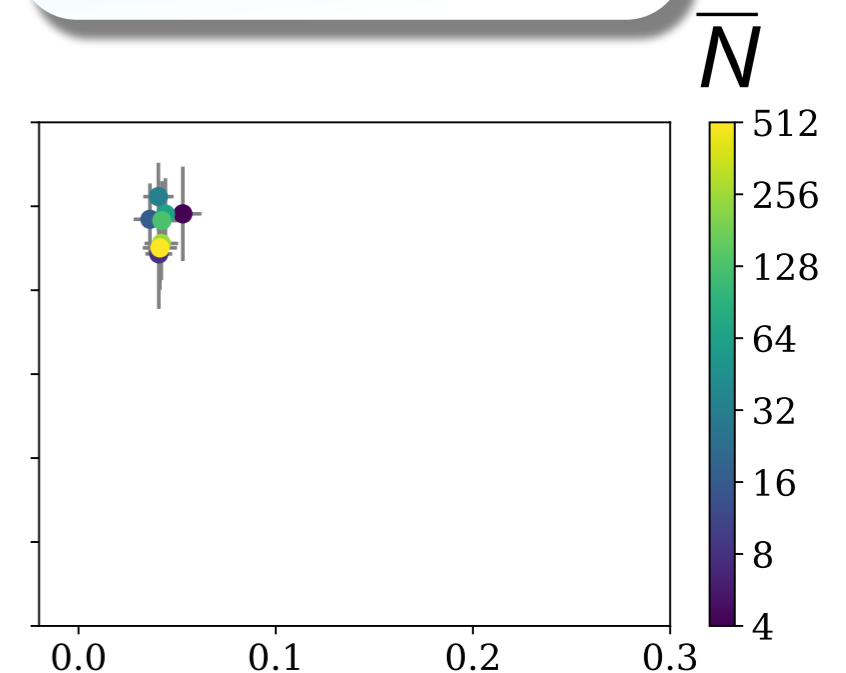
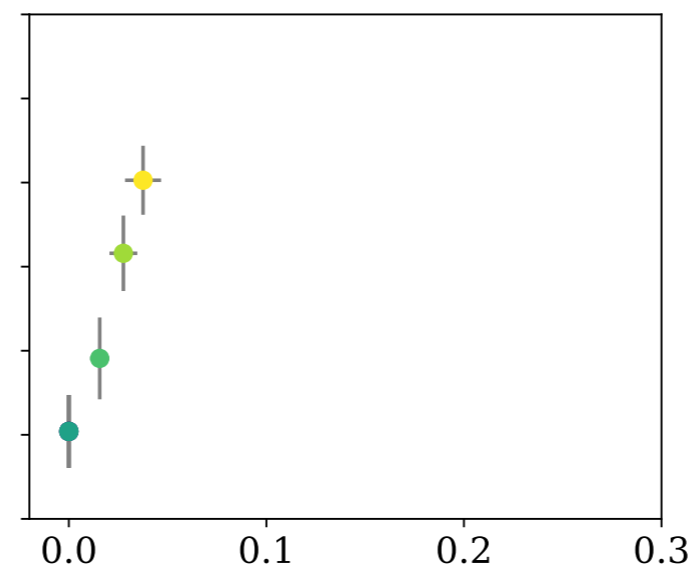
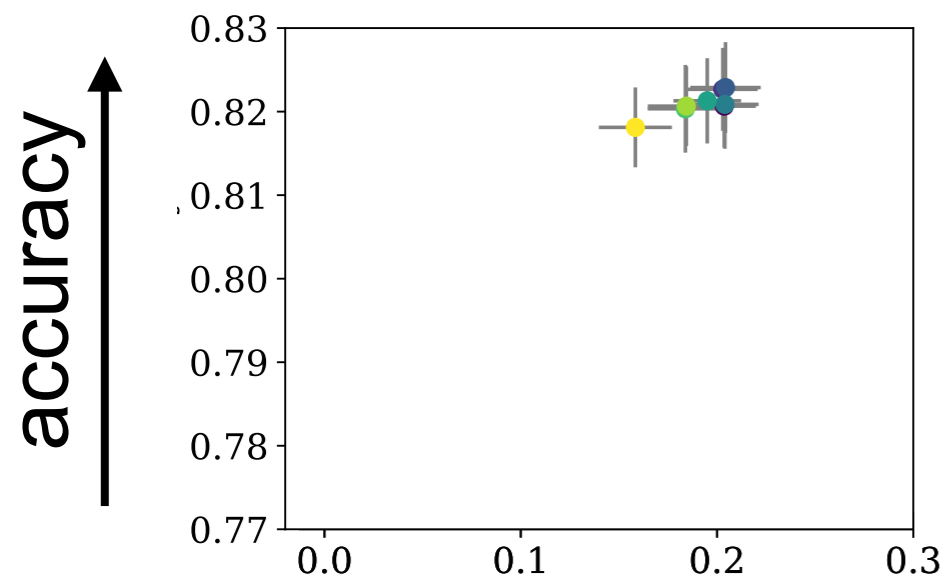
unfairness →

# Classification

no unfairness penalty

biased gradients

unbiased gradients



## Classical ERM:

- ▶ high accuracy
- ▶ high unfairness

## Fair ERM:

- ▶ accuracy sensitive to  $\bar{N}$
- ▶ low unfairness

## Our method:

- ▶ high accuracy
- ▶ low unfairness

# Conclusions

## ▶ Impact of SP constraints

- ▶  $Y$  has no bias in training  $\implies$  SP increases test error
- ▶  $Y$  has small bias in training &  $A$  is irrelevant for predicting  $Y$   
 $\implies$  SP decreases test error
- ▶ Good sensitive attribute: Any feature  $A$  with  $\mathbb{P}_{Y|X} \perp A$

## ▶ Learning problems with unfairness penalties

- ▶ Any IPM provides an unfairness measure
- ▶ Empirical estimator of unfairness penalty is biased
- ▶ Moore Aronszajn theorem  $\implies$  squared kernel distance admits unbiased estimator
- ▶ Fair learning problems can be solved with SGD

# This Talk is Based on...

Y. Rychener, B. Taşkesen, D. Kuhn. **Metrizing Fairness**. arXiv. 2024.



**Yves Rychener**



**Bahar Taşkesen**

