

---

---

---

---

---



# Bias Mitigation in Machine Learning (Revised 2025)

# Ways Things Go Wrong

- Have much less data on some group (fine if groups all "same")
- Different groups have different distributions
- Our features are less predictive on some group
- Some group inherently less predictable
- Our data is biased in the first place

# Fairness in ML

- Typically a property of a **model** (ML algo output)
- Exceptions: online decision-making, RL, bandit settings
- Multiple **types** of fairness definitions

# Types of Model Fairness

- Group fairness  
(most common)
- Individual fairness
- Interpolations between  
the two
- Others (causal, fair  
representations, ...)

# Group Fairness Notions

Start by identifying:

- groups or attributes we wish to "protect" (e.g. race, gender)
- what constitutes harm (e.g. error, false pos/neg)

---

Choices are subjective & domain-specific

Then seek to equalize rates of harm across groups.

Example:

- domain: consumer lending
- groups: male & female
- harm: false rejection (negs)

Want to find model  $h(x)$  s.t.

$$FN(h, \text{male}) \approx FN(h, \text{female})$$

↗  
↖ allows for optimization of overall error

If we are given a model  $h(x)$  & have access to group membership, easy to audit  $h(x)$  for fairness.

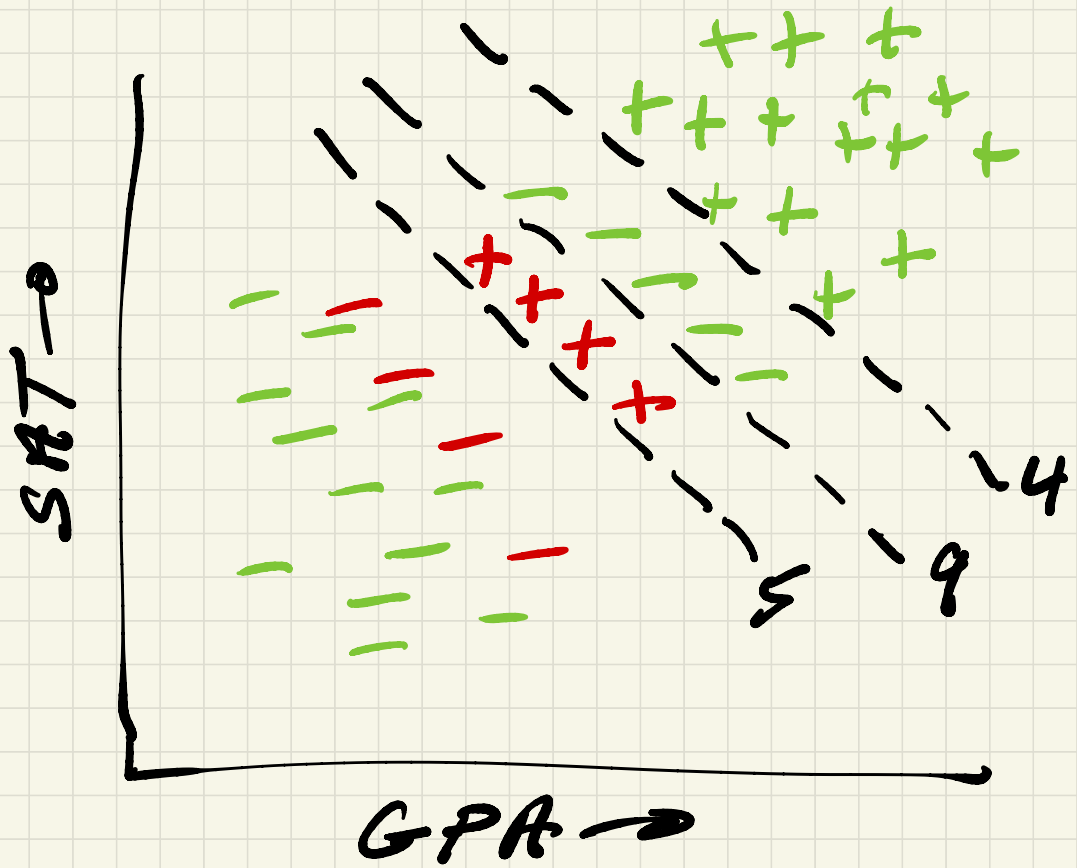
---

How can we learn a fair model  $h(x)$ ?

Why won't standard ML algos work? •



# A more subtle example



"Bolt-on" Bias Mitigation

- Suppose we are given a model  $h \in H$  trained to minimize overall error  $\epsilon = \epsilon(h)$  within  $H$
- But  $h$  is "unfair", meaning that
 
$$\epsilon_A = \epsilon_A(h) = P_{\langle x, y \rangle \sim P} [h(x) \neq y | x \in A]$$

$$\epsilon_B = \epsilon_B(h) = P_{\langle x, y \rangle \sim P} [h(x) \neq y | x \in B]$$
 are unequal - say  $\epsilon_A \neq \epsilon_B \neq 1/2$
- Here  $A \neq B$  are two disjoint & exhaustive subgroups, so  $A \cap B = \emptyset, A \cup B = X$
- Assume  $\forall x$  we know whether  $x \in A$  or  $x \in B$  (might be a feature)
- Goal: Modify  $h$  to get a fair model in which error rates on  $A \neq B$  are equal.

## Technique: Selective Randomization

- Define a new model  $\tilde{h}$  "on  $H_P$ " of  $h$ :
  - if  $x \in B$ , define  $\tilde{h}(x) = h(x)$
  - if  $x \in A$ : for  $g \in [0, 1]$  TBP, with prob.  $1-g$  let  $\tilde{h}(x) = h(x)$  with prob.  $g$  let  $\tilde{h}(x)$  be determined by the flip of a fair coin
- Let  $\tilde{\epsilon}, \tilde{\epsilon}_A, \tilde{\epsilon}_B$  denote the error of  $\tilde{h}$  overall & on subgroups  $A$  &  $B$
- Note  $\tilde{\epsilon}_B = \epsilon_B$  (didn't change predictions on  $B$ )

• Now

$$\tilde{\epsilon}_A = \underbrace{(1-q)}_{\text{prob. we use } h(x)} \epsilon_A + q \underbrace{\left(\frac{1}{2}\right)}_{\substack{\text{prob. we} \\ \text{flip coin}}}$$

$\epsilon_A$  ← error rate of coin flip

• Set  $\tilde{\epsilon}_A = \epsilon_B$  & solve for  $q$ :

$$(1-q)\epsilon_A + q/2 = \epsilon_B$$

$$q\left(\frac{1}{2} - \epsilon_A\right) + \epsilon_A = \epsilon_B$$

$$q = \frac{\epsilon_B - \epsilon_A}{\frac{1}{2} - \epsilon_A}$$

↳ At this value of  $q$ , we equalize the error rates of  $\tilde{h}$  on A & B to be  $\epsilon_B$ .  
The overall error  $\tilde{\epsilon}$  of  $\tilde{h}$  will also be  $\epsilon_B$ .

Let's do some examples.

• Suppose  $\epsilon_A = 0.05$ ,  $\epsilon_B = 0.2$ . Then

$$g = \frac{0.2 - 0.05}{0.5 - 0.05} = \frac{0.15}{0.45} = \frac{1}{3}$$

So for group B, we flip  
a coin  $\frac{1}{3}$  of the time

•  $\epsilon_A = 0.1$ ,  $\epsilon_B = 0.2$ :

$$g = \frac{0.2 - 0.1}{\frac{1}{2} - 0.1} = \frac{0.1}{0.4} = \frac{1}{4}$$

• So when  $\epsilon_A$  &  $\epsilon_B$  are closer  
we randomize less  
on group B.

- Now suppose we **weaken** our fairness demand - asking only that

$$\tilde{\epsilon}_A = \epsilon_B - \gamma \text{ for some } \gamma > 0$$

- Then we again solve:

$$\tilde{\epsilon}_A = (1-q)\epsilon_A + q/2 = \epsilon_B - \gamma$$

which gives

$$q = \frac{\epsilon_B - \epsilon_A - \gamma}{1/2 - \epsilon_A}$$

↗ Note that q decreases as  $\gamma$  increases - the less fairness we demand, the less we randomize on group B

• Also  $\gamma = \epsilon_B - \epsilon_A$  is largest sensible value for  $\gamma$ , and corresponds to not changing  $h$  at all

• And now the overall error of  $\tilde{h}$  will be:

$$\begin{aligned}\tilde{\epsilon} &= p_A(\epsilon_B - \gamma) + p_B \epsilon_B \\ &= (p_A + p_B) \epsilon_B - p_A \gamma \\ &= \epsilon_B - p_A \gamma\end{aligned}$$

Here  $p_A, p_B$  are probs. of groups  $A$  &  $B$ , so  $p_A + p_B = 1$



• Example:  $\epsilon_A = 0.05$ ,  $\epsilon_B = 0.2$ ,  
 $p_A = 0.7$ ,  $\gamma = 0.03$ :

$$q = \frac{\epsilon_B - \epsilon_A - \gamma}{\frac{1}{2} - \epsilon_A} = \frac{0.2 - 0.05 - 0.03}{\frac{1}{2} - 0.05}$$

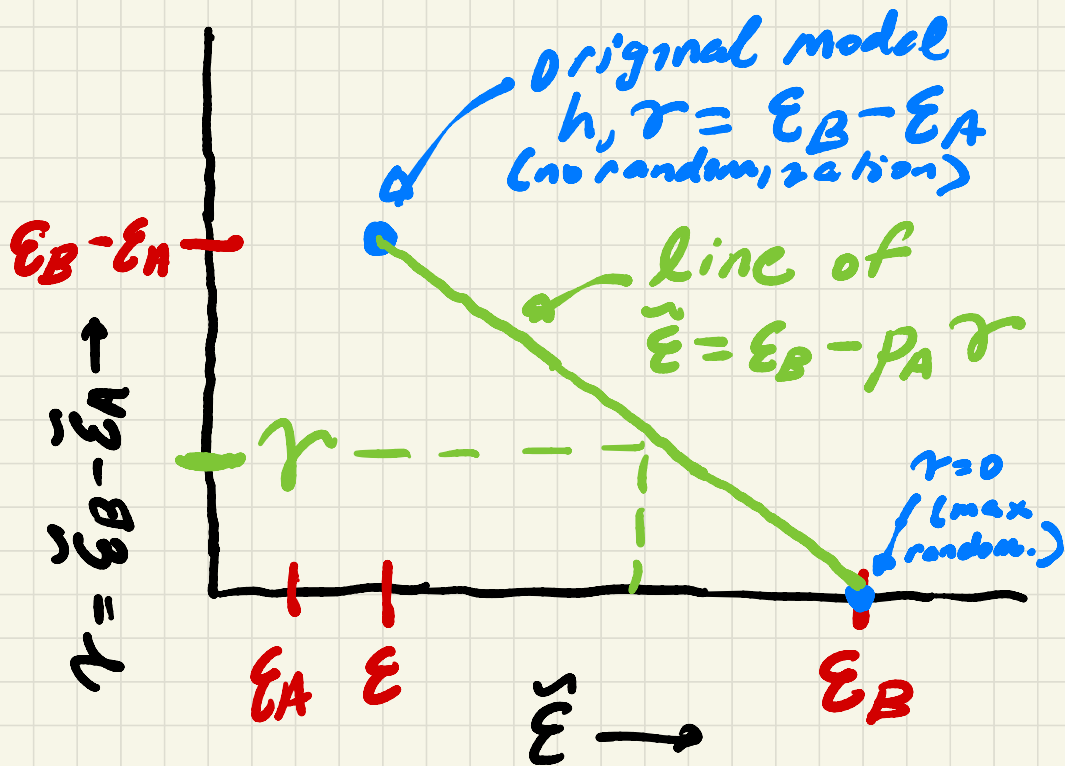
$$= \frac{0.12}{0.45} \approx 0.267$$

and  $\tilde{\epsilon} = \epsilon_B - p_A \gamma$

$$= 0.2 - 0.7 \times 0.03$$

$$= 0.179$$

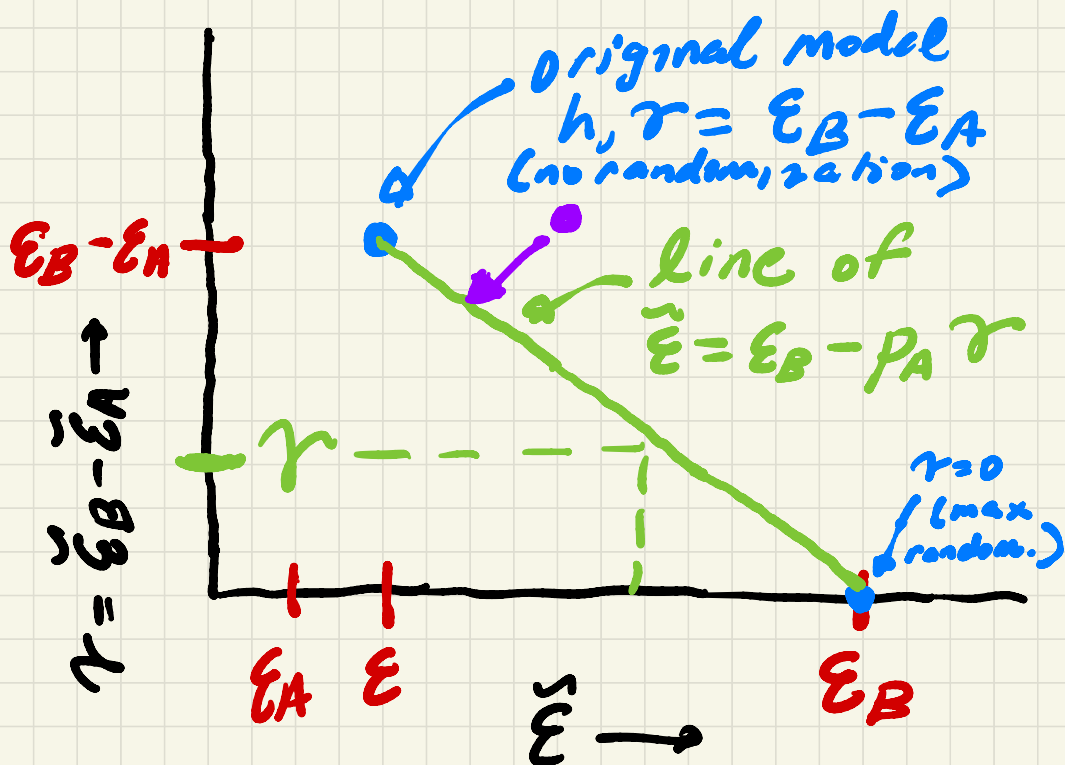
# Big Picture: Summary



$$\text{Slope of line} = \frac{\epsilon_B - \epsilon_A}{\epsilon_B - \epsilon}$$

This line gives the achievable trade-offs between error & unfairness in this scheme.

# Big Picture: Summary



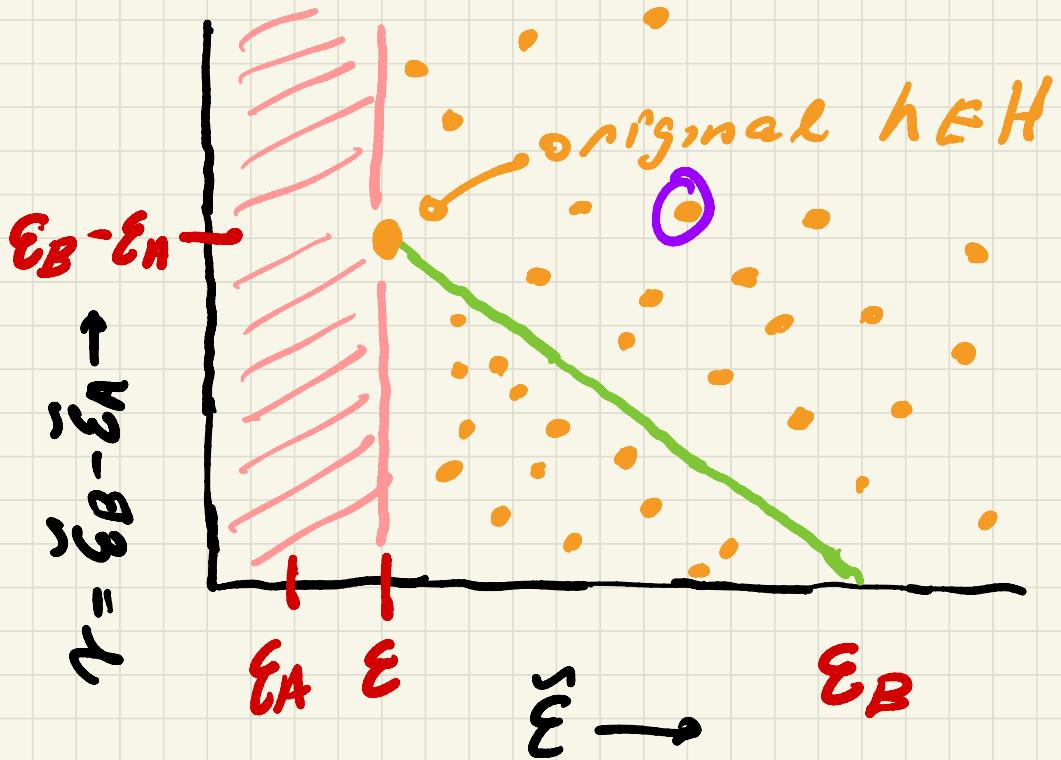
- If we are given a point / model "northeast" of this line, we're better off moving to the line along the perpendicular - line Pareto dominates the point

# Observations

- The bolt-on tradeoff can only improve fairness by deliberately raising the error on the advantaged group
- The bolt-on tradeoff is inherently **linear**

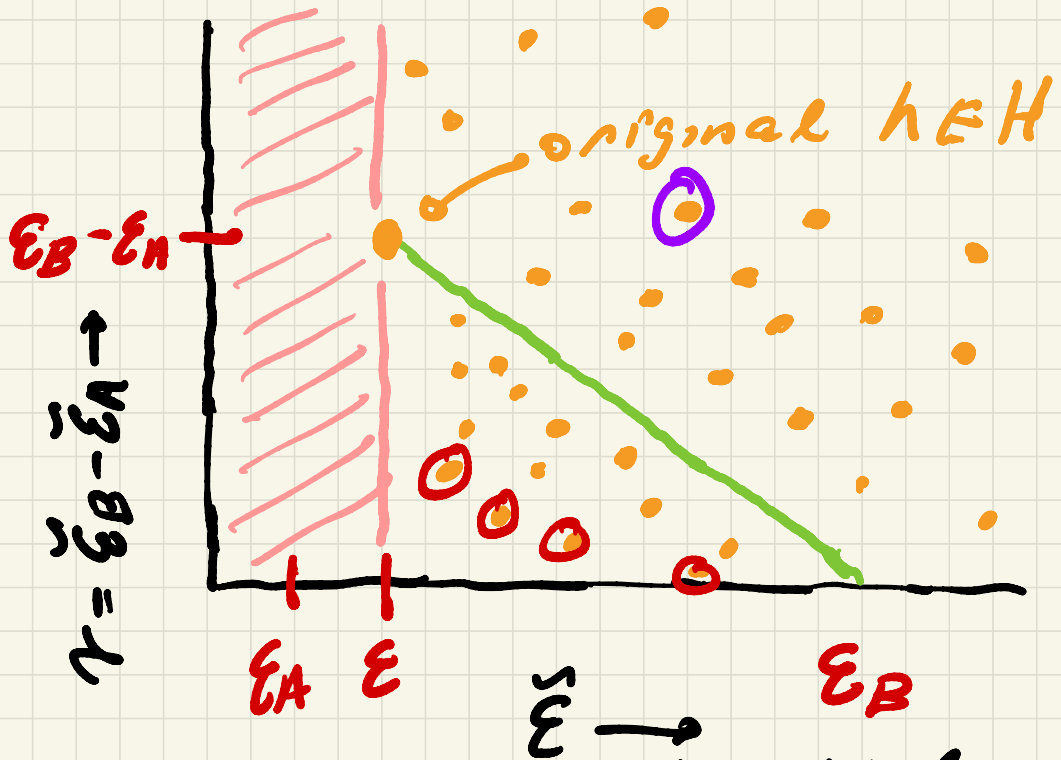
Can we do better?

# The H-Frontier



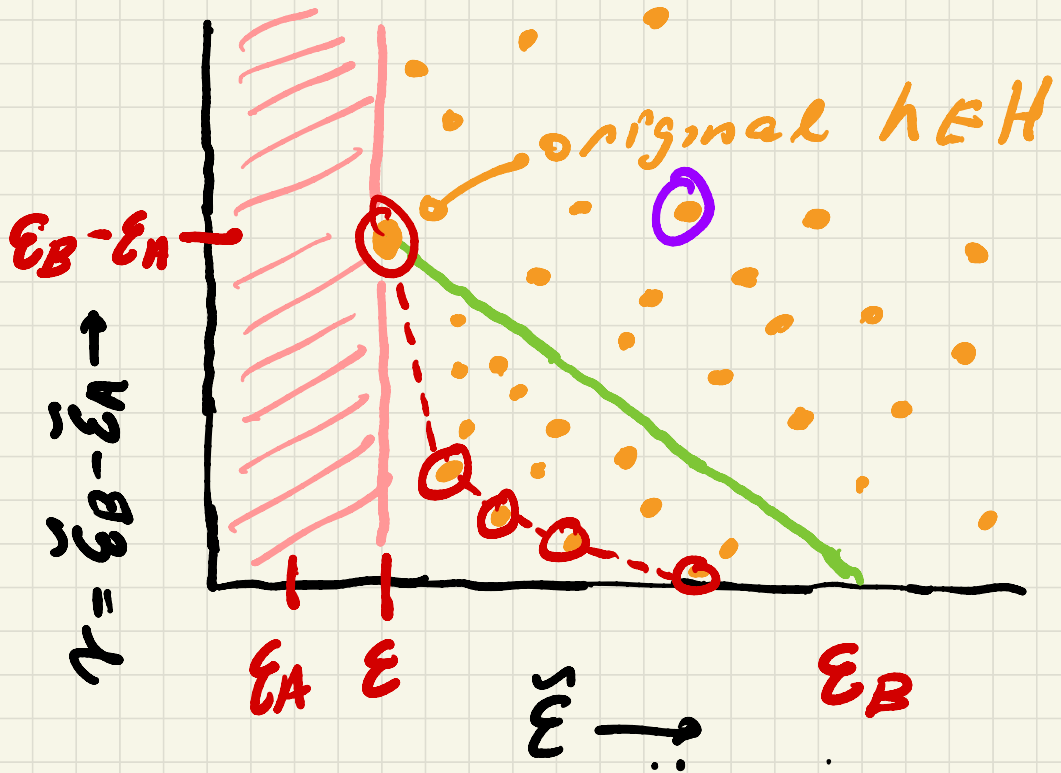
- Imagine we enumerated all the **models/points** in  $H$
- If original has optimal error in  $h$ , no points in **shade**
- And many points might be **Pareto dominated** by **best-on tradeoff**

# The H-Frontier



- But **some points** in  $H$  might lie below the bolt-on line & also have no points to their SW
- These points have better error & fairness than bolt-on & the other models in  $H$

# The H-Frontier



- If we "connect the dots" of these points we have a convex and better tradeoff than bolt-on
- The H-frontier

## Connecting the dots:

• Note that if we define

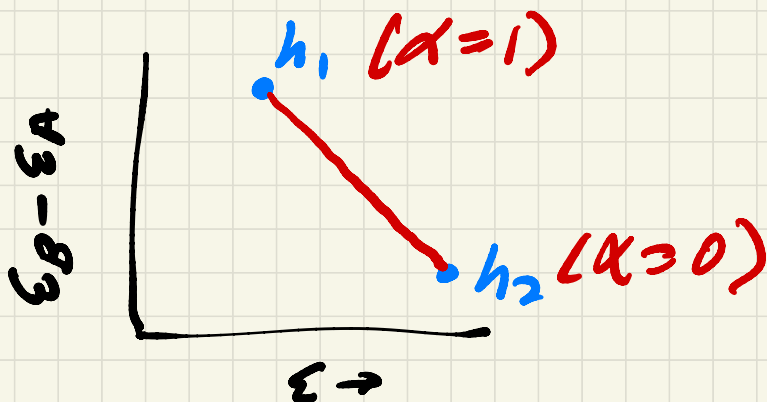
$$g(x) = \begin{cases} h_1(x) & \text{with prob. } \alpha \\ h_2(x) & \text{with prob. } 1-\alpha \end{cases}$$

then  $E(g) = \alpha E(h_1) + (1-\alpha)E(h_2)$

$$E_A(g) = \alpha E_A(h_1) + (1-\alpha)E_A(h_2)$$

similar for B

$$\begin{aligned} \Rightarrow E_B(g) - E_A(g) &= \alpha(E_B(h_1) - E_A(h_1)) \\ &+ (1-\alpha)(E_B(h_2) - E_A(h_2)) \end{aligned}$$





How can we  
find/compute  
the H-frontiers?

Well...

- even finding  $h^* \in H$  is intractable  
In worst case
- but we do have effective non-fair heuristics

# The Reductions/Oracle Approach

• Assume we have a black-box subroutine  $L$  for learning  $h \in \mathcal{H}$  w.r.t.  $\epsilon(h)$  only (non-fair)

• But  $L$  is "pretty good" & general (can solve weighted class. problems in  $\mathcal{H}$ )

Show we can use  $L$  for fair learning.

# Constrained optimization

$$\min_{h \in \Delta(H)} \{ \epsilon(h) \} \text{ s.t.}$$

fairness constraints:

- (1)  $|\epsilon(h, \text{white}) - \epsilon(h, \text{black})| \leq \gamma$
- (2)  $|\epsilon(h, \text{white}) - \epsilon(h, \text{hispanic})| \leq \gamma$
- (3)  $|\epsilon(h, \text{black}) - \epsilon(h, \text{hispanic})| \leq \gamma$
- $\vdots$
- (k) (usually small, but...)

Introduce variables for weights  
in  $\Delta(H)$  & constraints  $\Rightarrow$

huge LP.

- Learner picks  $h \in H$
- Regulator picks a pair of groups, e.g.  $A \neq B$
- Payoff to Regulator:

$$\hat{\epsilon}(h) + \max(0, |\hat{\epsilon}(h, A) - \hat{\epsilon}(h, B)| - \gamma)$$

- Payoff to Learner:  
- (Payoff to Regulator)

# Game Theory Formulation

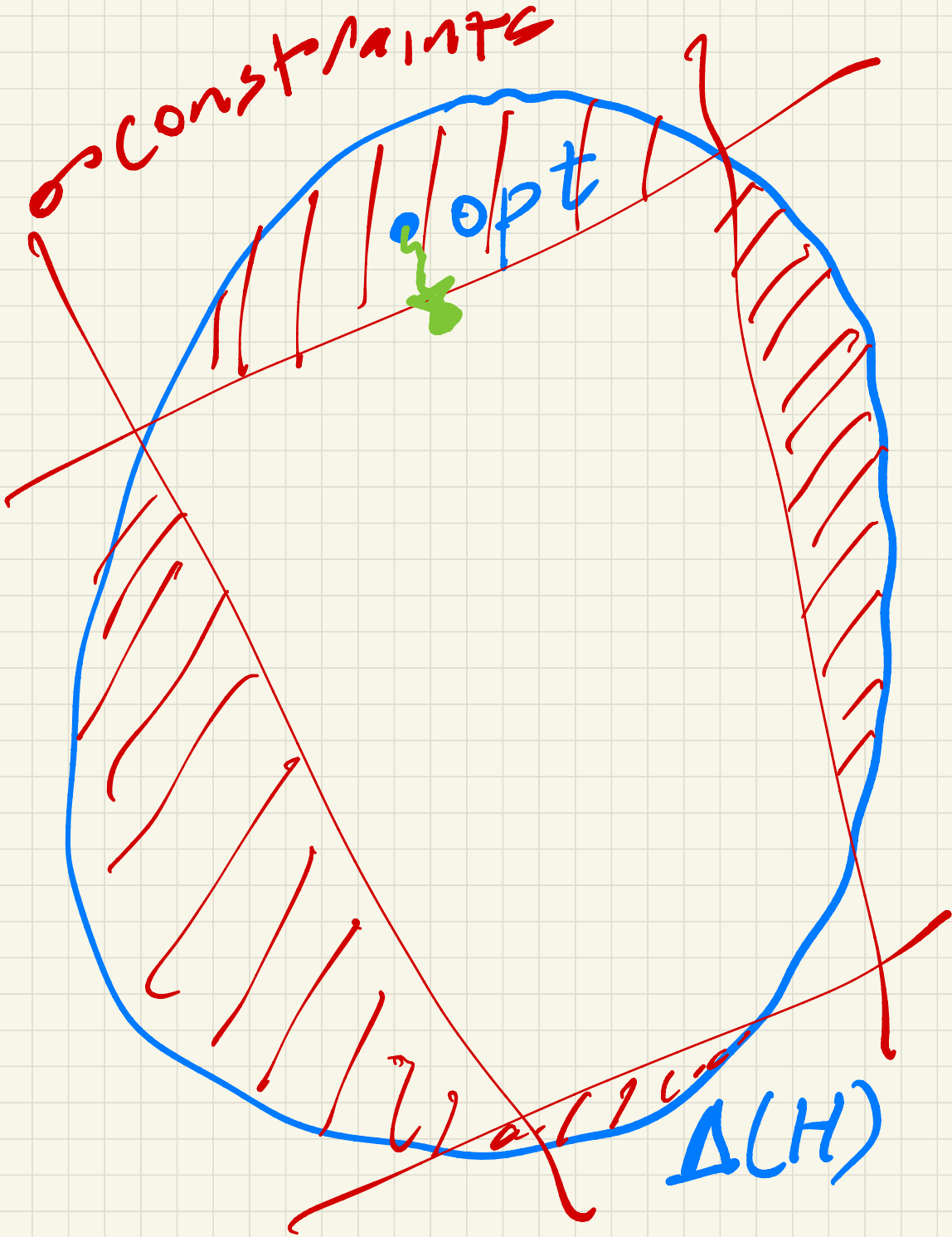
- Learner plays mixed strategy  $p \in \Delta(H)$
- Regulator plays mixed strategy  $q$  over fairness constraints

- Zero-sum game on:

$$u(p) + \text{constraint violations}(p, q)$$

$$= \text{payoff to Regulator} \\ = - \text{payoff to Learner}$$

Nash equil = constrained opt solution



# A Classic Theorem (Freund & Schapire)

If  $L$  &  $R$  play iteratively:

(1)  $L$  best responds to  $g_t$

(2)  $R$  updates  $g_{t+1}$  using  
no-regret algo

Then converge to

$1/\sqrt{t}$ -optimal  
solution.




(2) usually easy

(1) often reduces to  
weighted classification  
with wts. given by  $\delta_t$

$\Rightarrow$  "oracle" L.  
(Agarwal et al.)

---

Yields "principled  
heuristics" that  
are implementable.



Subgroup  
Fairness

- Suppose we ask for group fairness by all of race, gender, disability, age, income, ...

- Might still discriminate against disabled Hispanic women over age 55 making  $\leq 20K/year$

# Framework

- Model class  $H$
- Group membership class  $G$
- For  $g \in G$ ,  $g(x) \in \{0, 1\}$  indicates if  $x$  is in  $g$  (e.g. disabled Hispanic...)
- Now allowing  $G$  to be large or infinite

# Game Theory II

- Learner plays  $h \in H$
- Regulator plays  $g \in G$ ,  
finds **most violated**  
 $g$  (e.g.  $h$  has high  
error on  $g$ )

Reduce to non-fair  
case;  $L$  no-regret,  
 $R$  best response

Another Approach:

Average

Individual

Fairness

- Suppose we will make many decisions about  $x$  over time
- E.g. product rec's
- Then any  $h$  has error rate  $\underline{\epsilon}_x(h)$  across problems
- Ask that all  $\epsilon_x(h)$  be  $\approx$  equal across individuals  $x$
- Game Theory III •

Fairness

Elicitation



- What if fairness isn't "simple"...

- ...but we can elicit empirical fairness judgements.

- E.g.

"Alice & Bob should receive same treatment"

"Alice should be treated at least as well as Bob"

# Framework

- Outcome data  $S = \{(x_i, y_i)\}$
- Fairness data  $F$  of form  $x_i = x_j, x_i \geq x_j$
- Find  $h \in H$  that min's error on  $S$  subject to  $F$
- Generalize to dist's of  $S$  &  $F$
- Game Theory IV •

# Beyond Equalization

- Problem: may achieve by heedlessly inflating harm to advantaged
- Alternative: minimax group fairness:

$$\min_{h \in H} \max_{\text{groups } g} \{ E_g(h) \}$$

- Game Theory V •

$\pi_1, \pi_2 \rightarrow$  stat. dist  $P_1, P_2$   
& some mix times

---

• param class  $H$ ,  $\dim(H)$

•  $E_{P_1}[\underbrace{(h_1(x) - v_1(x))^2}_{\text{same for } P_2}] \leq \epsilon \checkmark$

• cross-train:

$$E_{P_1}[\underbrace{(h_2(x) - v_2(x))^2}_{\text{same for } P_2}] \leq \epsilon$$

4 conditions

Alternative Approach:

"Bias Bounties"

# Motivation

- AI activism & adversarial dynamics
- Can't defend against everything (?)
- Hard to anticipate groups harmed (gerrymandering)
- "Unexpected" data
- Bug bounties in software

# First Idea:

- Invite "crowd" to find problems in a trained model
- E.g. on subsets of training data
- On new data
- "Blurry pictures problem"

## Second Idea

- Invite crowd to propose improvements to a trained model
- Submit two classifiers:
  - group classifier  
 $g(x) \in \{0, 1\}$
  - model classifier  
 $h(x) \in \{0, 1\}$

Idea is that  $h(x)$  is more accurate on  $g(x) = 1$  than current model

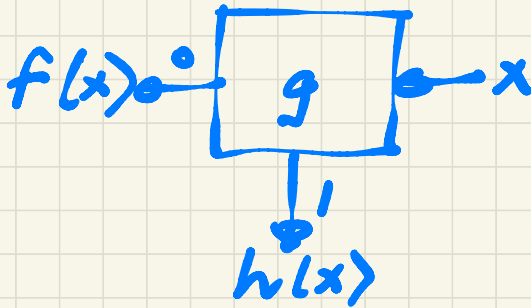


More formally:

- $f(x)$  is current model
- wrt  $P, \epsilon_g(h) \not\leq \epsilon_g(f)$

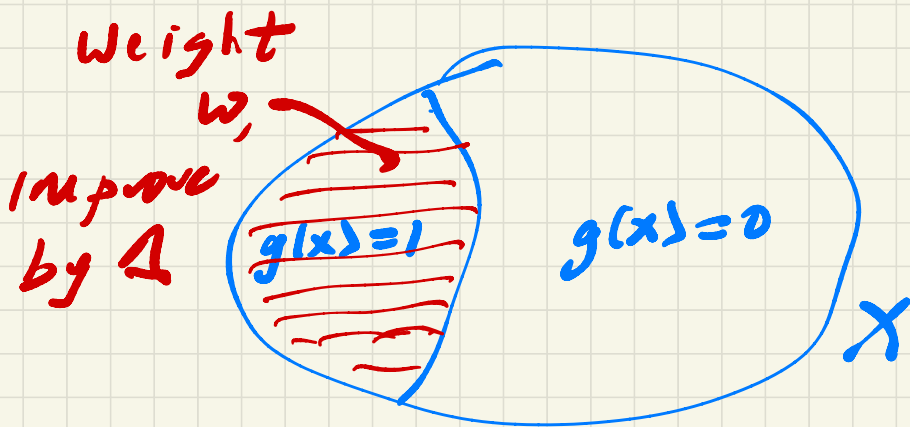
$$P_{x \sim P} [h(x) \neq f^*(x) | g(x) = 1]$$

- What to do?

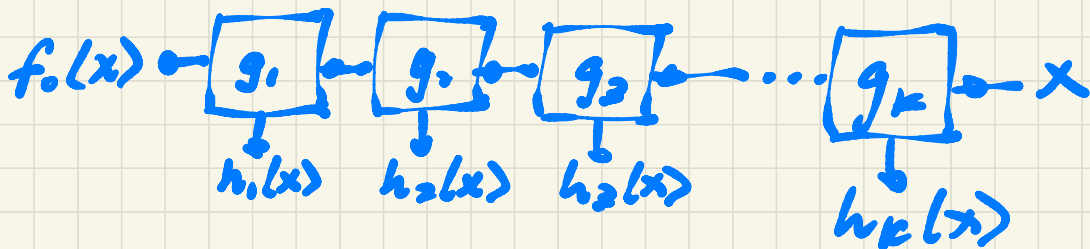


new model  
 $f_1(x)$

- $\mathcal{E}_g(f_1) = \mathcal{E}_g(h) \leq \mathcal{E}_g(f)$
- let  $\Delta = (\mathcal{E}_g(f) - \mathcal{E}_g(f_1)) \geq 0$
- let  $w = P_{x \sim p} [g(x) = 1]$
- $\mathcal{E}(f_1) = (1-w)\mathcal{E}_{\pi_g}(f)$  *no change*  
 $+ w\mathcal{E}_g(f_1)$  *better!*
- $\mathcal{E}(f) - \mathcal{E}(f_1) = w \cdot \Delta \geq 0$

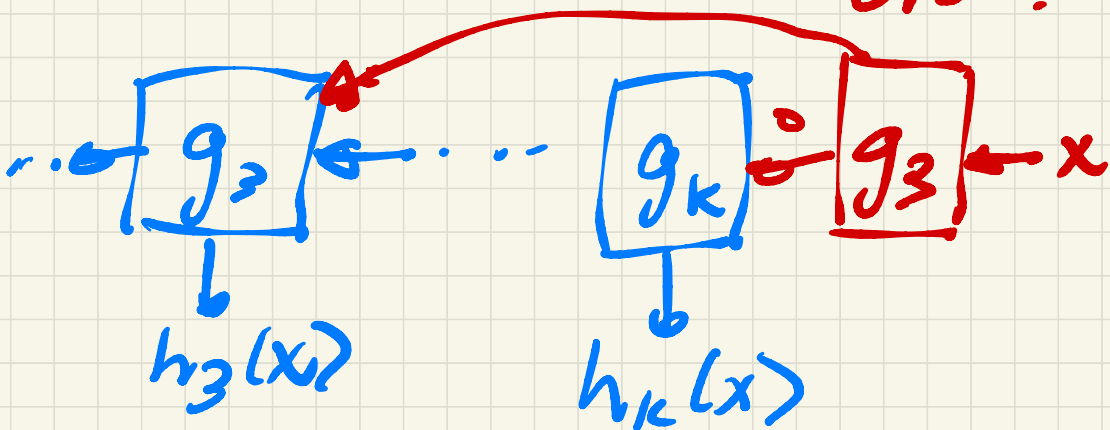


Now repeat/recurse:



Hmm... what if  $g_k$

makes e.g.  $g_3$  worse  
off?



"pointer decision list"  
(PDL)

Ends in one of three ways:

- $E(f_k) = 0$
- No improving  $(g, h)$  exists  
- then  $f_k = \text{Bayes optimal}$
- Nobody can find improving  $(g, h)$   
- effective Bayes optimal

# Practicalities

- Must give crowd **data**, not distribution
- What about overfitting?
- One solution:
  - only accept when  $w \cdot \Delta \geq \gamma \neq 0$
  - give **no info** when  $(g, h)$  rejected
  - can only be  $1/\gamma$  accepts!
  - use test set to accept/reject

# Algorithmic Framework

- Plug in your favorite algo for finding  $(g, h)$  improvements
- Ternary classification
- EM-style algo
- Others?

Trade-offs will return...