# Fairness in Machine Learning

# 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.

$$FN(h, male) \approx FN(h, female)$$

↳ allows for optimization
   of overall error

**Note:** We can achieve = FN rates by randomization.

If individual $x$, predict $\hat{y} = +$ with prob. $p$

If $y = -$, can't be a FN

If $y = +$, $\hat{y} = -$ w.p. $p$

$\therefore FN(p, *) = p$.

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?

# 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**

# Algos for Fair ML:

# Bias Mitigation

# A Post-Processing Approach ("bolt on")

- start with **non-fair** $h(x)$, want to ≈ M/F error rates
- build a **probabilistic classifier on top** of $h(x)$:

$$h(x): \begin{array}{c} + \\ \\ - \end{array} \quad \begin{array}{c|c} M & F \\ \hline P & q \\ \hline r & s \end{array}$$

$$\underbrace{\phantom{\begin{array}{cc} M & F \\ P & q \\ r & s \end{array}}}$$

prob. $\breve{h}(x) = +$

$\breve{h}(x)$

(closed under mixtures)

$$p = q = 1, \; r = s = 0:$$
$$\breve{h} \equiv h, \; \varepsilon(\breve{h}) = \varepsilon(h)$$

$$p = q = r = s = \tfrac{1}{2}:$$
$$\varepsilon(\breve{h}) = \tfrac{1}{2}, \quad \text{perfectly fair}$$

$$p = r = \tfrac{1}{2}, \; q = s = 1:$$
error on men $= \tfrac{1}{2}$
error on women $=$ same as $h$

etc.

Set of all $\langle p, q, r, s \rangle$ gives
Pareto frontier of $h$:



error on men / error on women

$\mathcal{E}(h)$

overall error $\mathcal{E}(\hat{h})$

$1/2$

# Algorithm

- Problem of finding $\check{h}$
  than minimizes $\varepsilon(\check{h})$

$$\text{subject to}$$
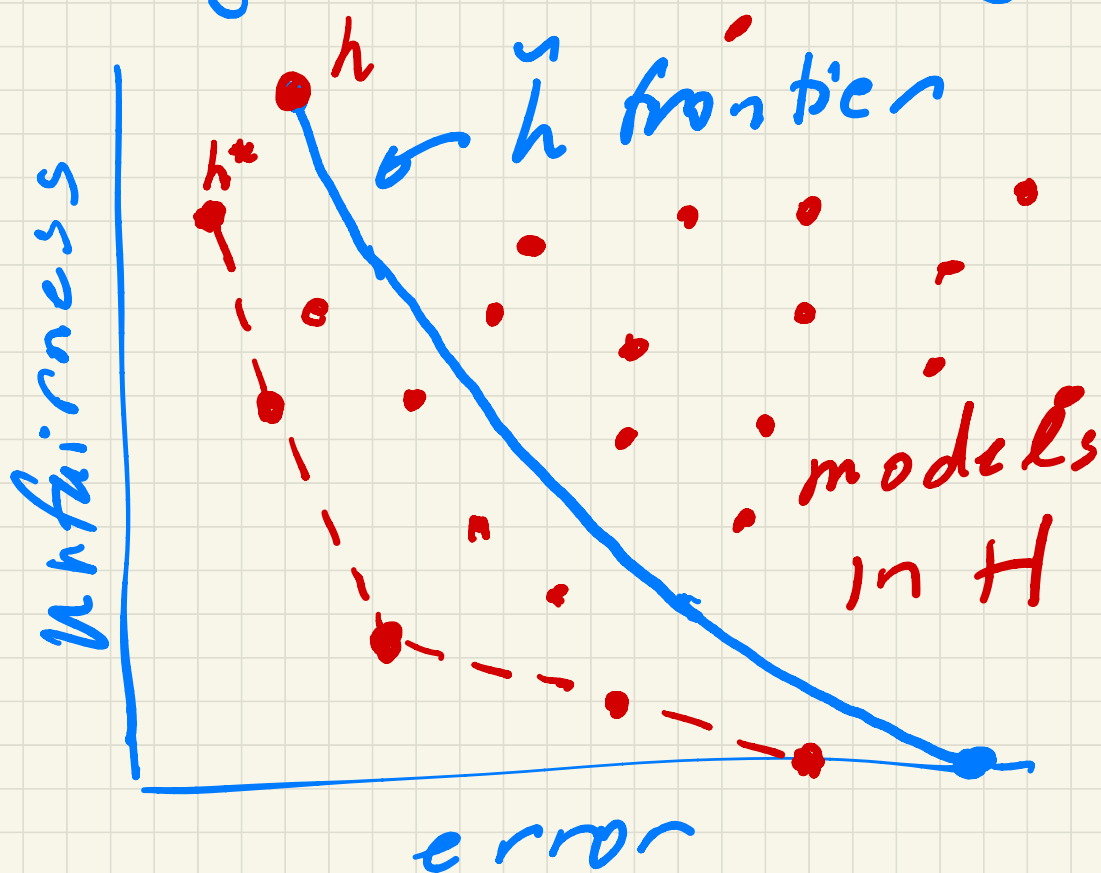
$$y\text{-axis} \leq \gamma$$

is a **linear program**

in $p, q, r, s$.

(Framework & result
due to Hardt, Price,
Srebro.)

# What more could we want?

- Imagine $h \in H$ (NNs, DTs, ...) by some learning algo



$h$

$h$ frontier

$h^*$

Unfairness

models in H

error

Can we find H-frontier?

## 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 $h \in H$ w.r.t. $\varepsilon(h)$ only $\binom{\text{non-}}{\text{fair}}$

- But $L$ is "pretty good" & general (can solve weighted class. problems in $H$)

Show we can use $L$ for fair learning.

# Constrained Optimization

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

fairness constraints:

(1) $|\varepsilon(h, white) - \varepsilon(h, black)| \leq \gamma$

(2) $|\varepsilon(h, white) - \varepsilon(h, hispanic)| \leq \gamma$

(3) $|\varepsilon(h, black) - \varepsilon(h, hispanic)| \leq \gamma$

$\vdots$

(k) (usually small, but...)

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

**huge LP.**

# Game Theory Formulation

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

- Zero-sum game on:

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

$$\underbrace{\phantom{\mathcal{E}(p) + \text{constraint violations}(p,q)}}_{\text{payoff to Regulator}}$$

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

Nash equil = constrained opt solution

# A Classic Theorem (Freund & Schapire)

If L & R play iteratively:

(1) L best responds to $q_t$

(2) R updates $q_{t+1}$ 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 $\beta t$

$\Rightarrow$ "oracle" L.

(Agarwal et al.)

Yields "principled heuristics" that are implementable.

# Towards Individual Fairness

**Q:** Why not treat each individual x as their own "group"?

**A:** Error (or FP, FN,...) "rate" on x is either 0 or 1.

But there are other approaches...

# Metric Fairness

- Posit a distance metric $d(x,x')$ between pairs of individuals

- $h(x)$ our real-valued prediction

- Then constrain $h(x)$ to obey $\forall x, x'$:

$$|h(x) - h(x')| \leq \alpha\, d(x,x')$$

# Difficulties

- Where do we get $d(x, x')$?
- Closed form?
- Usually want to **threshold** $h(x)$, lose fairness
- Practical challenges

# 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 $\varepsilon_x(h)$ across problems

- Ask that all $\varepsilon_x(h)$ be $\frac{1}{2}$ 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 $S = \{\langle x_i, y_i \rangle\}$

- Fairness data $F$ of form $x_i = x_j,\ x_i \gtrsim x_j$

- 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 **needlessly inflating** harm to advantaged

- Alternative: **minimax group fairness:**

$$\min_{h \in H} \max_{\substack{groups \\ g}} \left\{ E_g(h) \right\}$$

- Game Theory $\overline{\mathbb{V}}$ ●

# Other Learning Settings

# Fairness in Bandits

- Ground truth data

$$\langle x, y \rangle$$

loan ↗ ↘ $\in \mathbb{R}$, prob. of
app          repayment

- Unknown linear map
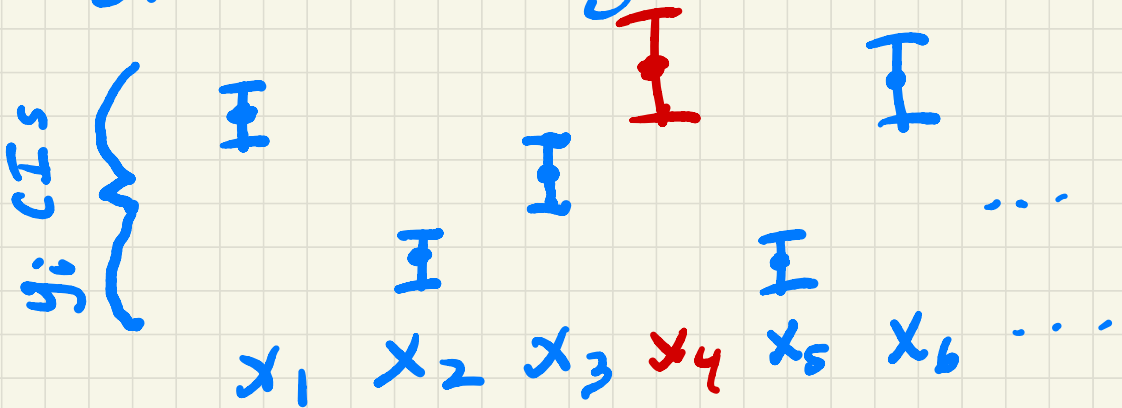
$$y = \theta \cdot x + noise$$
$$(linear\ regress.)$$

- Meritocratic fairness:

If $y_1 \geq y_2$, must have

prob. of    $\geq$    prob. of
loan to $x_1$        loan to $x_2$

- <u>Bandit setting</u>: each day $x_1, \dots x_K$ arrive, must choose loans **fairly**

- Standard algo: LIN-UCB

$y_i$ CIS $\left\{ \phantom{xx} \right.$



$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6 \quad \dots$
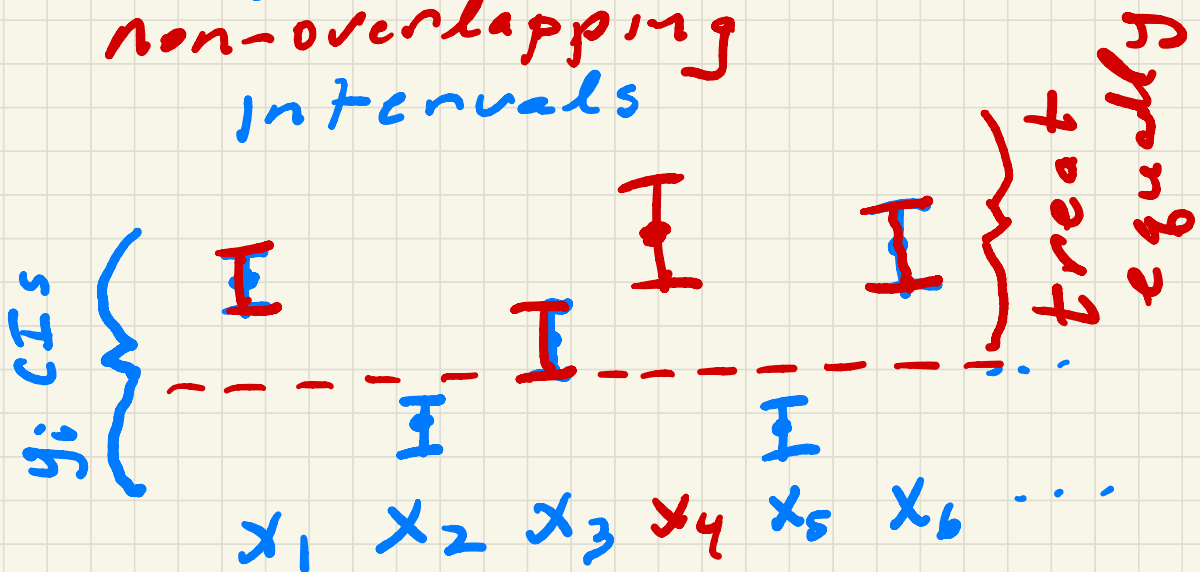
Give loan(s) to highest UCBs $\Rightarrow$ fast convergence to opt

<u>Not fair</u>

# Fair Modification

- **Interval chaining**
- May even choose **non-overlapping** intervals



$y_i$ CIs

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6 \quad \cdots$

treat equally

- choose interval
  $\Rightarrow$ more data
  $\Rightarrow$ chains fragment
  $\Rightarrow$ fast convergence
  $\qquad$ to opt

# Other Topics

- Fair RL
  (e.g. meritocratic
         wrt Q-values)

- Fair Representations

- Causal Approaches

- Fair Clustering

- Fair Rankings

  :

# Some Resources

- "Frontiers of Fairness in Machine Learning"
  Chouldechova & Roth

- "Fairness and ML"
  Barocas, Hardt, Narayanan
  fairmlbook.org

- "The Ethical Algorithm"
  Kearns & Roth
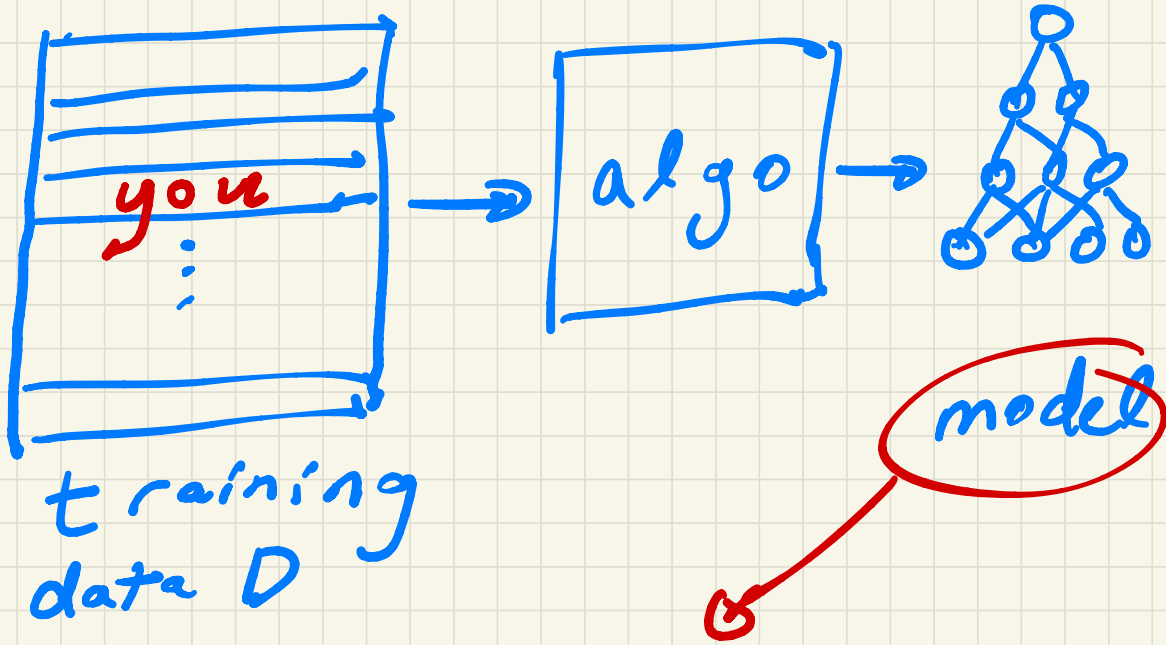
# Privacy in ML

# What Do We Want?

- **Not** addressing preventing data breaches, unwanted access, etc — domains of **cryptography** and **security**

- Rather, prevent **inferences and exfiltration** from trained model

# (Bad) Examples

- K-NN models
- SVMs
- Neural Networks
- Any model with confidence ratings

---

- Even **black-box** access problematic
- "Anonymizing" data **doesn't work**

# High-Level Idea



you

training data $D$

algo

model

Shouldn't reveal "anything" about your data - even with additional computation & data

# Differential Privacy

Say algo $A$ is $\varepsilon$-DP if
$\forall$ neighboring $D, D'$
$\forall$ set $S \subseteq range(A)$:

$$Pr[A(D') \in S] \le e^{\varepsilon} Pr[A(D) \in S]$$

$\hookrightarrow$ wrt randomization
of $A$ only



$D \searrow \swarrow D'$

$A$

$\leftarrow range(A) \rightarrow$