

$$\text{Min } \text{loss}(h_\theta(x), y)$$

s.t.

$$\text{Unfairness} \leq \epsilon, \text{ for some Notions of fairness}$$

→ is not Convex

↳ how to make it

Convex

more efficient

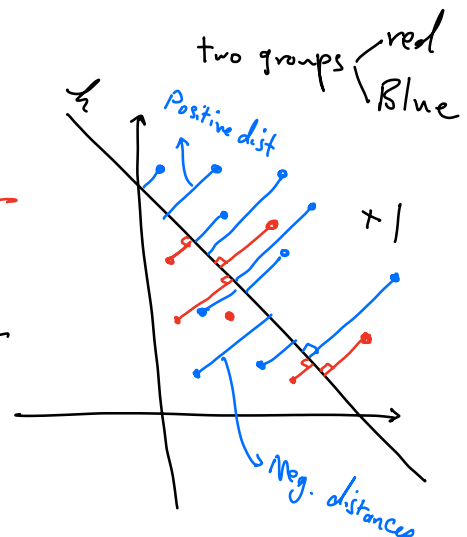
[Zafar, et al.]

The overall idea: Change the notion of fairness to make the opt. Convex

Binary Classification
Binary Sensitive Attribute
Fairness Notion: Demographic Parity

$$\text{rel-ratio} = \frac{2}{6}$$

$$\text{Blue-ratio} = \frac{7}{10}$$



Fairness: The distance of points from the decision Boundary based on

Independence: $S \perp\!\!\!\perp h_\theta(x)$

$$\text{Cov}(S, h_\theta) = 0$$

$\text{Cov}(S, d_\theta(x))$ is bounded

↳ distance to the decision boundary

$$|\text{Cov}(S, d_\theta(x))| \leq \tau$$

$$\begin{aligned} \text{Cov}(S, d_\theta(x)) &= E[(S - \bar{S}) d_\theta(x)] - E[S - \bar{S}] E[d_\theta(x)] \\ &= E[(S - \bar{S}) d_\theta(x)] \end{aligned}$$

for Linear classifiers

$$\approx \frac{1}{n} \sum_{i=1}^n (S_i - \bar{S}) d_\theta(x_i) \quad ?$$

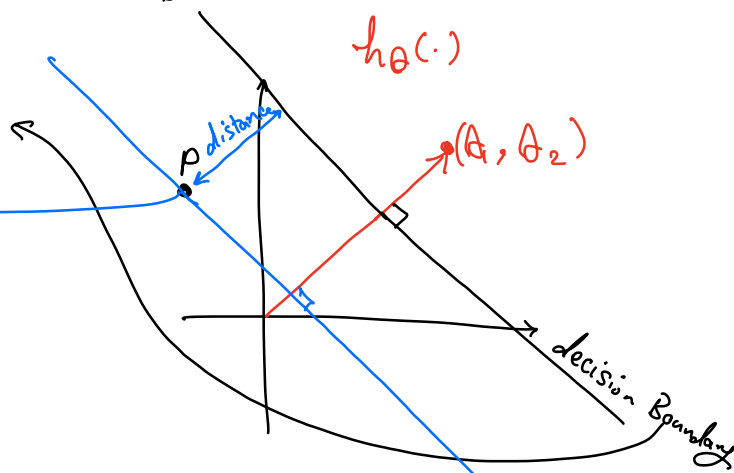
$$h_\theta(x) = \theta^T x = \sum_{j=0}^m \theta_j x_{ij}$$

$$\theta_1 x_1 + \theta_2 x_2 = -\theta_0$$

$$\Rightarrow \theta_0 + \theta_1 x_1 + \theta_2 x_2 = 0$$

$$\Rightarrow \theta^T x = 0$$

$$\Rightarrow \theta^T x^p = d_\theta(P)$$



$$\Rightarrow \text{Cov}(S, d_{\theta}(x)) \approx \frac{1}{n} \sum_{i=1}^n (s_i - \bar{s}) \theta^T X_i$$

$$|\text{Cov}(S, d_{\theta}(x))| \leq \tau$$

$$\Rightarrow \frac{1}{n} \sum (s_i - \bar{s}) \theta^T X_i \leq \tau$$

F and

$$\frac{1}{n} \sum (\bar{s} - s_i) \theta^T X_i \leq \tau$$

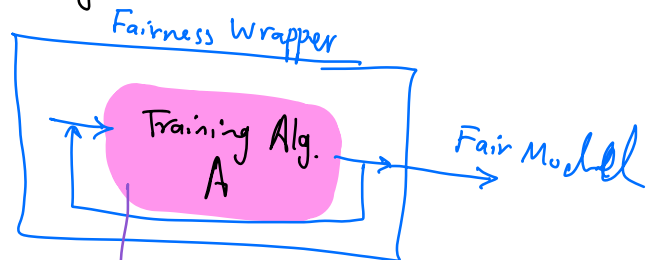
$$\min_{\theta} \text{loss}(h_{\theta}(x), y)$$

s.t.

$F \leftarrow \text{Convex.}$

[Alek Agarwal et al.] (A reduction approach to Fair Classification)

→ They Create A Fairness Wrapper on top of Black Box Training Alg.



⇒ Should work on Several Fairness notions

- "Non-Binary" Sensitive Attributes
- Not Limited to Binary classification

requirement → The Training Alg. should accept weighted Samples.

Step 1: Propose a unified Theme for fairness definitions.

✓ group j & definition F ,

$$\mu_j(\theta) = E[g(x, y, s, \theta) | \varepsilon(x, y, s)]$$

Performance

don't Support Sufficiency
SUFF

e.g., DP: $g_j = f_\theta(x)$, $\varepsilon_j = \{s = s_j\}$

$$\mu_j(\theta) = E[f_\theta(x) | s = s_j]$$

$$\mu_j^*(\theta) = E[f_\theta(x)]$$

$$\mu_j(\theta) - \mu_j^*(\theta) \leq \tau$$

for a set of groups & defi:

$$\mu_j(\theta) - \mu_j^*(\theta) \leq \tau$$

$$\dots$$

$$\mu_k(\theta) - \mu_k^*(\theta) \leq \tau$$

$$\begin{matrix} \text{Matrix} & \text{Vector} \\ \curvearrowright & \curvearrowright \end{matrix} M_\mu(\theta) \leq \tau$$

$$\Rightarrow M_\mu - \tau \leq 0$$

Step 2: Optimization
Reformulation

$$\text{Min } \ell(\theta)$$

$$\text{s.t. } M_\mu - \tau \leq 0$$

Optimize

$$L(\theta, \lambda) = \ell(\theta) - \lambda(M_\mu - \tau)$$

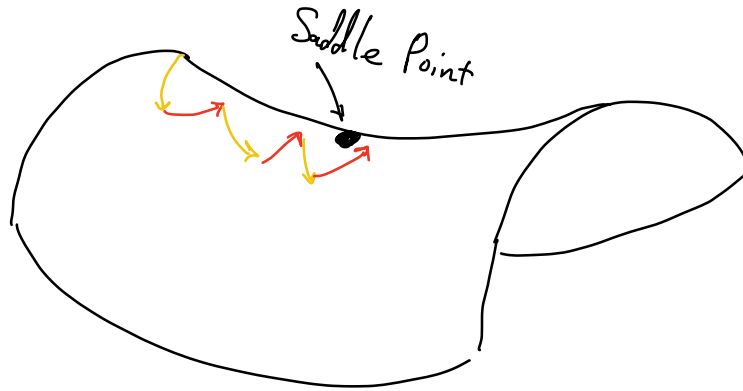
Review: Lagrangian Multipliers
= Dual Form

$$\begin{matrix} \text{Max } f_\theta \\ \text{s.t. } ax \leq b \end{matrix} \Rightarrow L(\theta, \lambda) = f_\theta - \lambda(ax - b)$$

optimize for L

$$\text{Max } f_\theta - \lambda(ax - b)$$

max ← f_θ → min
variables ← $\lambda(ax - b)$



Opt $L(\theta, \lambda)$

$$\max_{\lambda} \min_{\theta} h_{\theta}(x) - \lambda (M_{\theta} - \tau)$$

→ Reweighting

Existing Training Alg.

The # iterations is bounded by $\Theta(n^2)$

Price (Aka Cost) of Fairness

The diff. b/w the opt. value of unfair vs. Fair Solutions.

e.g.,

(A)

$$\min L(\theta) = L_A$$

(B)

$$\begin{aligned} \min L(\theta) &= L_B \\ \text{s.t.} & \text{Fairness} \end{aligned}$$

$$L_B \geq L_A \Rightarrow \text{Price of Fairness} = L_B - L_A$$

Pareto-Optimality of Fairness

Suppose there are groups $\{g_1, \dots, g_K\}$.

the Performance of the Alg/model for g_i is Per_i .

If Per_i^U is the Per. of unFair Alg. for g_i
and Per_i^F is the Per. of Fair Alg. for g_i

If $\forall g_i ; Per_i^U \geq Per_i^F$

$\exists g_j : Per_j^U > Per_j^F$ and

\Rightarrow the Fair Solution is **DOMINATED**
By the Unfair one.

Pareto-Optimal is the Set of Solutions
that are Not Pareto-dominated

\rightarrow AKA skyline