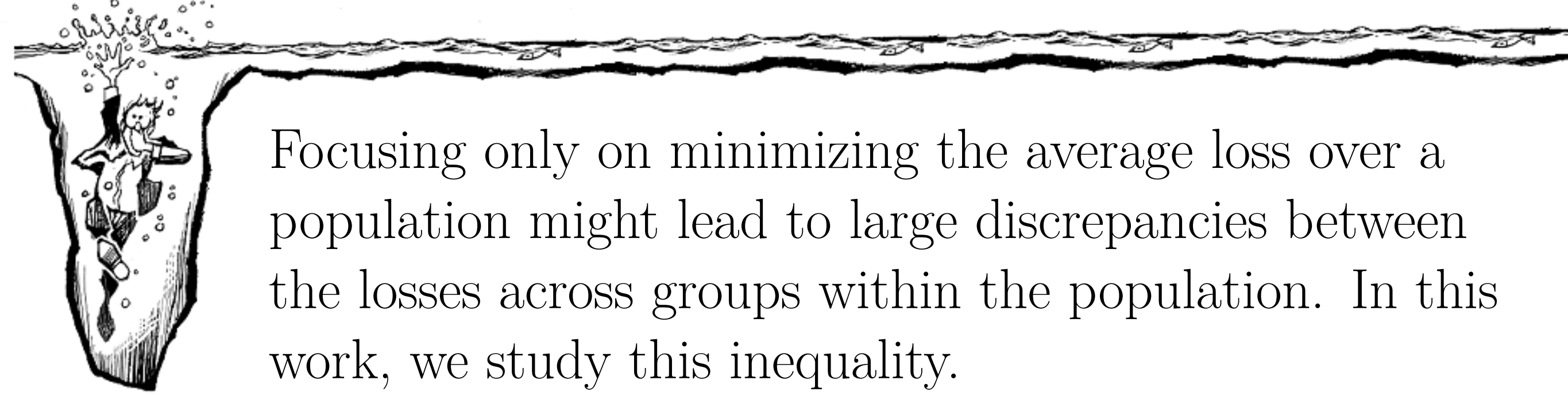


# Maximum Weighted Loss Discrepancy

Fereshte Khani, Aditi Raghunathan, Percy Liang

## Motivation



## Setup

- individuals:  $z = (x, y)$ , underlying distribution:  $p^*$
- group: a measurable function  $g : \mathcal{X} \times \mathcal{Y} \rightarrow \{0, 1\}$ , all groups:  $\mathcal{G}$
- predictor:  $h : \mathcal{X} \rightarrow \mathcal{Y}$
- bounded measurable loss function :  $\ell(h, z)$
- population loss:  $\mathbb{E}[\ell]$ , group loss:  $\mathbb{E}[\ell \mid g = 1]$

## Maximum Weighted Loss Discrepancy (MWLD)

$$\text{MWLD}(w, \ell, h) \stackrel{\text{def}}{=} \sup_{g \in \mathcal{G}} w(g) |\mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell]|$$

## Connection to Group Fairness

MWLD can be viewed as an upper bound for the loss of any group:

$$\mathbb{E}[\ell \mid g = 1] \leq \mathbb{E}[\ell] + \frac{\text{MWLD}(w, \ell, h)}{w(g)}.$$

- Existing statistical notions of fairness can be viewed as enforcing small  $\text{MWLD}(w)$  for different weighting functions.
- Equalized opportunity:** weighting function is 1 on sensitive groups (e.g., different races) and 0 on all other groups.

## Connection to AI safety (Robustness)

MWLD can be viewed as an upper bound for the loss on a population with shifted demographics:

$$\mathbb{E}_{z \sim q}[\ell] \leq \mathbb{E}_{z \sim p^*}[\ell] + \text{MWLD}(w, \ell, h)$$

where  $q(\cdot) \stackrel{\text{def}}{=} w(g)p^*(\cdot \mid g = 1) + (1 - w(g))p^*(\cdot \mid g = 0)$ .

- This is similar in spirit to distributionally robust optimization (DRO) using a max-norm metric.
- The difference is that the mixture coefficient is group-dependent.



## Questions?



- For what weighting functions we can estimate  $\text{MWLD}(w)$ ?

$$\forall \epsilon, \delta \in (0, \frac{1}{2}) : \mathbb{P} \left[ \left| \text{MWLD}(w) - \widehat{\text{MWLD}}_n(w) \right| \geq \epsilon \right] \leq \delta$$

- When can we compute  $\widehat{\text{MWLD}}_n(w)$  efficiently?
- Is there any connection between  $\text{MWLD}(w)$  and other notions?

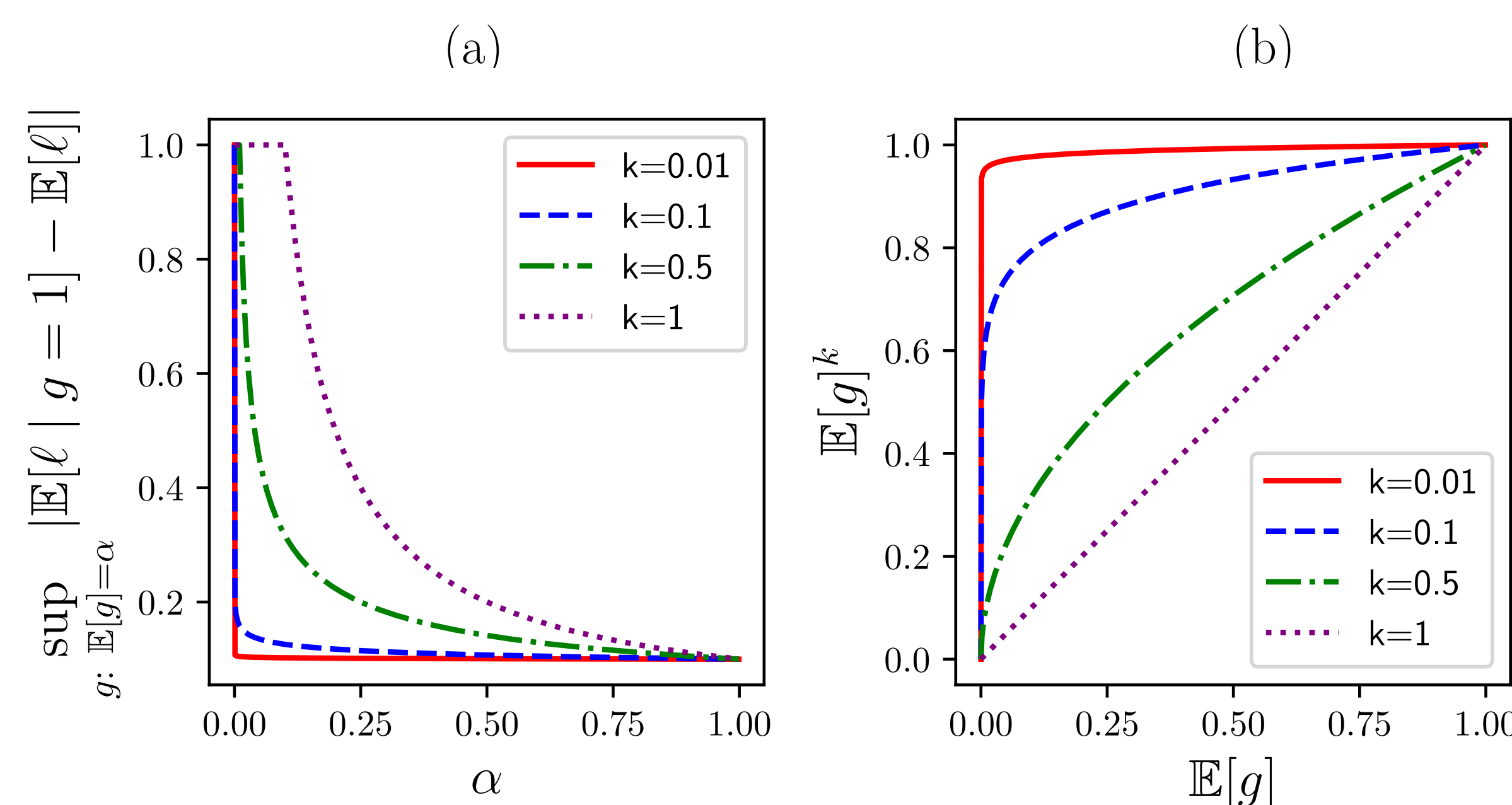
## Proposition

Let  $w^0(g) \stackrel{\text{def}}{=} \mathbb{I}[\mathbb{E}[g] > 0]$ . For non-degenerate  $(\ell, h)$ , it is impossible to estimate  $\text{MWLD}(w^0, \ell, h)$ .

## Theorem

For  $k \in (0, 1]$ , let  $w^k \stackrel{\text{def}}{=} \mathbb{E}[g]^k$ . Given  $n \geq \frac{C \log(1/\delta)}{\epsilon^{2+k}}$  i.i.d. samples from  $p^*$ , we can estimate  $\text{MWLD}(w^k)$  efficiently.

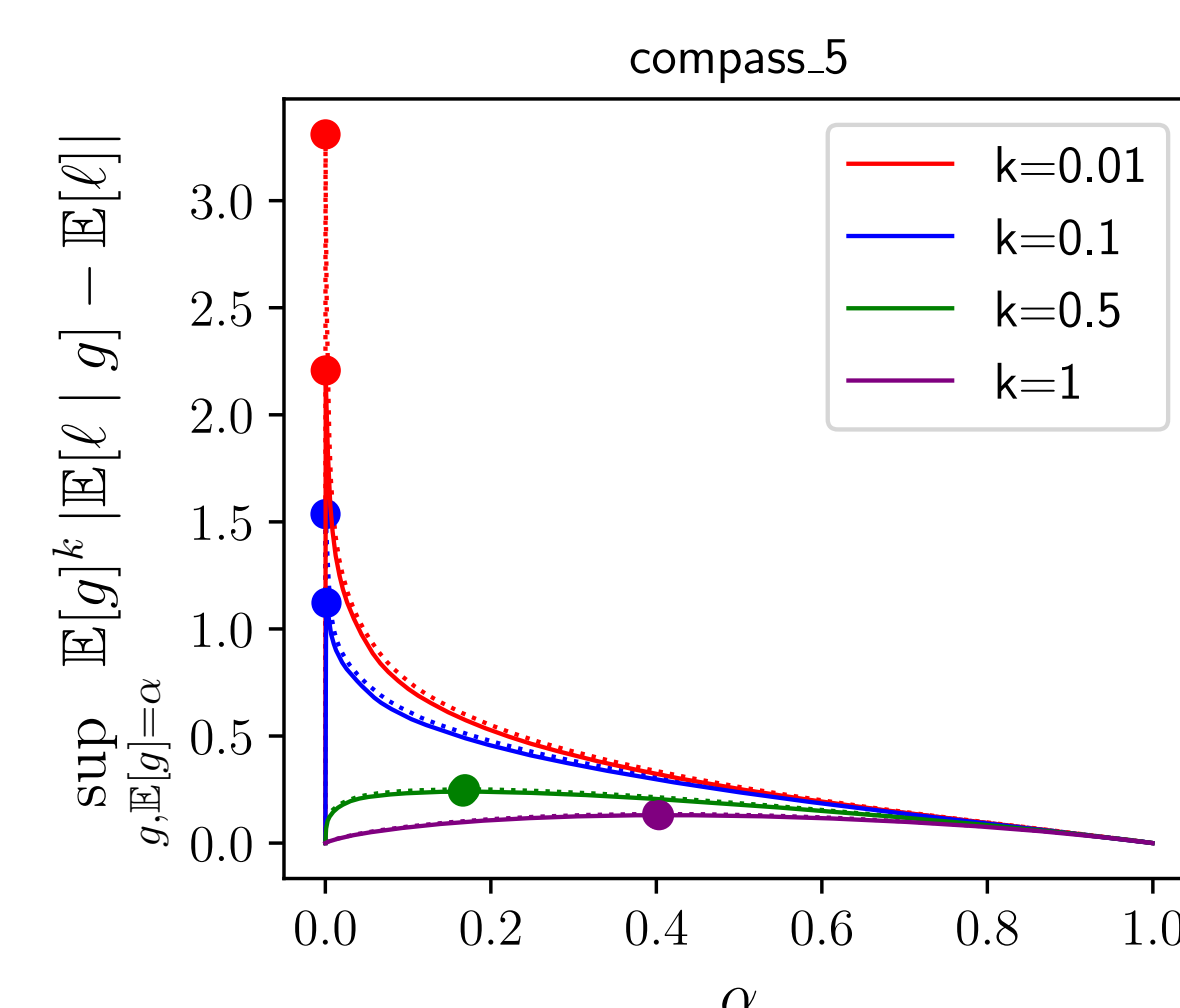
$$w^k(g) = \mathbb{E}[g]^k$$



The parameter  $k$  governs variation of (a) upper bound on loss discrepancy, (b) up-weighting factor across group sizes.

## Estimating $\text{MWLD}(w^k)$ on real world datasets

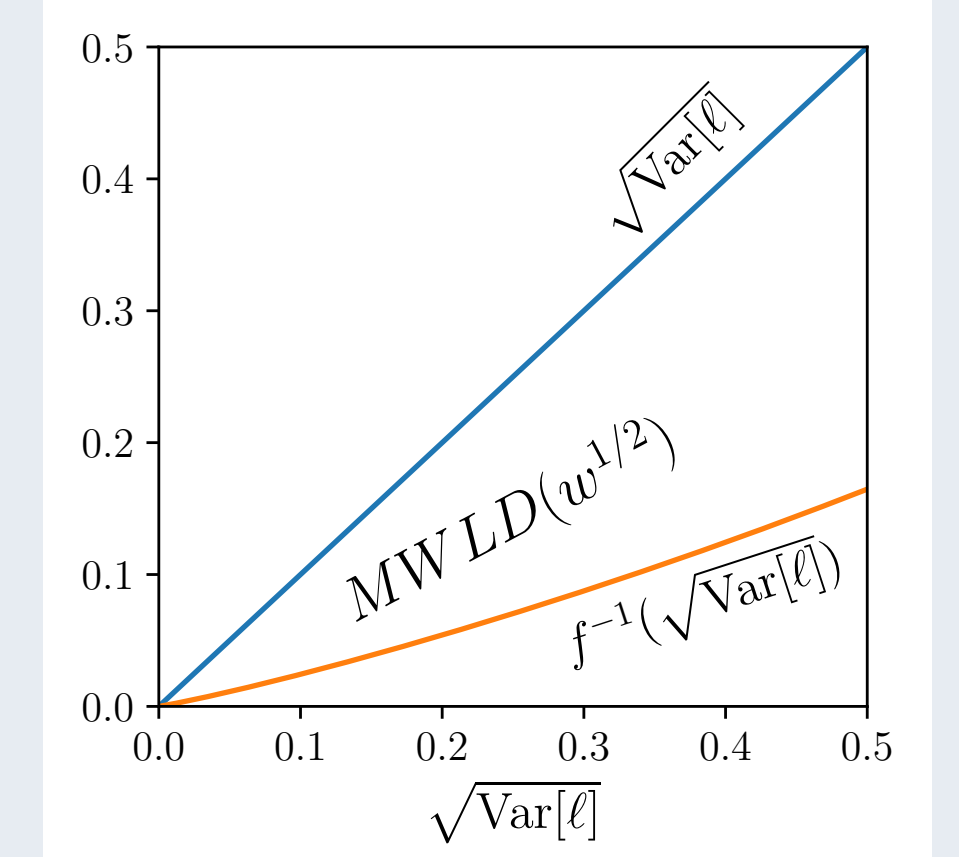
- Dots denote  $\widehat{\text{MWLD}}(w^k)$ .
- For smaller  $k$ , there is a bigger gap between values of  $\widehat{\text{MWLD}}(w^k)$  corresponding to train (dashed lines) and test (solid lines) set.



## Loss Variance

### Proposition

For  $f(x) = x\sqrt{2 - 4 \ln(x)}$ :



- Average individual discrepancy

$$\text{Var}[\ell] = \mathbb{E}[(\ell(z) - \mathbb{E}[\ell])^2]$$

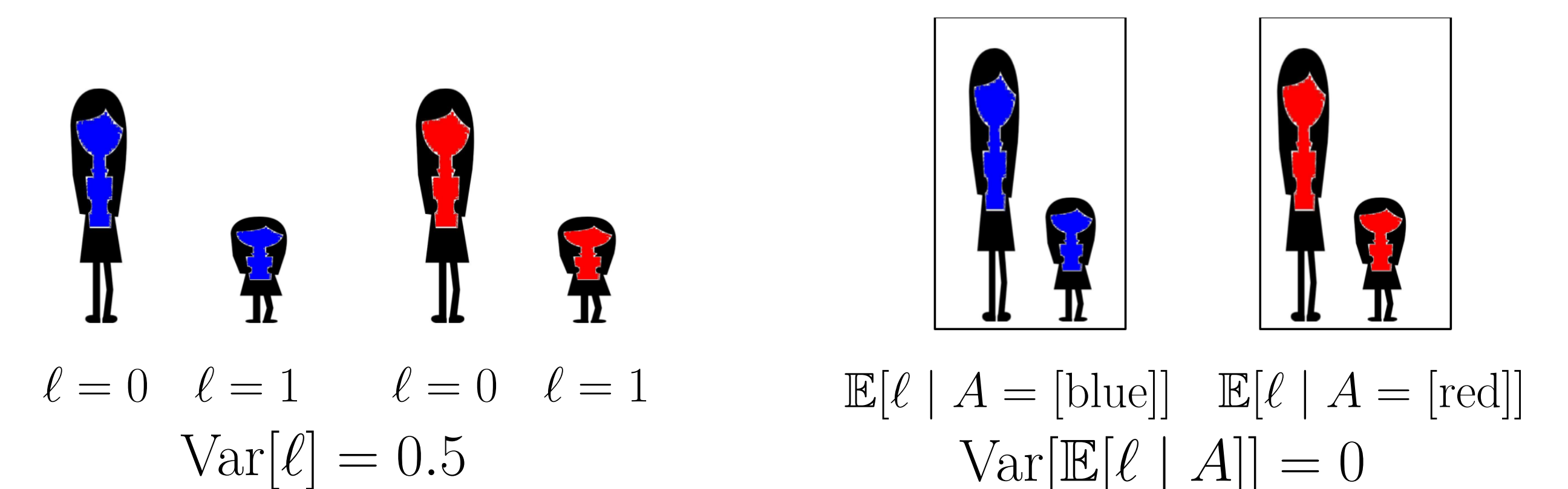
- Generalization bounds

$$|\hat{\mathbb{E}}[\ell] - \mathbb{E}[\ell]| \leq C_1 \sqrt{\frac{\hat{\text{Var}}[\ell]}{n}} + \frac{C_2}{n}$$

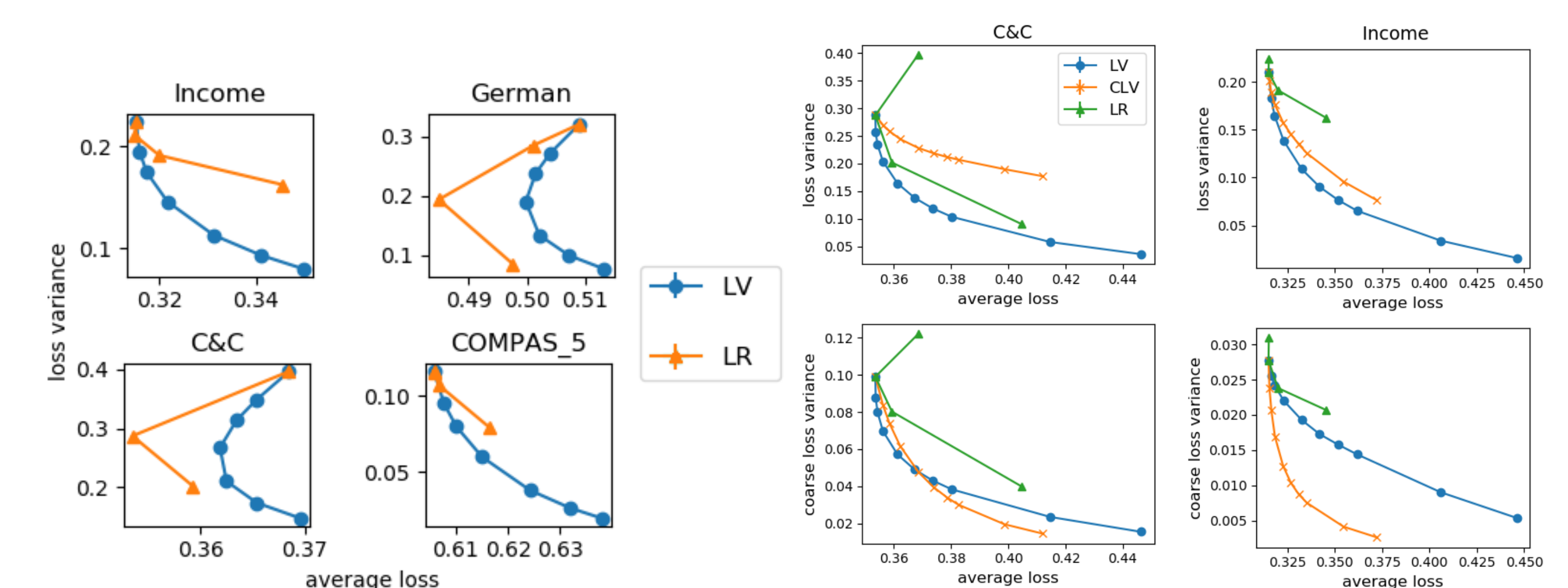
## Coarse Loss Variance

Let  $A$  denote the sensitive attributes; e.g.,  $A = [\text{race}, \text{gender}, \dots]$ .

$$\text{Var}[\mathbb{E}[\ell \mid A]] = \mathbb{E}[(\mathbb{E}[\ell \mid A] - \mathbb{E}[\ell])^2].$$



## (Coarse) Loss Variance Regularization



$$\mathcal{O}_{LR} \stackrel{\text{def}}{=} \hat{\mathbb{E}}[\ell] + \eta \|\theta\|_2^2$$

$$\mathcal{O}_{LV} \stackrel{\text{def}}{=} \mathcal{O}_{LR} + \lambda \hat{\mathbb{E}}[\hat{\text{Var}}[\ell \mid y]] \quad \mathcal{O}_{CLV} \stackrel{\text{def}}{=} \mathcal{O}_{LR} + \lambda \hat{\mathbb{E}}[\hat{\text{Var}}[\mathbb{E}[\ell \mid A, y] \mid y]]$$

- We halve the loss variance with only a small drop in the average loss.
- In some cases, using loss variance as a regularizer simultaneously reduces the classification loss and loss variance.