

# Fairness via Loss Variance Regularization



Fereshte Khani



Aditi Raghunathan



Percy Liang

## Microsoft, Google Beat Humans at Image Recognition

Deep learning algorithms compete at ImageNet challenge

By R. Colin Johnson, 02.18.15 □ 14

## Microsoft, Google Beat Humans at Image Recognition

Deep learning algorithms compete at ImageNet challenge

By R. Colin Johnson, 02.18.15 □ 14

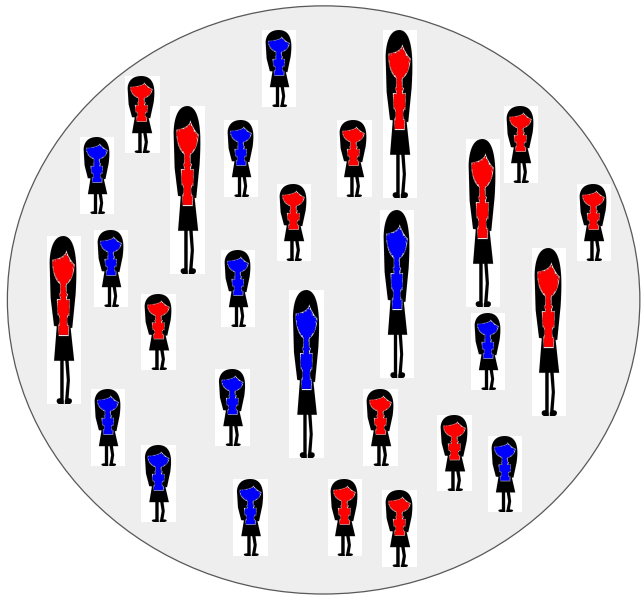
---

The New York Times

---

### *Facial Recognition Is Accurate, if You're a White Guy*

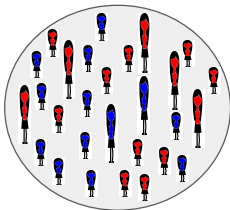
By Steve Lohr



$$\mathbb{E}[\ell] = 0.3$$

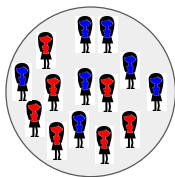
Sensitive attributes

$A = [\text{Height}, \text{Color}]$

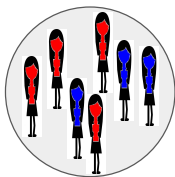


$$\mathbb{E}[\ell] = 0.3$$

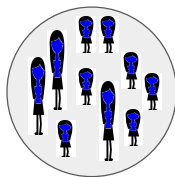
short



tall



blue

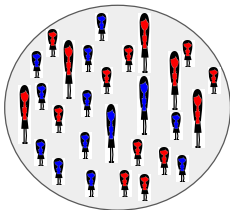


red



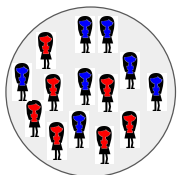
Sensitive attributes

$A = [\text{Height}, \text{Color}]$



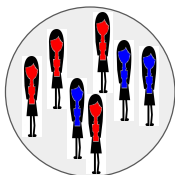
$$\mathbb{E}[\ell] = 0.3$$

short



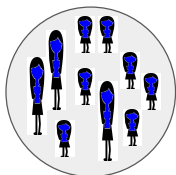
$$\mathbb{E}[\ell] = 0.3$$

tall



$$\mathbb{E}[\ell] = 0.3$$

blue



$$\mathbb{E}[\ell] = 0.3$$

red

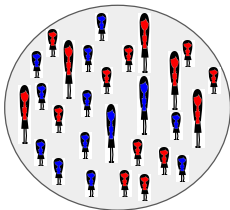


$$\mathbb{E}[\ell] = 0.3$$

**FAIR**

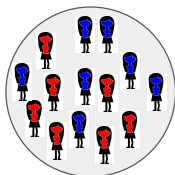
Sensitive attributes

$A = [\text{Height, Color}]$



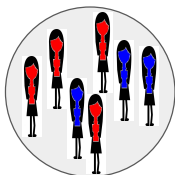
$$\mathbb{E}[\ell] = 0.3$$

short



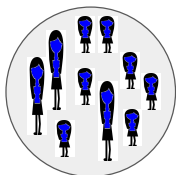
$$\mathbb{E}[\ell] = 0.3$$

tall



$$\mathbb{E}[\ell] = 0.3$$

blue



$$\mathbb{E}[\ell] = 0.3$$

red



$$\mathbb{E}[\ell] = 0.3$$

**FAIR**

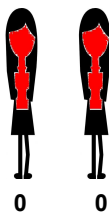
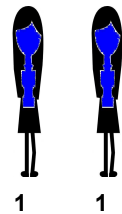
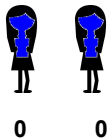
$$\mathbb{E}[\ell] = 0$$

**UNFAIR!**

$$\mathbb{E}[\ell] = 0.5$$

$$\mathbb{E}[\ell] = 0$$

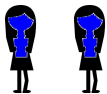
$$\mathbb{E}[\ell] = 0.5$$





**FAIR**

Loss blue: 0.5



0

0

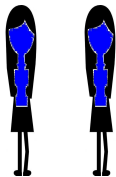
**FAIR**

Loss red: 0.5



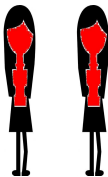
1

1



1

1

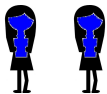


0

0

**FAIR**

Loss blue: 0.5



0

0

**FAIR**

Loss red: 0.5

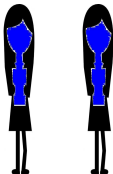


1

1

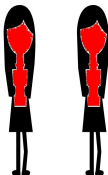
**FAIR**

Loss short: 0.5



1

1



0

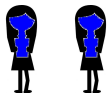
0

**FAIR**

Loss tall: 0.5

**FAIR**

Loss blue: 0.5



0

0

**FAIR**

Loss red: 0.5



1

1

**FAIR**

Loss short: 0.5

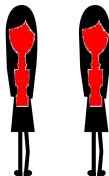


1

1

**FAIR**

Loss tall: 0.5

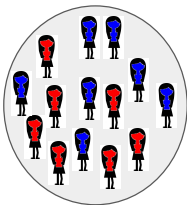


0

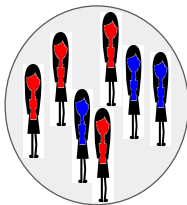
0

$G_{sub}$

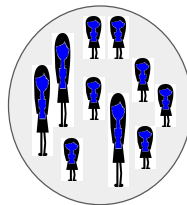
Short



Tall



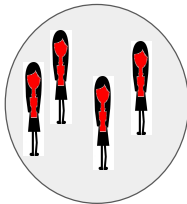
Blue



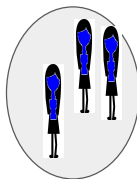
Red



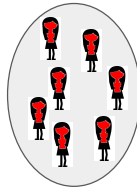
Red tall



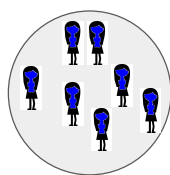
Blue tall



Short red



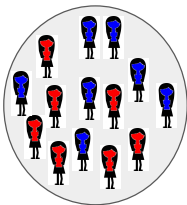
Short blue



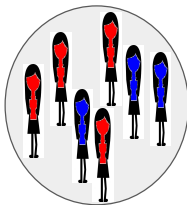
$G_{sub}$ 

$$VC(G_{sub}) \leq C$$

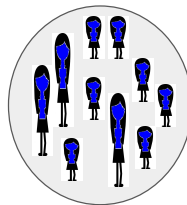
Short



Tall



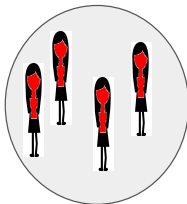
Blue



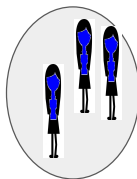
Red



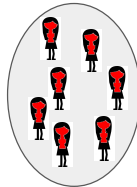
Red tall



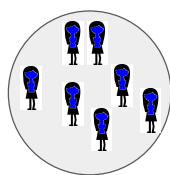
Blue tall



Short red



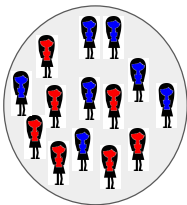
Short blue



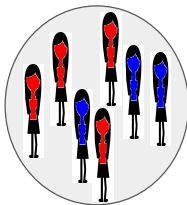
$G_{sub}$ 

$$VC(G_{sub}) \leq C$$

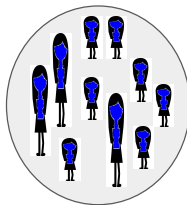
Short



Tall



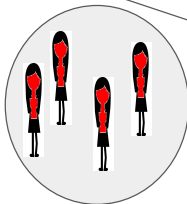
Blue



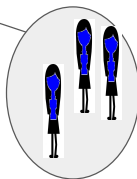
Red



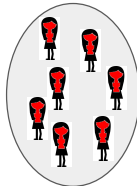
Red tall



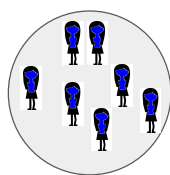
Blue tall



Short red



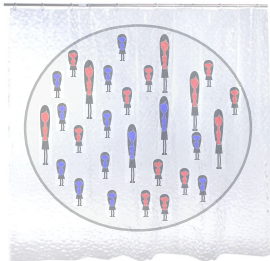
Short blue



Smaller groups  
can have higher  
loss

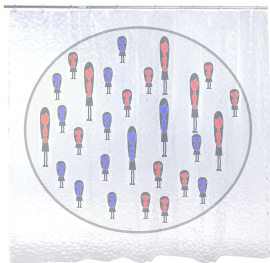
No sensitive attributes  
is available

$A = [?]$



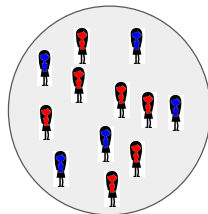
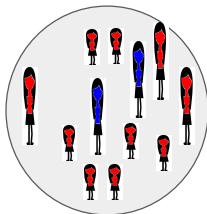
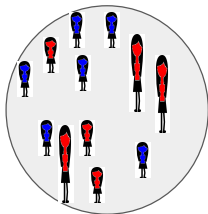
$$\mathbb{E}[\ell] = 0.3$$

No sensitive attributes  
is available  
 $A = [?]$



$$\mathbb{E}[\ell] = 0.3$$

All groups  
larger than  $\alpha$





# Maximum weighted loss discrepancy

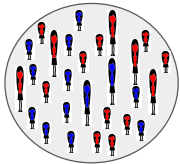
The diagram illustrates the formula for maximum weighted loss discrepancy,  $U(w)$ . The formula is  $U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell] \right|$ . Annotations in blue boxes with lines pointing to the formula components are as follows:

- Group weight**: Points to  $w(g)$ .
- Group loss**: Points to  $\mathbb{E}[\ell \mid g = 1]$ .
- population loss**: Points to  $\mathbb{E}[\ell]$ .
- All possible groups on population**: Points to the maximization set  $g \in \mathcal{G}$ .

$$U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell] \right|$$

$$U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell] \right|$$

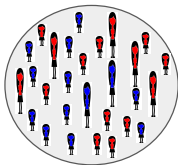
Normal



$$w(g) = \begin{cases} 1 & g = \text{all points} \\ 0 & \text{o. w.} \end{cases}$$

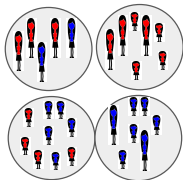
$$U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell] \right|$$

### Normal



$$w(g) = \begin{cases} 1 & g = \text{all points} \\ 0 & \text{o. w.} \end{cases}$$

### Group fairness



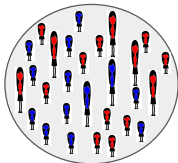
Hardt et al. 2016



$$w(g) = \begin{cases} 1 & \text{sensitive} \\ 0 & \text{o. w.} \end{cases}$$

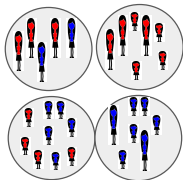
$$U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell] \right|$$

### Normal



$$w(g) = \begin{cases} 1 & g = \text{all points} \\ 0 & \text{o. w.} \end{cases}$$

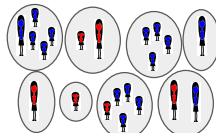
### Group fairness



Hardt et al. 2016

$$w(g) = \begin{cases} 1 & \text{sensitive} \\ 0 & \text{o. w.} \end{cases}$$

### Subgroup fairness

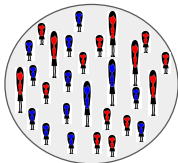


kearns et al. 2018

$$w(g) = \begin{cases} \mathbb{E}[g] & g \in G_{\text{sub}} \\ 0 & \text{o. w.} \end{cases}$$

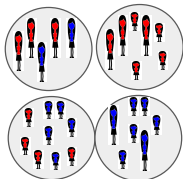
$$U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell] \right|$$

### Normal



$$w(g) = \begin{cases} 1 & g = \text{all points} \\ 0 & \text{o. w.} \end{cases}$$

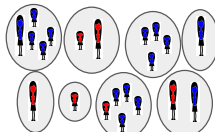
### Group fairness



Hardt et al. 2016

$$w(g) = \begin{cases} 1 & \text{sensitive} \\ 0 & \text{o. w.} \end{cases}$$

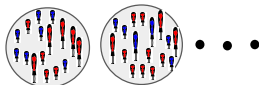
### Subgroup fairness



kearns et al. 2018

$$w(g) = \begin{cases} \mathbb{E}[g] & g \in G_{sub} \\ 0 & \text{o. w.} \end{cases}$$

### Large-group fairness



Hashimoto et al. 2018

$$w(g) = \begin{cases} 1 & \mathbb{E}[g] \geq \alpha \\ 0 & \text{o. w.} \end{cases}$$

## Interpretation

$$U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g] - \mathbb{E}[\ell] \right|$$

## Interpretation

$$U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g] - \mathbb{E}[\ell] \right|$$

For any  $g$ , we have:

$$\mathbb{E}[\ell] - \frac{U(w)}{w(g)} \leq \mathbb{E}[\ell \mid g] \leq \mathbb{E}[\ell] + \frac{U(w)}{w(g)}$$

## Interpretation

$$U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g] - \mathbb{E}[\ell] \right|$$

For any  $g$ , we have:

$$\mathbb{E}[\ell] - \frac{U(w)}{w(g)} \leq \mathbb{E}[\ell \mid g] \leq \mathbb{E}[\ell] + \frac{U(w)}{w(g)}$$

**Note.** If  $w(g) = 0$  then  $\mathbb{E}[\ell \mid g]$  can be arbitrary even when  $U(w)$  is small.



Without any prior information, what is an appropriate weighting function?

Without any prior information, what is an appropriate weighting function?

$$U(w) = \max_{g \in \mathcal{G}} w(g) |\mathbb{E}[\ell \mid g] - \mathbb{E}[\ell]|$$

$$w(g) = ?$$



$$w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$$

### Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

### Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

### Proposition

*There is no auditor for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$*

### Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

### Proposition

*There is no auditor for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$*

### Proof idea.

- Let  $\mathcal{P}_1$  be any distribution such that  $U(w) < \frac{1}{2}$  for  $\mathcal{P}_1$ .

## Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

## Proposition

*There is no auditor for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$*

## Proof idea.

- ▶ Let  $\mathcal{P}_1$  be any distribution such that  $U(w) < \frac{1}{2}$  for  $\mathcal{P}_1$ .
- ▶ Construct  $\mathcal{P}_2$  such that  $z \sim \mathcal{P}_1$  with probability  $1 - \eta$  and  $z = z_0$  and  $z = z_1$  each with probability  $\frac{\eta}{2} \implies U(w) \geq \frac{1}{2}$  for  $\mathcal{P}_2$

## Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

## Proposition

*There is no auditor for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$*

## Proof idea.

- ▶ Let  $\mathcal{P}_1$  be any distribution such that  $U(w) < \frac{1}{2}$  for  $\mathcal{P}_1$ .
- ▶ Construct  $\mathcal{P}_2$  such that  $z \sim \mathcal{P}_1$  with probability  $1 - \eta$  and  $z = z_0$  and  $z = z_1$  each with probability  $\frac{\eta}{2} \implies U(w) \geq \frac{1}{2}$  for  $\mathcal{P}_2$
- ▶ Choose a small  $\eta = 1 - \sqrt[n]{2\delta}$ .



### Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

### Proposition

*There is no auditor for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$*

### Take aways.

If we want to have positive weights for all groups then smaller groups should have smaller weights □

Without any prior information, what is an appropriate weighting function?

Without any prior information, what is an appropriate weighting function?

$$U(w) = \max_{g \in \mathcal{G}} w(g) |\mathbb{E}[\ell \mid g] - \mathbb{E}[\ell]|$$

$$w(g) = ?$$



$$w(g) = \mathbb{E}[g]^k \quad k \in (0, 1]$$

### Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

### Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

### Theorem

*Fix a bounded measurable loss. There is an auditor when  $w(g) = \mathbb{E}[g]^k$  for  $k \in (0, 1)$  which needs  $n = \mathcal{O}\left(\frac{\ln(1/\delta)}{\epsilon^{2+\frac{1}{k}}}\right)$  data points.*

## Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

## Theorem

*Fix a bounded measurable loss. There is an auditor when  $w(g) = \mathbb{E}[g]^k$  for  $k \in (0, 1)$  which needs  $n = \mathcal{O}\left(\frac{\ln(1/\delta)}{\epsilon^{2+\frac{1}{k}}}\right)$  data points.*

## Example

$$k = \frac{1}{2}$$

$$n = \mathcal{O}\left(\frac{\ln(1/\delta)}{\epsilon^4}\right)$$

## Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

## Theorem

*Fix a bounded measurable loss. There is an auditor when  $w(g) = \mathbb{E}[g]^k$  for  $k \in (0, 1)$  which needs  $n = \mathcal{O}\left(\frac{\ln(1/\delta)}{\epsilon^{2+\frac{1}{k}}}\right)$  data points.*

## Take aways.

For small  $k$ , we need more examples for convergence!





## Definition (Auditor)

Given any  $\epsilon, \delta \in (0, \frac{1}{2})$  an auditor returns an estimate  $\gamma$  for  $U(w)$  such that:

$$\mathbb{P}[|U(w) - \gamma| \leq \epsilon] \geq 1 - \delta$$

## Theorem

*Fix a bounded measurable loss. There is an auditor when  $w(g) = \mathbb{E}[g]^k$  for  $k \in (0, 1)$  which needs  $n = \mathcal{O}\left(\frac{\ln(1/\delta)}{\epsilon^{2+\frac{1}{k}}}\right)$  data points.*

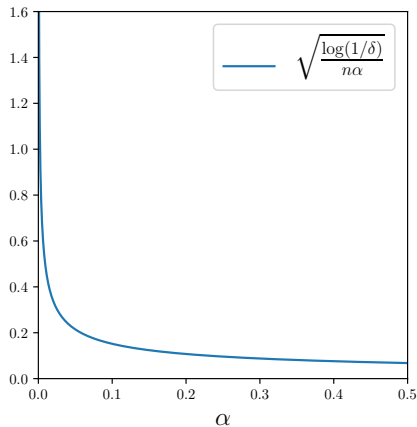
## Proof idea.

- ▶  $\mathbb{P}\left[\left|U(w) - \hat{U}(w)\right| \leq \epsilon\right] \geq 1 - \delta$
- ▶ We can compute  $\hat{U}(w)$  efficiently.



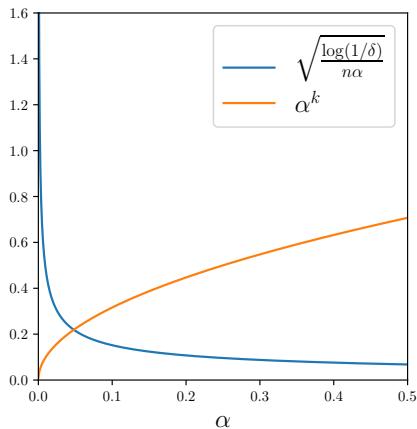
$D(g^*)$  converges to  $\hat{D}(g^*)$

- ▶  $D(g) = \mathbb{E}[g]^k |\mathbb{E}[\ell \mid g] - \mathbb{E}[\ell]|$   
(Weighted loss discrepancy)
- ▶  $g^* = \arg \max_{g \in \mathcal{G}} D(g)$
- ▶  $\alpha = \mathbb{E}[g^*]$
- ▶  $D(g^*) - \hat{D}(g^*) \leq \sqrt{\frac{\ln(1/\delta)}{n\alpha}}$



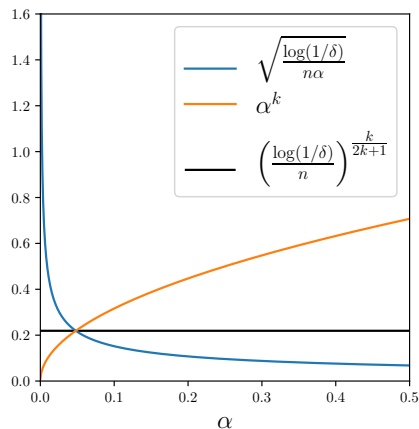
$D(g^*)$  converges to  $\hat{D}(g^*)$

- ▶  $D(g) = \mathbb{E}[g]^k |\mathbb{E}[\ell \mid g] - \mathbb{E}[\ell]|$   
(Weighted loss discrepancy)
- ▶  $g^* = \arg \max_{g \in \mathcal{G}} D(g)$
- ▶  $\alpha = \mathbb{E}[g^*]$
- ▶  $D(g^*) - \hat{D}(g^*) \leq \sqrt{\frac{\ln(1/\delta)}{n\alpha}}$
- ▶  $D(g^*) - \hat{D}(g^*) \leq \alpha^k$



$D(g^*)$  converges to  $\hat{D}(g^*)$

- ▶  $D(g) = \mathbb{E}[g]^k |\mathbb{E}[\ell \mid g] - \mathbb{E}[\ell]|$   
(Weighted loss discrepancy)
- ▶  $g^* = \arg \max_{g \in \mathcal{G}} D(g)$
- ▶  $\alpha = \mathbb{E}[g^*]$
- ▶  $D(g^*) - \hat{D}(g^*) \leq \sqrt{\frac{\ln(1/\delta)}{n\alpha}}$
- ▶  $D(g^*) - \hat{D}(g^*) \leq \alpha^k$
- ▶  $D(g^*) - \hat{D}(g^*) \leq \max_{\alpha} \min \left( \alpha^k, \sqrt{\frac{\ln(1/\delta)}{n\alpha}} \right)$



$\hat{D}(\hat{g})$  uniformly converges to  $D(\hat{g})$

Lemma

Let  $\hat{D}(g) = \hat{\mathbb{E}}[g]^k \left| \hat{\mathbb{E}}[\ell \mid g] - \hat{\mathbb{E}}[\ell] \right|$ . Given  $n$  data points, let  $\hat{g}$  be a group with maximum empirical weighted loss discrepancy,  $\hat{g} = \arg \max_{g \in \hat{\mathcal{G}}} \hat{D}(g)$ .

$\hat{D}(\hat{g})$  uniformly converges to  $D(\hat{g})$

**Lemma**

Let  $\hat{D}(g) = \hat{\mathbb{E}}[g]^k \left| \hat{\mathbb{E}}[\ell \mid g] - \hat{\mathbb{E}}[\ell] \right|$ . Given  $n$  data points, let  $\hat{g}$  be a group with maximum empirical weighted loss discrepancy,  $\hat{g} = \arg \max_{g \in \hat{\mathcal{G}}} \hat{D}(g)$ . Assume  $\hat{\mathbb{E}}[\ell] = 0$  and  $\hat{\mathbb{E}}[\ell \mid \hat{g}] \leq \hat{\mathbb{E}}[\ell]$ , then we can represent  $\hat{g}$ :

$$\underbrace{\ell_1 \leq \ell_2 \leq \cdots \leq \ell_t}_{\hat{g}} < \ell_{t+1} \leq \cdots \leq \ell_n$$

$\hat{D}(\hat{g})$  uniformly converges to  $D(\hat{g})$

Lemma

Let  $\hat{D}(g) = \hat{\mathbb{E}}[g]^k \left| \hat{\mathbb{E}}[\ell \mid g] - \hat{\mathbb{E}}[\ell] \right|$ . Given  $n$  data points, let  $\hat{g}$  be a group with maximum empirical weighted loss discrepancy,  $\hat{g} = \arg \max_{g \in \hat{\mathcal{G}}} \hat{D}(g)$ . Assume  $\hat{\mathbb{E}}[\ell] = 0$  and  $\hat{\mathbb{E}}[\ell \mid \hat{g}] \leq \hat{\mathbb{E}}[\ell]$ , then we can represent  $\hat{g}$ :

$$\underbrace{\ell_1 \leq \ell_2 \leq \cdots \leq \ell_t}_{\hat{g}} < \ell_{t+1} \leq \cdots \leq \ell_n$$

$$g_u = \{\text{all points with loss less than } u\}$$

$$\sup_{g \in \mathcal{G}} |D(\hat{g}) - \hat{D}(\hat{g})| = \sup_{u \in [0,1]} |D(\hat{g}_u) - \hat{D}(\hat{g}_u)|$$



Maximum weighted loss discrepancy,  $U(w)$





1

Maximum weighted loss discrepancy,  $U(w)$

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$



1

Maximum weighted loss discrepancy,  $U(w)$

2

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$



Possible to audit  $U(w)$  for  $w(g) = \mathbb{E}[g]^k$

1

Maximum weighted loss discrepancy,  $U(w)$

2

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$

3

Possible to audit  $U(w)$  for  $w(g) = \mathbb{E}[g]^k$

DETOUR



Loss variance

## Hoeffding's inequality

$$\left| \underbrace{\hat{\mathbb{E}}[\ell]}_{\text{training loss}} - \underbrace{\mathbb{E}[\ell]}_{\text{population loss}} \right| \leq \sqrt{\frac{C_1 \ln(1/\delta)}{n}}$$

## Hoeffding's inequality

$$\left| \underbrace{\hat{\mathbb{E}}[\ell]}_{\text{training loss}} - \underbrace{\mathbb{E}[\ell]}_{\text{population loss}} \right| \leq \sqrt{\frac{C_1 \ln(1/\delta)}{n}}$$

## Bennett's inequality

$$\left| \underbrace{\hat{\mathbb{E}}[\ell]}_{\text{training loss}} - \underbrace{\mathbb{E}[\ell]}_{\text{population loss}} \right| \leq \sqrt{\frac{C_1 \ln(1/\delta) \text{Var}[\ell]}{n}} + \frac{C_2 \ln(1/\delta)}{n}$$

## Bennett's inequality

$$\left| \underbrace{\hat{\mathbb{E}}[\ell]}_{\text{training loss}} - \underbrace{\mathbb{E}[\ell]}_{\text{population loss}} \right| \leq \sqrt{\frac{C_1 \text{Var}[\ell]}{n}} + \frac{C_2}{n}$$

[Bennett, 1962]

[Maurer and Pontil, 2009]

[Mnih et al., 2008]

[Audibert et al., 2009]

[Shivaswamy and Jebara, 2010]

[Namkoong and Duchi, 2017]

## Bennett's inequality

$$\left| \underbrace{\hat{\mathbb{E}}[\ell]}_{\text{training loss}} - \underbrace{\mathbb{E}[\ell]}_{\text{population loss}} \right| \leq \sqrt{\frac{C_1 \text{Var}[\ell]}{n}} + \frac{C_2}{n}$$

[Bennett, 1962]

[Maurer and Pontil, 2009]  $\implies$  Empirical Bernstein Bounds and Sample Variance Penalization

[Mnih et al., 2008]

[Audibert et al., 2009]

[Shivaswamy and Jebara, 2010]

[Namkoong and Duchi, 2017]

## Bennett's inequality

$$\left| \underbrace{\hat{\mathbb{E}}[\ell]}_{\text{training loss}} - \underbrace{\mathbb{E}[\ell]}_{\text{population loss}} \right| \leq \sqrt{\frac{C_1 \text{Var}[\ell]}{n}} + \frac{C_2}{n}$$

[Bennett, 1962]

[Maurer and Pontil, 2009]  $\implies$  Empirical Bernstein Bounds and Sample Variance Penalization

[Mnih et al., 2008]

[Audibert et al., 2009]

[Shivaswamy and Jebara, 2010]

[Namkoong and Duchi, 2017]  $\implies$  Variance-based regularization with convex objectives



1

Maximum weighted loss discrepancy,  $U(w)$

2

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$

3

Possible to audit  $U(w)$  for  $w(g) = \mathbb{E}[g]^k$

DETOUR



Loss variance

Test error

## Average individual discrepancy

$$\begin{aligned}\text{Var}[\ell] &= \mathbb{E} \left[ (\ell(z) - \mathbb{E}[\ell])^2 \right] \\ &= \frac{1}{2} \mathbb{E}_{z, z' \sim p^*} \left[ (\ell(z) - \ell(z'))^2 \right]\end{aligned}$$

1

Maximum weighted loss discrepancy,  $U(w)$

2

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$

3

Possible to audit  $U(w)$  for  $w(g) = \mathbb{E}[g]^k$

DETOUR



Loss variance

Test error

Average individual discrepancy

1

Maximum weighted loss discrepancy,  $U(w)$

2

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$

3

Possible to audit  $U(w)$  for  $w(g) = \mathbb{E}[g]^k$



DETOUR

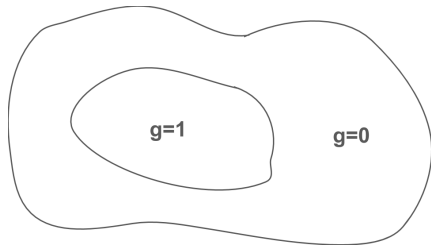
Loss variance

Test error

Average individual discrepancy

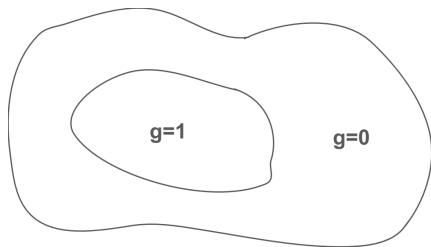
$U(w)$  for  $w(g) = \mathbb{E}[g]^{\frac{1}{2}}$   $\longleftrightarrow$  Loss variance

## law of total variance



$$\text{Var}[\mathbb{E}[\ell \mid g]] \leq \text{Var}[\ell]$$

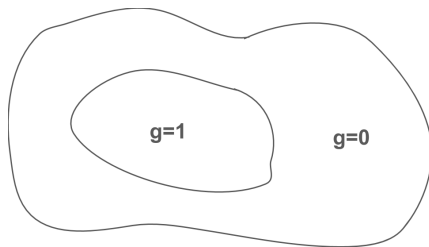
## law of total variance



$$\text{Var}[\mathbb{E}[\ell \mid g]] \leq \text{Var}[\ell]$$

$$\mathbb{E}[g](\mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell])^2 + \mathbb{E}[1 - g](\mathbb{E}[\ell \mid g = 0] - \mathbb{E}[\ell])^2 \leq \text{Var}[\ell]$$

## law of total variance

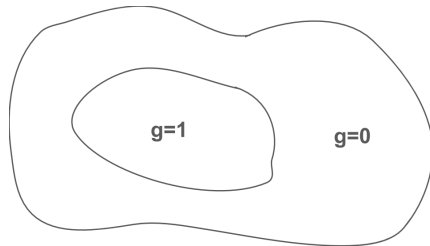


$$\text{Var}[\mathbb{E}[\ell \mid g]] \leq \text{Var}[\ell]$$

$$\mathbb{E}[g](\mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell])^2 + \mathbb{E}[1 - g](\mathbb{E}[\ell \mid g = 0] - \mathbb{E}[\ell])^2 \leq \text{Var}[\ell]$$

$$\underbrace{\sqrt{\mathbb{E}[g]}}_{\text{weighted}} \underbrace{|\mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell]|}_{\text{loss discrepancy}} \leq \sqrt{\text{Var}[\ell]}$$

## law of total variance



$$\text{Var}[\mathbb{E}[\ell \mid g]] \leq \text{Var}[\ell]$$

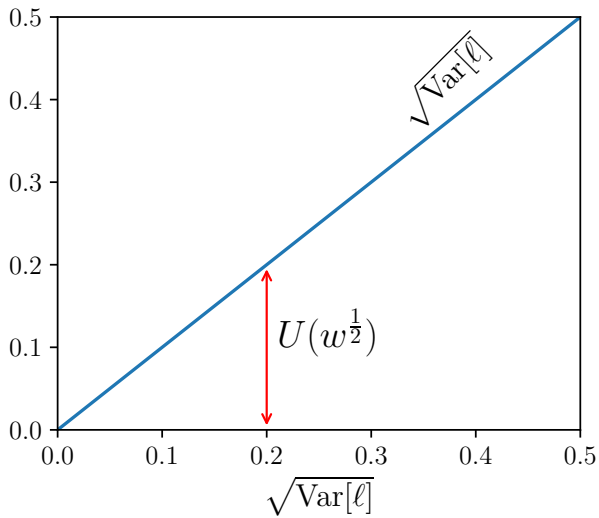
$$\mathbb{E}[g](\mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell])^2 + \mathbb{E}[1 - g](\mathbb{E}[\ell \mid g = 0] - \mathbb{E}[\ell])^2 \leq \text{Var}[\ell]$$

$$\underbrace{\sqrt{\mathbb{E}[g]}}_{\text{weighted}} \underbrace{|\mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell]|}_{\text{loss discrepancy}} \leq \sqrt{\text{Var}[\ell]}$$

### Proposition

Let  $w^{\frac{1}{2}}(g) = \mathbb{E}[g]^{\frac{1}{2}}$  then  $U(w^{\frac{1}{2}}) \leq \sqrt{\text{Var}[\ell]}$ .





## Proposition

Let  $w^{\frac{1}{2}}(g) = \mathbb{E}[g]^{\frac{1}{2}}$  then  $\sqrt{\text{Var}[\ell]} \leq U(w^{\frac{1}{2}}) \sqrt{2 - 4 \ln(U(w^{\frac{1}{2}}))}$

## Proposition

Let  $w^{\frac{1}{2}}(g) = \mathbb{E}[g]^{\frac{1}{2}}$  then  $\sqrt{\text{Var}[\ell]} \leq U(w^{\frac{1}{2}}) \sqrt{2 - 4 \ln(U(w^{\frac{1}{2}}))}$

Proof idea.

$$\blacktriangleright \mathbb{P}[\ell \geq u] \leq \frac{U(w^{\frac{1}{2}})}{u^2}$$

## Proposition

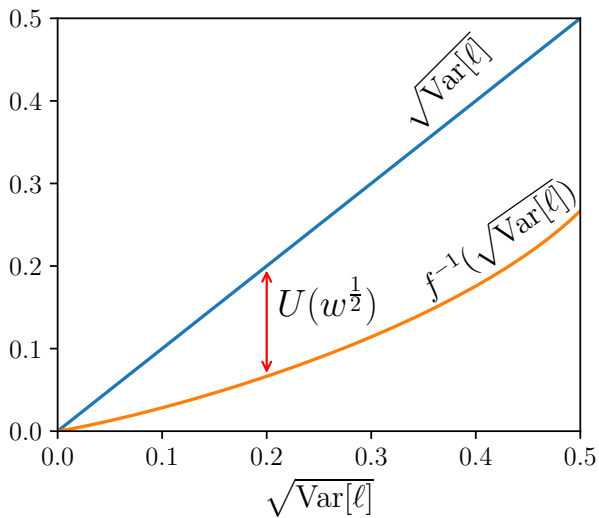
Let  $w^{\frac{1}{2}}(g) = \mathbb{E}[g]^{\frac{1}{2}}$  then  $\sqrt{\text{Var}[\ell]} \leq U(w^{\frac{1}{2}}) \sqrt{2 - 4 \ln(U(w^{\frac{1}{2}}))}$

Proof idea.

$$\blacktriangleright \mathbb{P}[\ell \geq u] \leq \frac{U(w^{\frac{1}{2}})}{u^2}$$

$$\blacktriangleright \text{Var}[\ell] = \int u \mathbb{P}[\ell \geq u]$$





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

1

Maximum weighted loss discrepancy,  $U(w)$

2

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$

3

Possible to audit  $U(w)$  for  $w(g) = \mathbb{E}[g]^k$


 DETOUR

Loss variance

Test error

Average individual discrepancy

$U(w)$  for  $w(g) = \mathbb{E}[g]^{\frac{1}{2}}$   $\longleftrightarrow$  Loss variance

1

Maximum weighted loss discrepancy,  $U(w)$

2

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$

3

Possible to audit  $U(w)$  for  $w(g) = \mathbb{E}[g]^k$



DETOUR

Loss variance

Test error

Average individual discrepancy

4

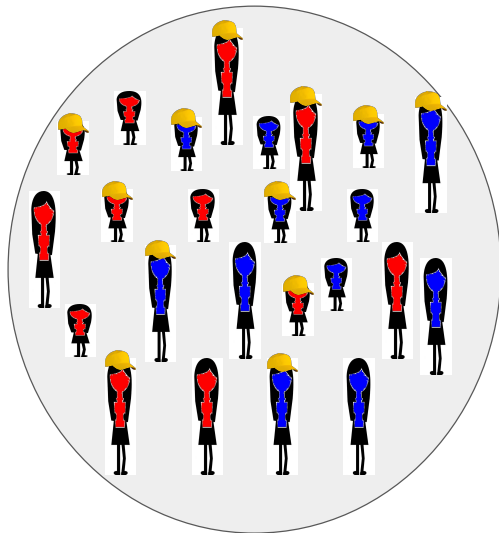
$U(w)$  for  $w(g) = \mathbb{E}[g]^{\frac{1}{2}}$   $\longleftrightarrow$  Loss variance

No prior  
information

# Handling prior information

A is given

A = [height, color]



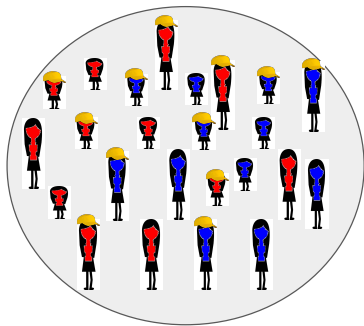


$$U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell] \right|$$

$$U(w) = \max_{g \in \mathcal{G}} w(g) \left| \mathbb{E}[\ell \mid g = 1] - \mathbb{E}[\ell] \right|$$

$$w(g) = \begin{cases} > 0 & g \in \mathcal{G}_A \\ 0 & o. w. \end{cases}$$

All possible groups on  $A = [\text{height, color}]$



# Coarse loss variance



$$\ell = 1$$

$$\ell = 0$$

$$\ell = 1$$

$$\ell = 0$$



$$\ell = 1$$

$$\ell = 0$$

$$\ell = 1$$

$$\ell = 0$$

$$\text{Var}[\ell] = 0.25$$

# Coarse loss variance



$$\ell = 1$$

$$\ell = 0$$

$$\ell = 1$$

$$\ell = 0$$



$$\ell = 1$$

$$\ell = 0$$

$$\ell = 1$$

$$\ell = 0$$

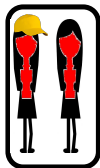
$$\text{Var}[\ell] = 0.25$$



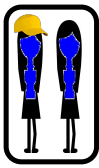
$$\mathbb{E}[\ell \mid A = \begin{bmatrix} \text{short}, \\ \text{red} \end{bmatrix}] = 0.5$$



$$\mathbb{E}[\ell \mid A = \begin{bmatrix} \text{short}, \\ \text{blue} \end{bmatrix}] = 0.5$$



$$\mathbb{E}[\ell \mid A = \begin{bmatrix} \text{tall}, \\ \text{red} \end{bmatrix}] = 0.5$$



$$\mathbb{E}[\ell \mid A = \begin{bmatrix} \text{tall}, \\ \text{blue} \end{bmatrix}] = 0.5$$

$$\text{Var}[\mathbb{E}[\ell \mid A]] = 0$$

## Coarse loss variance

$$U(w) \text{ for } w(g) = \begin{cases} \mathbb{E}[g]^{\frac{1}{2}} & g \in \mathcal{G}_A \\ 0 & \text{o. w.} \end{cases} \iff \text{Var}[\mathbb{E}[\ell \mid A]]$$

All possible groups on A = [height, color]

1

Maximum weighted loss discrepancy,  $U(w)$

2

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$

3

Possible to audit  $U(w)$  for  $w(g) = \mathbb{E}[g]^k$



DETOUR

Loss variance

Test error

Average individual discrepancy

4

$U(w)$  for  $w(g) = \mathbb{E}[g]^{\frac{1}{2}}$   $\longleftrightarrow$  Loss variance

$U(w)$  for  $w(g) = \begin{cases} \mathbb{E}[g]^{\frac{1}{2}} & g \in \mathcal{G}_A \\ 0 & o.w. \end{cases}$   $\longleftrightarrow$  Coarse loss variance

1

Maximum weighted loss discrepancy,  $U(w)$

2

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$

3

Possible to audit  $U(w)$  for  $w(g) = \mathbb{E}[g]^k$



DETOUR

Loss variance

Test error

Average individual discrepancy

4

$U(w)$  for  $w(g) = \mathbb{E}[g]^{\frac{1}{2}}$   $\longleftrightarrow$  Loss variance

5

$U(w)$  for  $w(g) = \begin{cases} \mathbb{E}[g]^{\frac{1}{2}} & g \in \mathcal{G}_A \\ 0 & o.w. \end{cases}$   $\longleftrightarrow$  Coarse loss variance



What should be the weighting function?

# Questions?

1

Maximum weighted loss discrepancy,  $U(w)$

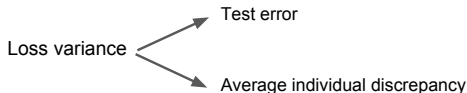
2

Impossible to audit  $U(w)$  for  $w(g) = \mathbb{I}[\mathbb{E}[g] > 0]$

3

Possible to audit  $U(w)$  for  $w(g) = \mathbb{E}[g]^k$

DETOUR



4

$U(w)$  for  $w(g) = \mathbb{E}[g]^{\frac{1}{2}}$   $\longleftrightarrow$  Loss variance

5

$U(w)$  for  $w(g) = \begin{cases} \mathbb{E}[g]^{\frac{1}{2}} & g \in \mathcal{G}_A \\ 0 & o.w. \end{cases}$   $\longleftrightarrow$  Coarse loss variance

6

What should be the weighting function?



## References

- J. Audibert, R. Munos, and C. Szepesv'ari. Exploration–exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19): 1876–1902, 2009.
- G. Bennett. Probability inequalities for the sum of independent random variables. *Journal of the American Statistical Association (JASA)*, 57(297):33–45, 1962.
- A. Maurer and M. Pontil. Empirical bernstein bounds and sample variance penalization. *arXiv preprint arXiv:0907.3740*, 2009.
- V. Mnih, C. Szepesv'ari, and J. Audibert. Empirical berstein stopping. In *International Conference on Machine Learning (ICML)*, 2008.
- H. Namkoong and J. Duchi. Variance regularization with convex objectives. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- P. Shivaswamy and T. Jebara. Empirical Bernstein boosting. In *Artificial Intelligence and Statistics (AISTATS)*, pages 733–740, 2010.