# An Axiomatic Theory of Provably-Fair Welfare-Centric Machine Learning
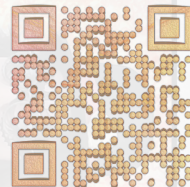
## Cyrus Cousins

Brown University
Department of Computer Science
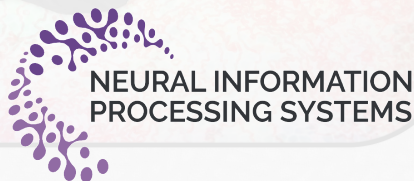
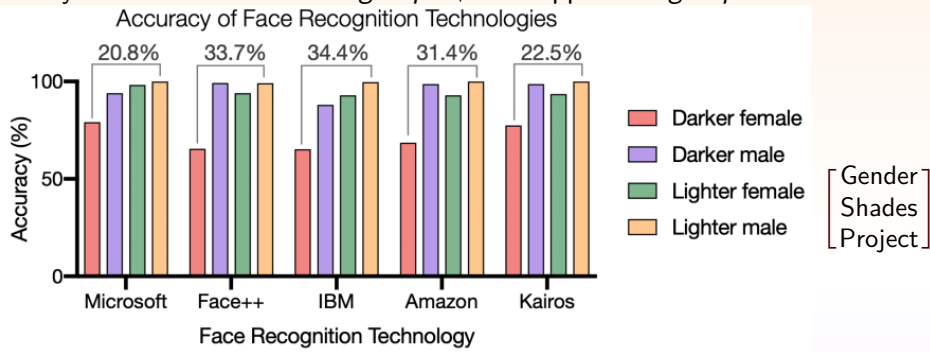December 2021

http://cs.brown.edu/people/ccousins/

BROWN
Computer Science

NEURAL INFORMATION
PROCESSING SYSTEMS

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
Welfare
Malfare
Axiomatic Characterization

Estimation and Inference
Linear Classifiers
Statistical Estimation

Fair PAC Learning
Computational Learnability

In Conclusion

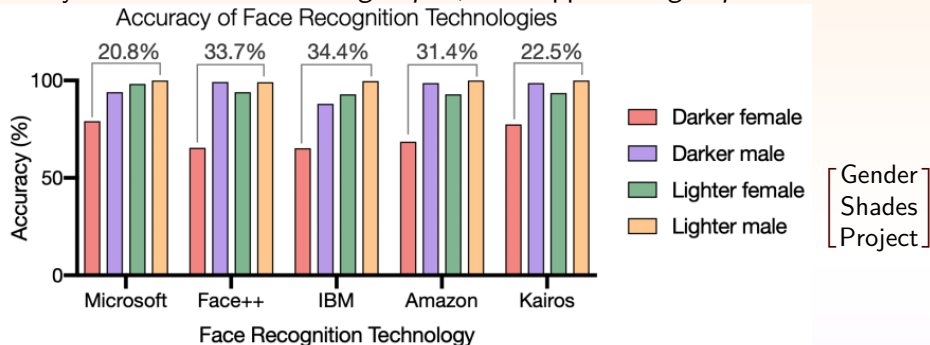# Fairness in Machine Learning (or Lack Thereof)

- ML systems often trained on *group A*, then applied to *group B*



Accuracy of Face Recognition Technologies

[ Gender Shades Project ]

- ML systems often trained on *group A*, then applied to *group B*



Accuracy of Face Recognition Technologies

- Differential performance $\implies$ algorithmic discrimination
  - Facial recognition and policing
  - Speech recognition and accessibility
  - Many more examples

- ML systems often trained on *group A*, then applied to *group B*



Accuracy of Face Recognition Technologies

$$\begin{bmatrix} \text{Gender} \\ \text{Shades} \\ \text{Project} \end{bmatrix}$$

- Differential performance $\implies$ algorithmic discrimination
  - Facial recognition and policing
  - Speech recognition and accessibility
  - Many more examples

- What has gone wrong? Is the problem:
  1. *that* a machine is learning;
  2. *from what* a machine is learning; or
  3. *how* a machine is learning?

- Focus on *differential accuracy* between *protected groups*
  - Want *group level* fairness
  - Learn from sample of *many individuals* drawn from *each group*

- Focus on *differential accuracy* between *protected groups*
  - Want *group level* fairness
  - Learn from sample of *many individuals* drawn from *each group*
- Claim: current ML systems are trained:
  - On the wrong data (well-known)
  - In the wrong way

- Focus on *differential accuracy* between *protected groups*
  - Want *group level* fairness
  - Learn from sample of *many individuals* drawn from *each group*
- Claim: current ML systems are trained:
  - On the wrong data (well-known)
  - In the wrong way
- Optimize models sensitive to performance on *protected-groups*
  - Introduce malfare learning target
  - Consider all groups (possibly nonlinearly)

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
Welfare
Malfare
Axiomatic Characterization

Estimation and Inference
Linear Classifiers
Statistical Estimation

Fair PAC Learning
Computational Learnability

In Conclusion

- Focus on *differential accuracy* between *protected groups*
  - Want *group level* fairness
  - Learn from sample of *many individuals* drawn from *each group*
- Claim: current ML systems are trained:
  - On the wrong data (well-known)
  - In the wrong way
- Optimize models sensitive to performance on *protected-groups*
  - Introduce malfare learning target
  - Consider all groups (possibly nonlinearly)
- Theoretical treatment of learning and statistics
  - Overfitting and statistical estimation
  - Computational complexity issues in learning
  - Introduce *fair-PAC-learning* to theoretically treat these issues

# First Canto: The Philosophy of Welfare and Malfare

## Fair machine learning and the social planning problem

- *Utility*: $\mathrm{U}(\cdot) : \mathcal{X} \to \mathbb{R}_{0+}$ Subjective measurement of *positive attribute*
  - Happiness, satisfaction, resource ownership

- *Utility*: $\mathrm{U}(\cdot) : \mathcal{X} \to \mathbb{R}_{0+}$ Subjective measurement of *positive attribute*
  - Happiness, satisfaction, resource ownership

- *Welfare* summarizes *population-level* utility across $\boldsymbol{x} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_g)$
  "*The subjective, measured objectively.*"

- *Utility*: $\mathrm{U}(\cdot) : \mathcal{X} \to \mathbb{R}_{0+}$ Subjective measurement of *positive attribute*
  - Happiness, satisfaction, resource ownership

- *Welfare* summarizes *population-level* utility across $\boldsymbol{x} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_g)$
  *"The subjective, measured objectively."*

- *Utilitarian* welfare: *average utility* $\mathrm{U}(\cdot)$

$$\mathrm{W}_{\mathrm{Util}}\big(\mathrm{U}(\boldsymbol{x}_1), \ldots, \mathrm{U}(\boldsymbol{x}_g)\big) \doteq \frac{1}{g}\sum_{i=1}^{g} \mathrm{U}(\boldsymbol{x}_i)$$

- *Utility*: $\mathrm{U}(\cdot) : \mathcal{X} \to \mathbb{R}_{0+}$ Subjective measurement of *positive attribute*
  - Happiness, satisfaction, resource ownership
- *Welfare* summarizes *population-level* utility across $\boldsymbol{x} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_g)$ "The subjective, measured objectively."
- *Utilitarian* welfare: *average utility* $\mathrm{U}(\cdot)$

$$\mathrm{W}_{\mathrm{Util}}\big(\mathrm{U}(\boldsymbol{x}_1), \ldots, \mathrm{U}(\boldsymbol{x}_g)\big) \doteq \frac{1}{g} \sum_{i=1}^{g} \mathrm{U}(\boldsymbol{x}_i)$$

- *Egalitarian* welfare: *worst-case utility* $\mathrm{U}(\cdot)$

$$\mathrm{W}_{\mathrm{Egal}}\big(\mathrm{U}(\boldsymbol{x}_1), \ldots, \mathrm{U}(\boldsymbol{x}_g)\big) \doteq \min_{i \in 1, \ldots, g} \mathrm{U}(\boldsymbol{x}_i)$$

  - A fair society should have *equality*
  - Incentivize aiding the most needy first

|              | Pop. A | Pop. B | Pop. C |
|--------------|--------|--------|--------|
| Utilitarian  |        |        |        |
| Egalitarian  |        |        |        |

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# What is Welfare (cont.)



|  | Pop. A | Pop. B | Pop. C |
|---|---|---|---|
| Utilitarian | | | |
| Egalitarian | | | |

- Limitations
  - Nonnegativity
  - *Positive* directedness (utility is desirable)

|  | Pop. A | Pop. B | Pop. C |
|---|---|---|---|
| Utilitarian |  |  |  |
| Egalitarian |  |  |  |

- Limitations
  - Nonnegativity
  - *Positive* directedness (utility is desirable)
- Which welfare function to use?
  - Analogy: *worst-case* vs *average case* bounds
  - Analogy: tail bounds vs expectation

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare

Welfare

Malfare

Axiomatic
Characterization

Estimation
and Inference

Linear Classifiers

Statistical Estimation

Fair PAC
Learning

Computational
Learnability

In Conclusion

# The Power Mean

Suppose vector $\ell = (\ell_1, \ldots, \ell_g)$ representing *utility* or *loss* across a population

Suppose vector $\boldsymbol{\ell} = (\boldsymbol{\ell}_1, \ldots, \boldsymbol{\ell}_g)$ representing *utility* or *loss* across a population

$$\mathrm{M}_p(\boldsymbol{\ell}) \doteq \begin{cases} p \in \mathbb{R} \setminus \{0\} & \sqrt[p]{\dfrac{1}{n} \sum_{i=1}^{n} \boldsymbol{\ell}_i^p} \end{cases}$$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# The Power Mean

Suppose vector $\boldsymbol{\ell} = (\boldsymbol{\ell}_1, \ldots, \boldsymbol{\ell}_g)$ representing *utility* or *loss* across a population

$$
\mathrm{M}_p(\boldsymbol{\ell}) \doteq
\begin{cases}
p \in \mathbb{R} \setminus \{0\} & \sqrt[p]{\dfrac{1}{n} \sum_{i=1}^{n} \boldsymbol{\ell}_i^p} \\[2ex]
p = -\infty & \min_{i \in 1,\ldots,g} \boldsymbol{\ell}_i \\[2ex]
p = 0 & \sqrt[n]{\prod_{i=1}^{n} \boldsymbol{\ell}_i} \\[2ex]
p = \infty & \max_{i \in 1,\ldots,g} \boldsymbol{\ell}_i
\end{cases}
$$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# The Power Mean

Suppose vector $\boldsymbol{\ell} = (\boldsymbol{\ell}_1, \dots, \boldsymbol{\ell}_g)$ representing *utility* or *loss* across a population

$$
\mathrm{M}_p(\boldsymbol{\ell}) \doteq
\begin{cases}
p \in \mathbb{R} \setminus \{0\} & \sqrt[p]{\dfrac{1}{n} \sum_{i=1}^{n} \boldsymbol{\ell}_i^p} \\[2em]
p = -\infty & \min_{i \in 1, \dots, g} \boldsymbol{\ell}_i \\[1.5em]
p = 0 & \sqrt[n]{\prod_{i=1}^{n} \boldsymbol{\ell}_i} \\[2em]
p = \infty & \max_{i \in 1, \dots, g} \boldsymbol{\ell}_i
\end{cases}
$$



$$\mathrm{M}_p(1,2,3)$$
$$\mathrm{M}_{-\infty}(1,2,3)$$
$$\mathrm{M}_{\infty}(1,2,3)$$

$p$

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
Welfare
Malfare
Axiomatic Characterization

Estimation and Inference
Linear Classifiers
Statistical Estimation

Fair PAC Learning
Computational Learnability

In Conclusion

# The Power Mean

Suppose vector $\boldsymbol{\ell} = (\boldsymbol{\ell}_1, \ldots, \boldsymbol{\ell}_g)$ representing *utility* or *loss* across a population

$$
\mathrm{M}_p(\boldsymbol{\ell}) \doteq
\begin{cases}
p \in \mathbb{R} \setminus \{0\} & \sqrt[p]{\dfrac{1}{n} \sum_{i=1}^{n} \boldsymbol{\ell}_i^p} \\[2ex]
p = -\infty & \min_{i \in 1, \ldots, g} \boldsymbol{\ell}_i \\[2ex]
p = 0 & \sqrt[n]{\prod_{i=1}^{n} \boldsymbol{\ell}_i} \\[2ex]
p = \infty & \max_{i \in 1, \ldots, g} \boldsymbol{\ell}_i
\end{cases}
$$



- *Smooth interpolation* between *min*, *arithmetic mean*, and *max*
  - Other special cases: *geometric*, *harmonic*, and *quadratic* means
- Monotonic in $p$: *interpolate between* utilitarian and egalitarian

- Why do we care about cardinal welfare?
  - *Welfare* $\mathrm{W}(\cdot)$ encodes an *ideal notion* of *societal wellbeing* (fairness)
  - Utilitarian versus Egalitarian

- Why do we care about cardinal welfare?
  - *Welfare* $\mathrm{W}(\cdot)$ encodes an *ideal notion* of *societal wellbeing* (fairness)
  - Utilitarian versus Egalitarian
- The *social planning problem*
  - Select allocation of goods and services to *maximize welfare*
  - *Fair ML* is *learning* an *optimal allocation* from data?
  - Learn *policy* to maximize *welfare* of *per-group* utilities

- Why do we care about cardinal welfare?
  - *Welfare* $\mathrm{W}(\cdot)$ encodes an *ideal notion* of *societal wellbeing* (fairness)
  - Utilitarian versus Egalitarian
- The *social planning problem*
  - Select allocation of goods and services to *maximize welfare*
  - *Fair ML* is *learning* an *optimal allocation* from data?
  - Learn *policy* to maximize *welfare* of *per-group* utilities

Is it really that easy?

- Why do we care about cardinal welfare?
  - *Welfare* $W(\cdot)$ encodes an *ideal notion* of *societal wellbeing* (fairness)
  - Utilitarian versus Egalitarian
- The *social planning problem*
  - Select allocation of goods and services to *maximize welfare*
  - *Fair ML* is *learning* an *optimal allocation* from data?
  - Learn *policy* to maximize *welfare* of *per-group* utilities

Is it really that easy?

What if we want to *minimize* a *loss function*?

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare

Welfare

Malfare

Axiomatic
Characterization

Estimation
and Inference

Linear Classifiers

Statistical Estimation

Fair PAC
Learning

Computational
Learnability

In Conclusion

# Introducing Malfare

- Standard: *maximize* a *welfare measure* of *societal wellbeing* $\mathrm{W}(\cdot)$
- This work: *minimize* a *malfare measure* of *societal suffering* $\mathrm{M}(\cdot)$
  - Generically termed *aggregator functions* $\mathrm{M}(\cdot)$

- Standard: *maximize* a *welfare measure* of *societal wellbeing* $\mathrm{W}(\cdot)$
- This work: *minimize* a *malfare measure* of *societal suffering* $\mathrm{M}(\cdot)$
  - Generically termed *aggregator functions* $\mathrm{M}(\cdot)$
- Is this really a new idea?
  - Everybody knows $\forall i : \underset{h \in \mathcal{H}}{\mathrm{argmax}} -\boldsymbol{\ell}_i(h) = \underset{h \in \mathcal{H}}{\mathrm{argmin}} \, \boldsymbol{\ell}_i(h)$
  - But we don't have $\underset{h \in \mathcal{H}}{\mathrm{argmax}} \, \mathrm{W}_p(-\boldsymbol{\ell}(h)) = \underset{h \in \mathcal{H}}{\mathrm{argmin}} \, \mathrm{M}_p(\boldsymbol{\ell}(h))$
  - Welfare is a *multivariate optimality concept*
  - Intuition from *univariate optimization* breaks down

- Standard: *maximize* a *welfare measure* of *societal wellbeing* $\mathrm{W}(\cdot)$
- This work: *minimize* a *malfare measure* of *societal suffering* $\mathrm{M}(\cdot)$
  - Generically termed *aggregator functions* $\mathrm{M}(\cdot)$
- Is this really a new idea?
  - Everybody knows $\forall i : \underset{h \in \mathcal{H}}{\operatorname{argmax}} -\boldsymbol{\ell}_i(h) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \boldsymbol{\ell}_i(h)$
  - But we don't have $\underset{h \in \mathcal{H}}{\operatorname{argmax}} \mathrm{W}_p(-\boldsymbol{\ell}(h)) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathrm{M}_p(\boldsymbol{\ell}(h))$
  - Welfare is a *multivariate optimality concept*
  - Intuition from *univariate optimization* breaks down
- We shall see *equal axiomatic justification* for welfare and malfare

- Standard: *maximize* a *welfare measure* of *societal wellbeing* $\mathrm{W}(\cdot)$
- This work: *minimize* a *malfare measure* of *societal suffering* $\mathrm{M}(\cdot)$
  - Generically termed *aggregator functions* $\mathrm{M}(\cdot)$
- Is this really a new idea?
  - Everybody knows $\forall i : \underset{h \in \mathcal{H}}{\operatorname{argmax}} -\boldsymbol{\ell}_i(h) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \boldsymbol{\ell}_i(h)$
  - But we don't have $\underset{h \in \mathcal{H}}{\operatorname{argmax}} \mathrm{W}_p(-\boldsymbol{\ell}(h)) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathrm{M}_p(\boldsymbol{\ell}(h))$
  - Welfare is a *multivariate optimality concept*
  - Intuition from *univariate optimization* breaks down
- We shall see *equal axiomatic justification* for welfare and malfare

| | 😀 | 🙂 | 😐 | 🙁 | 😡 |
|---|---|---|---|---|---|
| Utility | 5 | 4 | 3 | 2 | 1 |
| Loss | 1 | 2 | 3 | 4 | 5 |

- *Malfare* extends the concept of *welfare* to *undesirable quantities* (disutility)
- Direct correspondence only for $p \in \{-\infty, 1, \infty\}$

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
  Welfare
  **Malfare**
  Axiomatic Characterization

Estimation and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC Learning
  Computational Learnability

In Conclusion

Introducing Malfare

- Standard: *maximize* a *welfare measure* of *societal wellbeing* $\mathrm{W}(\cdot)$
- This work: *minimize* a *malfare measure* of *societal suffering* $\mathrm{M}(\cdot)$
  - Generically termed *aggregator functions* $\mathrm{M}(\cdot)$
- Is this really a new idea?
  - Everybody knows $\forall i: \underset{h \in \mathcal{H}}{\operatorname{argmax}} -\boldsymbol{\ell}_i(h) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \boldsymbol{\ell}_i(h)$
  - But we don't have $\underset{h \in \mathcal{H}}{\operatorname{argmax}} \mathrm{W}_p(-\boldsymbol{\ell}(h)) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathrm{M}_p(\boldsymbol{\ell}(h))$
  - Welfare is a *multivariate optimality concept*
  - Intuition from *univariate optimization* breaks down
- We shall see *equal axiomatic justification* for welfare and malfare

|        | 😆 | 🙂 | 😐 | 🙁 | 😡 |
|--------|----|----|----|----|----|
| Utility | 5  | 4  | 3  | 2  | 1  |
| Loss    | 1  | 2  | 3  | 4  | 5  |

- *Malfare* extends the concept of *welfare* to *undesirable quantities* (disutility)
- Direct correspondence only for $p \in \{-\infty, 1, \infty\}$



$\mathrm{M}_p(1;2;3)$

$4 - \mathrm{M}_{2-p}(4-(1,2,3))$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# Fair Machine Learning with Malfare Minimization

- Standard machine learning over *instance distribution* $\mathcal{D}$
  - Machine learning tasks often cast as *risk minimization* w.r.t. *loss function* $\ell$

$$h^* \doteq \operatorname*{argmin}_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(h(x), y)]$$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# Fair Machine Learning with Malfare Minimization

- Standard machine learning over *instance distribution* $\mathcal{D}$
  - Machine learning tasks often cast as *risk minimization* w.r.t. *loss function* $\ell$

$$h^* \doteq \underset{h \in \mathcal{H}}{\operatorname{argmin}} \ \underset{(x,y) \sim \mathcal{D}}{\mathbb{E}} [\ell(h(x), y)]$$

- In fair machine learning, *different groups* have *different needs and preferences*
  - Need to consider *multiple distributions* $\mathcal{D}_1, \ldots, \mathcal{D}_g$ over *multiple groups*
  - Analogously cast learning tasks as *malfare minimization*

$$h^* \doteq \underset{h \in \mathcal{H}}{\operatorname{argmin}} \ \mathrm{M} \left( \underset{(x,y) \sim \mathcal{D}_1}{\mathbb{E}} \big[\ell(h(x), y)\big], \ldots, \underset{(x,y) \sim \mathcal{D}_g}{\mathbb{E}} \big[\ell(h(x), y)\big] \right)$$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
  Welfare
  Malfare
  Axiomatic
  Characterization

Estimation
and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC
Learning
  Computational
  Learnability

In Conclusion

# Fair Machine Learning with Malfare Minimization

- Standard machine learning over *instance distribution* $\mathcal{D}$
  - Machine learning tasks often cast as *risk minimization* w.r.t. *loss function* $\ell$

$$h^* \doteq \underset{h \in \mathcal{H}}{\operatorname{argmin}} \underset{(x,y) \sim \mathcal{D}}{\mathbb{E}} [\ell(h(x), y)]$$

- In fair machine learning, *different groups* have *different needs and preferences*
  - Need to consider *multiple distributions* $\mathcal{D}_1, \ldots, \mathcal{D}_g$ over *multiple groups*
  - Analogously cast learning tasks as *malfare minimization*

$$h^* \doteq \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathbb{M} \left( \underset{(x,y) \sim \mathcal{D}_1}{\mathbb{E}} [\ell(h(x), y)], \ldots, \underset{(x,y) \sim \mathcal{D}_g}{\mathbb{E}} [\ell(h(x), y)] \right)$$

- Contrast with welfare maximization:
  ⚠ Need to define a *utility function* $\mathrm{U}(\cdot)$

$$h^* \doteq \underset{h \in \mathcal{H}}{\operatorname{argmax}} \mathrm{W} \left( \underset{(x,y) \sim \mathcal{D}_1}{\mathbb{E}} [\mathrm{U}(h(x), y)], \ldots, \underset{(x,y) \sim \mathcal{D}_g}{\mathbb{E}} [\mathrm{U}(h(x), y)] \right)$$

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
  Welfare
  **Malfare**
  Axiomatic Characterization

Estimation and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC Learning
  Computational Learnability

In Conclusion

# Fair Machine Learning with Malfare Minimization

- Standard machine learning over *instance distribution* $\mathcal{D}$
  - Machine learning tasks often cast as *risk minimization* w.r.t. *loss function* $\ell$

$$h^* \doteq \operatorname*{argmin}_{h \in \mathcal{H}} \mathop{\mathbb{E}}_{(x,y) \sim \mathcal{D}} [\ell(h(x), y)]$$

- In fair machine learning, *different groups* have *different needs and preferences*
  - Need to consider *multiple distributions* $\mathcal{D}_1, \ldots, \mathcal{D}_g$ over *multiple groups*
  - Analogously cast learning tasks as *malfare minimization*

$$h^* \doteq \operatorname*{argmin}_{h \in \mathcal{H}} \mathrm{M} \left( \mathop{\mathbb{E}}_{(x,y) \sim \mathcal{D}_1} \big[\ell(h(x), y)\big], \ldots, \mathop{\mathbb{E}}_{(x,y) \sim \mathcal{D}_g} \big[\ell(h(x), y)\big] \right)$$

- Contrast with welfare maximization:
  ⚠ Need to define a *utility function* $\mathrm{U}(\cdot)$

$$h^* \doteq \operatorname*{argmax}_{h \in \mathcal{H}} \mathrm{W} \left( \mathop{\mathbb{E}}_{(x,y) \sim \mathcal{D}_1} \big[\mathrm{U}(h(x), y)\big], \ldots, \mathop{\mathbb{E}}_{(x,y) \sim \mathcal{D}_g} \big[\mathrm{U}(h(x), y)\big] \right)$$

⚠ Don't know $\mathcal{D}_{1:g}$; have to work from *training samples*
  - Select $\hat{h}$ to optimize empirical *risk / malfare / welfare* estimates

Axiomatic Fair
Learning

Axioms of Cardinal Welfare

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare

Welfare

Malfare

Axiomatic
Characterization

Estimation
and Inference

Linear Classifiers

Statistical Estimation

Fair PAC
Learning

Computational
Learnability

In Conclusion

❶ Strict Monotonicity: $\forall \varepsilon \succeq 0$ s.t. $\varepsilon \neq 0$: $\mathrm{M}(\ell) < \mathrm{M}(\ell + \varepsilon)$

- Adding utility never harms welfare

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare

Welfare

Malfare

Axiomatic
Characterization

Estimation
and Inference

Linear Classifiers

Statistical Estimation

Fair PAC
Learning

Computational
Learnability

In Conclusion

# Axioms of Cardinal Welfare

❶ Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$
  - Adding utility never harms welfare
❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$
  - No *exceptionalism*; welfare is *identity blind*
  - Inherent tenet of *fairness* and *equality*

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare

Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference

Linear Classifiers
Statistical Estimation

Fair PAC
Learning

Computational
Learnability

In Conclusion

# Axioms of Cardinal Welfare

❶ Strict Monotonicity: $\forall \varepsilon \succeq 0$ s.t. $\varepsilon \neq 0$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \varepsilon)$
  • Adding utility never harms welfare
❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$
  • No *exceptionalism*; welfare is *identity blind*
  • Inherent tenet of *fairness* and *equality*
❸ Continuity: $\forall \boldsymbol{\ell}$ : $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# Axioms of Cardinal Welfare

❶ Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$

- Adding utility never harms welfare

❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$

- No *exceptionalism*; welfare is *identity blind*
- Inherent tenet of *fairness* and *equality*

❸ Continuity: $\forall \boldsymbol{\ell} : \{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

❹ Independence of Unconcerned Agents (IUA):
$\forall a, b \in \mathbb{R}_+ : \mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$

- Compartmentalization and analysis of *subgroups*

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
Welfare
Malfare
Axiomatic Characterization

Estimation and Inference
Linear Classifiers
Statistical Estimation

Fair PAC Learning
Computational Learnability

In Conclusion

# Axioms of Cardinal Welfare

❶ Strict Monotonicity: $\forall \varepsilon \succeq 0$ s.t. $\varepsilon \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \varepsilon)$
- Adding utility never harms welfare

❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$
- No *exceptionalism*; welfare is *identity blind*
- Inherent tenet of *fairness* and *equality*

❸ Continuity: $\forall \boldsymbol{\ell} : \{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

❹ Independence of Unconcerned Agents (IUA):
$\forall a, b \in \mathbb{R}_+ : \mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$
- Compartmentalization and analysis of *subgroups*

❺ Independence of Common Scale (ICS):
$\forall \alpha \in \mathbb{R}_+ : \mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha \boldsymbol{\ell}) \leq \mathrm{M}(\alpha \boldsymbol{\ell}')$
- Relative value *invariant* under (absolute) unit conversion: \$ versus ₽

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# Axioms of Cardinal Welfare

❶ Strict Monotonicity: $\forall \varepsilon \succeq 0$ s.t. $\varepsilon \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \varepsilon)$
   - Adding utility never harms welfare

❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$
   - No *exceptionalism*; welfare is *identity blind*
   - Inherent tenet of *fairness* and *equality*

❸ Continuity: $\forall \boldsymbol{\ell}$: $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

❹ Independence of Unconcerned Agents (IUA):
   $\forall a, b \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$
   - Compartmentalization and analysis of *subgroups*

❺ Independence of Common Scale (ICS):
   $\forall \alpha \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha\boldsymbol{\ell}) \leq \mathrm{M}(\alpha\boldsymbol{\ell}')$
   - Relative value *invariant* under (absolute) unit conversion: $ versus ₽

## Theorem (Debreu-Gorman (1959))

*For some strictly increasing $F$, $p \in \mathbb{R}$, all welfare functions satisfying 1-5 take form*
$$\mathrm{M}(\boldsymbol{\ell}) = F\left(\mathrm{sgn}(p) \sum_{i=1}^{g} \boldsymbol{\ell}_i^p\right) \ \text{ or } \ \mathrm{M}(\boldsymbol{\ell}) = F\left(\prod_{i=1}^{g} \boldsymbol{\ell}_i\right)$$

⚠️ Can be *non-Lipschitz*, arbitrarily hard to *estimate* from *sample*

❶ Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$

❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$

❸ Continuity: $\forall \boldsymbol{\ell}$: $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

❹ IUA: $\forall a, b \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$

❺ ICS: $\forall \alpha \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha \boldsymbol{\ell}) \leq \mathrm{M}(\alpha \boldsymbol{\ell}')$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# Extended Axioms of Cardinal Welfare

**❶** Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$

**❷** Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$

**❸** Continuity: $\forall \boldsymbol{\ell} : \{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

**❹** IUA: $\forall a, b \in \mathbb{R}_+ : \mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$

**❺** ICS: $\forall \alpha \in \mathbb{R}_+ : \mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha\boldsymbol{\ell}) \leq \mathrm{M}(\alpha\boldsymbol{\ell}')$

**❻** Multiplicative Linearity: $\mathrm{M}(\alpha\boldsymbol{\ell}) = \alpha\mathrm{M}(\boldsymbol{\ell})$;
  - Implies ICS
  - Units of $\mathrm{M}(\cdot)$ match units of $\boldsymbol{\ell}$

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
  Welfare
  Malfare
  **Axiomatic Characterization**

Estimation and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC Learning
  Computational Learnability

In Conclusion

# Extended Axioms of Cardinal Welfare

**❶** Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$

**❷** Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$

**❸** Continuity: $\forall \boldsymbol{\ell}$: $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

**❹** IUA: $\forall a, b \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$

**❺** ICS: $\forall \alpha \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha\boldsymbol{\ell}) \leq \mathrm{M}(\alpha\boldsymbol{\ell}')$

**❻** Multiplicative Linearity: $\mathrm{M}(\alpha\boldsymbol{\ell}) = \alpha\mathrm{M}(\boldsymbol{\ell})$;
  - Implies ICS
  - Units of $\mathrm{M}(\cdot)$ match units of $\boldsymbol{\ell}$

**❼** Unit Scale: $\mathrm{M}(\mathbf{1}) = 1$
  - Scale of $\mathrm{M}$ matches units of $\boldsymbol{\ell}$
  - Absolute comparison: "My income is $x\%$ of average / maximum / minimum"

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
  Welfare
  Malfare
  Axiomatic
  Characterization

Estimation
and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC
Learning
  Computational
  Learnability

In Conclusion

Extended Axioms of Cardinal Welfare

❶ Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$

❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$

❸ Continuity: $\forall \boldsymbol{\ell}$: $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

❹ IUA: $\forall a, b \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$

❺ ICS: $\forall \alpha \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha\boldsymbol{\ell}) \leq \mathrm{M}(\alpha\boldsymbol{\ell}')$

❻ Multiplicative Linearity: $\mathrm{M}(\alpha\boldsymbol{\ell}) = \alpha\mathrm{M}(\boldsymbol{\ell})$;
  - Implies ICS
  - Units of $\mathrm{M}(\cdot)$ match units of $\boldsymbol{\ell}$

❼ Unit Scale: $\mathrm{M}(\mathbf{1}) = 1$
  - Scale of $\mathrm{M}$ matches units of $\boldsymbol{\ell}$
  - Absolute comparison: "My income is $x\%$ of average / maximum / minimum"

## Theorem (Axiomatic Characterization of Welfare and Malfare)

*For any aggregator function $\mathrm{M}(\cdot)$ satisfying 1-7, $\exists p \in \mathbb{R}$ s.t.*

$$\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}_p(\boldsymbol{\ell}) = \sqrt[p]{\frac{1}{g} \sum_{i=1}^{g} \boldsymbol{\ell}_i^p} \ \text{ or } \ \mathrm{M}(\boldsymbol{\ell}) = \sqrt[g]{\prod_{i=1}^{g} \boldsymbol{\ell}_i} \ ,$$

*which are Lipschitz-continuous in $\boldsymbol{\ell}$ for $p \in (-\infty, 0) \cup [1, \infty)$.*

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
Welfare
Malfare
**Axiomatic Characterization**

Estimation and Inference
Linear Classifiers
Statistical Estimation

Fair PAC Learning
Computational Learnability

In Conclusion

# Transfer Axioms of Cardinal Welfare

**❶** Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$

**❷** Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$

**❸** Continuity: $\forall \boldsymbol{\ell}$: $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

**❹** IUA: $\forall a, b \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$

**❺** ICS: $\forall \alpha \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha\boldsymbol{\ell}) \leq \mathrm{M}(\alpha\boldsymbol{\ell}')$

**❻** Multiplicative Linearity: $\mathrm{M}(\alpha\boldsymbol{\ell}) = \alpha\mathrm{M}(\boldsymbol{\ell})$

**❼** Unit Scale: $\mathrm{M}(\mathbf{1}) = 1$

❶ Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$

❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$

❸ Continuity: $\forall \boldsymbol{\ell}$: $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

❹ IUA: $\forall a, b \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$

❺ ICS: $\forall \alpha \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha\boldsymbol{\ell}) \leq \mathrm{M}(\alpha\boldsymbol{\ell}')$

❻ Multiplicative Linearity: $\mathrm{M}(\alpha\boldsymbol{\ell}) = \alpha\mathrm{M}(\boldsymbol{\ell})$

❼ Unit Scale: $\mathrm{M}(\mathbf{1}) = 1$

All apply *equally well* to welfare and malfare

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
  Welfare
  Malfare
  Axiomatic Characterization

Estimation and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC Learning
  Computational Learnability

In Conclusion

# Transfer Axioms of Cardinal Welfare

❶ Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$

❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$

❸ Continuity: $\forall \boldsymbol{\ell}$: $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

❹ IUA: $\forall a, b \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$

❺ ICS: $\forall \alpha \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha\boldsymbol{\ell}) \leq \mathrm{M}(\alpha\boldsymbol{\ell}')$

❻ Multiplicative Linearity: $\mathrm{M}(\alpha\boldsymbol{\ell}) = \alpha\mathrm{M}(\boldsymbol{\ell})$

❼ Unit Scale: $\mathrm{M}(\mathbf{1}) = 1$

❽ Pigou-Dalton Principle: Suppose $\boldsymbol{\ell}, \boldsymbol{\ell}'$ s.t. $\mathrm{M}_1(\boldsymbol{\ell}) = \mathrm{M}_1(\boldsymbol{\ell}') = \mu$. Then
$$\bigwedge_{i \in 1,\ldots,g} \left( |\boldsymbol{\ell}_i' - \mu| \geq |\boldsymbol{\ell}_i - \mu| \right) \implies \mathrm{W}(\boldsymbol{\ell}') \leq \mathrm{W}(\boldsymbol{\ell})$$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

❶ Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$

❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$

❸ Continuity: $\forall \boldsymbol{\ell}$: $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

❹ IUA: $\forall a, b \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$

❺ ICS: $\forall \alpha \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha \boldsymbol{\ell}) \leq \mathrm{M}(\alpha \boldsymbol{\ell}')$

❻ Multiplicative Linearity: $\mathrm{M}(\alpha \boldsymbol{\ell}) = \alpha \mathrm{M}(\boldsymbol{\ell})$

❼ Unit Scale: $\mathrm{M}(\mathbf{1}) = 1$

❽ Pigou-Dalton Principle: Suppose $\boldsymbol{\ell}, \boldsymbol{\ell}'$ s.t. $\mathrm{M}_1(\boldsymbol{\ell}) = \mathrm{M}_1(\boldsymbol{\ell}') = \mu$. Then

$$\bigwedge_{i \in 1, \ldots, g} \left(|\boldsymbol{\ell}'_i - \mu| \geq |\boldsymbol{\ell}_i - \mu|\right) \implies \mathrm{W}(\boldsymbol{\ell}') \leq \mathrm{W}(\boldsymbol{\ell})$$

❾ Anti-Pigou-Dalton Principle: Suppose as in (8). Then require *inverse*

$$\bigwedge_{i \in 1, \ldots, g} \left(|\boldsymbol{\ell}'_i - \mu| \geq |\boldsymbol{\ell}_i - \mu|\right) \implies \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})$$

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
  Welfare
  Malfare
  Axiomatic Characterization

Estimation and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC Learning
  Computational Learnability

In Conclusion

Transfer Axioms of Cardinal Welfare

❶ Strict Monotonicity: $\forall \boldsymbol{\varepsilon} \succeq 0$ s.t. $\boldsymbol{\varepsilon} \neq \mathbf{0}$: $\mathrm{M}(\boldsymbol{\ell}) < \mathrm{M}(\boldsymbol{\ell} + \boldsymbol{\varepsilon})$

❷ Symmetry: $\forall$ permutations $\pi$: $\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}(\pi(\boldsymbol{\ell}))$

❸ Continuity: $\forall \boldsymbol{\ell}$: $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \leq \mathrm{M}(\boldsymbol{\ell})\}$ and $\{\boldsymbol{\ell}' \mid \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})\}$ are *closed sets*

❹ IUA: $\forall a, b \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}, a) \leq \mathrm{M}(\boldsymbol{\ell}', a) \Leftrightarrow \mathrm{M}(\boldsymbol{\ell}, b) \leq \mathrm{M}(\boldsymbol{\ell}', b)$

❺ ICS: $\forall \alpha \in \mathbb{R}_+$: $\mathrm{M}(\boldsymbol{\ell}) \leq \mathrm{M}(\boldsymbol{\ell}') \implies \mathrm{M}(\alpha\boldsymbol{\ell}) \leq \mathrm{M}(\alpha\boldsymbol{\ell}')$

❻ Multiplicative Linearity: $\mathrm{M}(\alpha\boldsymbol{\ell}) = \alpha\mathrm{M}(\boldsymbol{\ell})$

❼ Unit Scale: $\mathrm{M}(\mathbf{1}) = 1$

❽ Pigou-Dalton Principle: Suppose $\boldsymbol{\ell}, \boldsymbol{\ell}'$ s.t. $\mathrm{M}_1(\boldsymbol{\ell}) = \mathrm{M}_1(\boldsymbol{\ell}') = \mu$. Then

$$\bigwedge_{i \in 1, \dots, g} \left( |\boldsymbol{\ell}'_i - \mu| \geq |\boldsymbol{\ell}_i - \mu| \right) \implies \mathrm{W}(\boldsymbol{\ell}') \leq \mathrm{W}(\boldsymbol{\ell})$$

❾ Anti-Pigou-Dalton Principle: Suppose as in (8). Then require *inverse*

$$\bigwedge_{i \in 1, \dots, g} \left( |\boldsymbol{\ell}'_i - \mu| \geq |\boldsymbol{\ell}_i - \mu| \right) \implies \mathrm{M}(\boldsymbol{\ell}') \geq \mathrm{M}(\boldsymbol{\ell})$$

| | 😐😐😐😐 | 😄😐😐🙁 |
|---|---|---|
| Welfare | $w$ | $\leq w$ |
| Malfare | $m$ | $\geq m$ |

1: Monotonicity
2: Symmetry
3: Continuity
4: IUA

Mono. in G. Mean
$\mathrm{M}(\boldsymbol{\ell}) = (F \circ \mathrm{M}_f)(\boldsymbol{\ell})$

1: Monotonicity
2: Symmetry
3: Continuity
4: IUA

5: ICS

Mono. in G. Mean
$\mathrm{M}(\boldsymbol{\ell}) = (F \circ \mathrm{M}_f)(\boldsymbol{\ell})$

Mono. in $p$-Mean
$\mathrm{M}(\boldsymbol{\ell}) = (F \circ \mathrm{M}_p)(\boldsymbol{\ell})$

1: Monotonicity
2: Symmetry
3: Continuity
4: IUA

5: ICS

7: Unit-Scale

6: $\times$ Linearity

Mono. in G. Mean
$\mathrm{M}(\boldsymbol{\ell}) = (F \circ \mathrm{M}_f)(\boldsymbol{\ell})$

Mono. in $p$-Mean
$\mathrm{M}(\boldsymbol{\ell}) = (F \circ \mathrm{M}_p)(\boldsymbol{\ell})$

$p$-Mean
$\mathrm{M}(\boldsymbol{\ell}) = \mathrm{M}_p(\boldsymbol{\ell})$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare

Welfare

Malfare

Axiomatic
Characterization

Estimation
and Inference

Linear Classifiers

Statistical Estimation

Fair PAC
Learning

Computational
Learnability

In Conclusion

# Properties of Welfare and Malfare



1: Monotonicity
2: Symmetry
3: Continuity
4: IUA

5: ICS

7: Unit-Scale

6: $\times$ Linearity

Mono. in G. Mean
$\text{M}(\boldsymbol{\ell}) = (F \circ \text{M}_f)(\boldsymbol{\ell})$

Mono. in $p$-Mean
$\text{M}(\boldsymbol{\ell}) = (F \circ \text{M}_p)(\boldsymbol{\ell})$

$p$-Mean
$\text{M}(\boldsymbol{\ell}) = \text{M}_p(\boldsymbol{\ell})$

8: PD

9: APD

$p \leq 1$
Fair Welfare

$p \geq 1$
Fair Malfare

# Second Canto: Estimation, Inference, and Fair Machine Learning

## The Statistics of Fair Machine Learning as Malfare Minimization

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

Learning Linear Classifiers

❶ *How much* training data do we need?
- Concentration inequalities
- Vapnik-Chervonenkis dimension
- Rademacher averages

❶ *How much* training data do we need?
- Concentration inequalities
- Vapnik-Chervonenkis dimension
- Rademacher averages

❶ *How much* training data do we need?
  - Concentration inequalities
  - Vapnik-Chervonenkis dimension
  - Rademacher averages
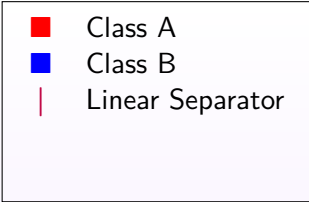
❷ How can we *learn* from the data?
  - *Empirical Risk Minimization*: select $h \in \mathcal{H}$ to optimize

$$\hat{\mathrm{R}}(h) \doteq \frac{1}{m} \sum_{j=1}^{m} \ell(h(\boldsymbol{x}_j), \boldsymbol{y}_j)$$



| | |
|---|---|
| ● | Class A |
| ● | Class B |

❶ *How much* training data do we need?
  - Concentration inequalities
  - Vapnik-Chervonenkis dimension
  - Rademacher averages

❷ How can we *learn* from the data?
  - *Empirical Risk Minimization*: select $h \in \mathcal{H}$ to optimize
  $$\hat{\mathrm{R}}(h) \doteq \frac{1}{m} \sum_{j=1}^{m} \ell(h(\boldsymbol{x}_j), \boldsymbol{y}_j)$$



| | |
|---|---|
| • | Class A |
| • | Class B |
| \| | Linear Separator |

- What changes with *multiple groups* $(\boldsymbol{x}_{1:g,1:m}, \boldsymbol{y}_{1:g,1:m})$?

| ■ | Class A |
|---|---------|
| ■ | Class B |
| | | Linear Separator |

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare

Welfare

Malfare

Axiomatic
Characterization

Estimation
and Inference

Linear Classifiers

Statistical Estimation

Fair PAC
Learning

Computational
Learnability

In Conclusion

- What changes with *multiple groups* $(\boldsymbol{x}_{1:g,1:m}, \boldsymbol{y}_{1:g,1:m})$?
- We can handle each group individually:

$$\hat{\mathrm{R}}(h; \boldsymbol{x}_i, \boldsymbol{y}_i) \doteq \frac{1}{m}\sum_{j=1}^{m}\ell\big(h(\boldsymbol{x}_{i,j}), \boldsymbol{y}_{i,j}\big); \quad \forall i: \ \hat{h}_i \doteq \operatorname*{argmin}_{h \in \mathcal{H}} \hat{\mathrm{R}}(h; \boldsymbol{x}_i, \boldsymbol{y}_i)$$



| | |
|---|---|
| 🟥 | Class A |
| 🟦 | Class B |
| \| | Linear Separator |
| ☼ | Group 1 |

- What changes with *multiple groups* $(\boldsymbol{x}_{1:g,1:m}, \boldsymbol{y}_{1:g,1:m})$?
- We can handle each group individually:

$$\hat{\mathrm{R}}(h; \boldsymbol{x}_i, \boldsymbol{y}_i) \doteq \frac{1}{m}\sum_{j=1}^{m}\ell\big(h(\boldsymbol{x}_{i,j}), \boldsymbol{y}_{i,j}\big); \quad \forall i: \ \hat{h}_i \doteq \underset{h \in \mathcal{H}}{\operatorname{argmin}} \ \hat{\mathrm{R}}(h; \boldsymbol{x}_i, \boldsymbol{y}_i)$$



| ■ | Class A |
| --- | --- |
| ■ | Class B |
| | | Linear Separator |
| ☾ | Group 2 |

- What changes with *multiple groups* $(\boldsymbol{x}_{1:g,1:m}, \boldsymbol{y}_{1:g,1:m})$?
- We can handle each group individually:

$$\hat{\mathrm{R}}(h; \boldsymbol{x}_i, \boldsymbol{y}_i) \doteq \frac{1}{m}\sum_{j=1}^{m}\ell\big(h(\boldsymbol{x}_{i,j}), \boldsymbol{y}_{i,j}\big); \quad \forall i: \ \hat{h}_i \doteq \operatorname*{argmin}_{h \in \mathcal{H}} \hat{\mathrm{R}}(h; \boldsymbol{x}_i, \boldsymbol{y}_i)$$

- What is the best classifier *overall*?



| ■ | Class A |
| | | Class B |
| | | Linear Separator |
| ☼ | Group 1 |
| ☾ | Group 2 |

- What changes with *multiple groups* $(\boldsymbol{x}_{1:g,1:m}, \boldsymbol{y}_{1:g,1:m})$?
- We can handle each group individually:

$$\hat{\mathrm{R}}(h; \boldsymbol{x}_i, \boldsymbol{y}_i) \doteq \frac{1}{m}\sum_{j=1}^{m}\ell\big(h(\boldsymbol{x}_{i,j}), \boldsymbol{y}_{i,j}\big); \quad \forall i: \ \hat{h}_i \doteq \operatorname*{argmin}_{h \in \mathcal{H}} \hat{\mathrm{R}}(h; \boldsymbol{x}_i, \boldsymbol{y}_i)$$

- What is the best classifier *overall*?
  - *Empirical malfare minimization* $\hat{h} \doteq \operatorname*{argmin}_{h \in \mathcal{H}} \mathrm{M}\left(\hat{\mathrm{R}}(h; \boldsymbol{x}_1, \boldsymbol{y}_1), \hat{\mathrm{R}}(h; \boldsymbol{x}_2, \boldsymbol{y}_2)\right)$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare

Welfare

Malfare

Axiomatic
Characterization

Estimation
and Inference

Linear Classifiers

Statistical Estimation

Fair PAC
Learning

Computational
Learnability

In Conclusion

- Suppose *sample mean* $\hat{\boldsymbol{\ell}}_i \doteq \frac{1}{m} \sum_{j=1}^{m} \ell(\boldsymbol{x}_{i,j})$, *true mean* $\boldsymbol{\ell}_i \doteq \mathop{\mathbb{E}}\limits_{x \sim \mathcal{D}_i}[\ell(x)]$

- Suppose *sample mean* $\hat{\ell}_i \doteq \frac{1}{m} \sum_{j=1}^{m} \ell(\boldsymbol{x}_{i,j})$, *true mean* $\boldsymbol{\ell}_i \doteq \mathop{\mathbb{E}}_{x \sim \mathcal{D}_i}[\ell(x)]$

- By *continuity* and the *law of large numbers*:

$$\lim_{m \to \infty} \mathbb{M}(\hat{\boldsymbol{\ell}}) = \mathbb{M}(\boldsymbol{\ell})$$

- Suppose *sample mean* $\hat{\boldsymbol{\ell}}_i \doteq \dfrac{1}{m} \sum_{j=1}^{m} \ell(\boldsymbol{x}_{i,j})$, *true mean* $\boldsymbol{\ell}_i \doteq \mathop{\mathbb{E}}_{x \sim \mathcal{D}_i}[\ell(x)]$

- By *continuity* and the *law of large numbers*:

$$\lim_{m \to \infty} \mathbb{M}(\hat{\boldsymbol{\ell}}) = \mathbb{M}(\boldsymbol{\ell})$$

- For *finite sample size* $m$
  - $\mathbb{E}[\mathbb{M}(\hat{\boldsymbol{\ell}})] \neq \mathbb{M}(\boldsymbol{\ell})$
  - $\mathbb{M}(\hat{\boldsymbol{\ell}})$ is a *biased estimator* of $\mathbb{M}(\boldsymbol{\ell})$!

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
  Welfare
  Malfare
  Axiomatic
  Characterization

Estimation
and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC
Learning
  Computational
  Learnability

In Conclusion

- Suppose *sample mean* $\hat{\boldsymbol{\ell}}_i \doteq \frac{1}{m} \sum_{j=1}^{m} \ell(\boldsymbol{x}_{i,j})$, *true mean* $\boldsymbol{\ell}_i \doteq \mathop{\mathbb{E}}_{x \sim \mathcal{D}_i}[\ell(x)]$

- By *continuity* and the *law of large numbers*:

$$\lim_{m \to \infty} \mathbb{M}(\hat{\boldsymbol{\ell}}) = \mathbb{M}(\boldsymbol{\ell})$$

- For finite sample size $m$
  - $\mathbb{E}[\mathbb{M}(\hat{\boldsymbol{\ell}})] \neq \mathbb{M}(\boldsymbol{\ell})$
  - $\mathbb{M}(\hat{\boldsymbol{\ell}})$ is a *biased estimator* of $\mathbb{M}(\boldsymbol{\ell})$!

### Theorem (A Hoeffding-Type Malfare-Estimation Bound)

*Suppose fair malfare* $\mathbb{M}_p(\cdot)$ *($p \geq 1$), $g$ groups, and loss range $r$. Then with probability at least $1 - \delta$*

$$\left| \mathbb{M}(\boldsymbol{\ell}) - \mathbb{M}(\hat{\boldsymbol{\ell}}) \right| \leq r \sqrt{\frac{\ln \frac{2g}{\delta}}{2m}}$$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
  Welfare
  Malfare
  Axiomatic
  Characterization

Estimation
and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC
Learning
  Computational
  Learnability

In Conclusion

# Statistical Estimation (contd.)

## Theorem (A Bernstein-Type Malfare-Estimation Bound)

*Suppose fair malfare $\Lambda_p(\cdot)$ ($p \geq 1$), $g$ groups, loss range $r$, and maximum variance $v_{\max}$. Then with probability at least $1 - \delta$ over sampling, we have*

$$\left| \Lambda(\boldsymbol{\ell}) - \Lambda(\hat{\boldsymbol{\ell}}) \right| \leq \underbrace{\frac{r \ln \frac{2g}{\delta}}{3m}}_{\text{SCALE TERM}} + \underbrace{\sqrt{\frac{v_{\max} \ln \frac{2g}{\delta}}{2m}}}_{\text{VARIANCE TERM}}$$

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
 Welfare
 Malfare
 Axiomatic Characterization

Estimation and Inference
 Linear Classifiers
 Statistical Estimation

Fair PAC Learning
 Computational Learnability

In Conclusion

# Statistical Estimation (contd.)

## Theorem (A Bernstein-Type Malfare-Estimation Bound)

*Suppose fair malfare $\Lambda_p(\cdot)$ ($p \geq 1$), $g$ groups, loss range $r$, and maximum variance $v_{\max}$. Then with probability at least $1 - \delta$ over sampling, we have*

$$\left| \Lambda(\boldsymbol{\ell}) - \Lambda(\hat{\boldsymbol{\ell}}) \right| \leq \underbrace{\frac{r \ln \frac{2g}{\delta}}{3m}}_{\text{SCALE TERM}} + \underbrace{\sqrt{\frac{v_{\max} \ln \frac{2g}{\delta}}{2m}}}_{\text{VARIANCE TERM}}$$

- Can show similar bounds for *any concentration inequality*
- *Uniform bounds* for a learnable family $\mathcal{H}$ with R a d e m a c h e r  a v e r a g e s

$$\sup_{h \in \mathcal{H}} \left| \Lambda(\boldsymbol{\ell}(h)) - \Lambda(\hat{\boldsymbol{\ell}}(h)) \right| \leq \max_{i \in 1, \ldots, g} 2\mathfrak{R}_m(\mathcal{F}, \mathcal{D}_i) + \varepsilon \in \Theta\left( \frac{r \ln \frac{g}{\delta}}{m} + \sqrt{\frac{v_{\max} \ln \frac{g}{\delta}}{m}} \right)$$

## Theorem (A Bernstein-Type Malfare-Estimation Bound)

*Suppose fair malfare $\Lambda_p(\cdot)$ ($p \geq 1$), $g$ groups, loss range $r$, and maximum variance $v_{\max}$. Then with probability at least $1 - \delta$ over sampling, we have*

$$\left| \Lambda(\boldsymbol{\ell}) - \Lambda(\hat{\boldsymbol{\ell}}) \right| \leq \underbrace{\frac{r \ln \frac{2g}{\delta}}{3m}}_{\text{SCALE TERM}} + \underbrace{\sqrt{\frac{v_{\max} \ln \frac{2g}{\delta}}{2m}}}_{\text{VARIANCE TERM}}$$

- Can show similar bounds for *any concentration inequality*
- *Uniform bounds* for a learnable family $\mathcal{H}$ with R a d e m a c h e r   a v e r a g e s

$$\sup_{h \in \mathcal{H}} \left| \Lambda(\boldsymbol{\ell}(h)) - \Lambda(\hat{\boldsymbol{\ell}}(h)) \right| \leq \max_{i \in 1, \dots, g} 2\mathfrak{R}_m(\mathcal{F}, \mathcal{D}_i) + \varepsilon \in \Theta \left( \frac{r \ln \frac{g}{\delta}}{m} + \sqrt{\frac{v_{\max} \ln \frac{g}{\delta}}{m}} \right)$$

- Control *overfitting* in *machine learning*
  - Finite $\mathcal{H}$, bounded Lipschitz families
  - Bounded linear regression, finite-dimensional linear classifiers,

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
    Welfare
    Malfare
    Axiomatic Characterization

Estimation and Inference
    Linear Classifiers
    **Statistical Estimation**

Fair PAC Learning
    Computational Learnability

In Conclusion

Statistical Estimation (contd.)

## Theorem (A Bernstein-Type Malfare-Estimation Bound)

*Suppose fair malfare* $\Lambda_p(\cdot)$ *($p \geq 1$), $g$ groups, loss range $r$, and maximum variance* $v_{\max}$. *Then with probability at least $1 - \delta$ over sampling, we have*

$$\left| \Lambda(\boldsymbol{\ell}) - \Lambda(\hat{\boldsymbol{\ell}}) \right| \leq \underbrace{\frac{r \ln \frac{2g}{\delta}}{3m}}_{\text{SCALE TERM}} + \underbrace{\sqrt{\frac{v_{\max} \ln \frac{2g}{\delta}}{2m}}}_{\text{VARIANCE TERM}}$$

- Can show similar bounds for *any concentration inequality*
- *Uniform bounds* for a learnable family $\mathcal{H}$ with R a d e m a c h e r   a v e r a g e s

$$\sup_{h \in \mathcal{H}} \left| \Lambda(\boldsymbol{\ell}(h)) - \Lambda(\hat{\boldsymbol{\ell}}(h)) \right| \leq \max_{i \in 1, \ldots, g} 2\mathfrak{R}_m(\mathcal{F}, \mathcal{D}_i) + \varepsilon \in \Theta \left( \frac{r \ln \frac{g}{\delta}}{m} + \sqrt{\frac{v_{\max} \ln \frac{g}{\delta}}{m}} \right)$$

- Control *overfitting* in *machine learning*
    - Finite $\mathcal{H}$, bounded Lipschitz families
    - Bounded linear regression, finite-dimensional linear classifiers, Generalized linear models, support vector machines, multiple kernel learning, bounded depth decision trees, rank-constrained matrix factorization, neural networks,

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
  Welfare
  Malfare
  Axiomatic Characterization
Estimation and Inference
  Linear Classifiers
  Statistical Estimation
Fair PAC Learning
  Computational Learnability
In Conclusion

Statistical Estimation (contd.)

## Theorem (A Bernstein-Type Malfare-Estimation Bound)

*Suppose fair malfare $\Lambda_p(\cdot)$ ($p \geq 1$), $g$ groups, loss range $r$, and maximum variance $v_{\max}$. Then with probability at least $1 - \delta$ over sampling, we have*

$$\left| \Lambda(\boldsymbol{\ell}) - \Lambda(\hat{\boldsymbol{\ell}}) \right| \leq \underbrace{\frac{r \ln \frac{2g}{\delta}}{3m}}_{\text{SCALE TERM}} + \underbrace{\sqrt{\frac{v_{\max} \ln \frac{2g}{\delta}}{2m}}}_{\text{VARIANCE TERM}}$$

- Can show similar bounds for *any concentration inequality*
- *Uniform bounds* for a learnable family $\mathcal{H}$ with R a d e m a c h e r  a v e r a g e s

$$\sup_{h \in \mathcal{H}} \left| \Lambda(\boldsymbol{\ell}(h)) - \Lambda(\hat{\boldsymbol{\ell}}(h)) \right| \leq \max_{i \in 1, \ldots, g} 2\mathfrak{R}_m(\mathcal{F}, \mathcal{D}_i) + \varepsilon \in \Theta\left( \frac{r \ln \frac{g}{\delta}}{m} + \sqrt{\frac{v_{\max} \ln \frac{g}{\delta}}{m}} \right)$$

- - Control *overfitting* in *machine learning*
    - Finite $\mathcal{H}$, bounded Lipschitz families
    - Bounded linear regression, finite-dimensional linear classifiers, Generalized linear models, support vector machines, multiple kernel learning, bounded depth decision trees, rank-constrained matrix factorization, neural networks, (constrained) Boolean formulae, boosting methods, convex ensemble methods, learning distance metrics, . . .

# Third Canto: Fair Probably Approximately Correct Learning

### A generic theory of fair statistical and computational learning

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
Welfare
Malfare
Axiomatic Characterization

Estimation and Inference
Linear Classifiers
Statistical Estimation

Fair PAC Learning
Computational Learnability

In Conclusion

# Classical Statistical Learning Theory

- Consider *linear classification*: $\mathcal{H}_d \doteq \left\{ \vec{x} \mapsto \mathrm{sgn}(\vec{w} \cdot \vec{x}) \,\middle|\, \vec{w} \in \mathbb{R}^d \right\}$
  - Optimize *risk* $\mathbb{E}_{(x,y) \sim \mathcal{D}}\big[\ell(y, h(x))\big]$, for 0-1 loss $\ell(y, \hat{y}) = 1 - \mathbb{1}_y(\hat{y})$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
  Welfare
  Malfare
  Axiomatic
  Characterization

Estimation
and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC
Learning
  Computational
  Learnability

In Conclusion

# Classical Statistical Learning Theory

- Consider *linear classification*: $\mathcal{H}_d \doteq \left\{ \vec{x} \mapsto \mathrm{sgn}(\vec{w} \cdot \vec{x}) \,\middle|\, \vec{w} \in \mathbb{R}^d \right\}$
  - Optimize *risk* $\mathbb{E}_{(x,y)\sim\mathcal{D}}\big[\ell(y, h(x))\big]$, for 0-1 loss $\ell(y, \hat{y}) = 1 - \mathbb{1}_y(\hat{y})$
- Can this class be *efficiently learned*?       What does that even mean?

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare

Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# Classical Statistical Learning Theory

- Consider *linear classification*: $\mathcal{H}_d \doteq \left\{ \vec{x} \mapsto \text{sgn}(\vec{w} \cdot \vec{x}) \,\middle|\, \vec{w} \in \mathbb{R}^d \right\}$
  - Optimize *risk* $\mathbb{E}_{(x,y) \sim \mathcal{D}}\big[\ell(y, h(x))\big]$, for 0-1 loss $\ell(y, \hat{y}) = 1 - \mathbb{1}_y(\hat{y})$
- Can this class be *efficiently learned*?     What does that even mean?

## Definition (PAC Learning)

Suppose

❶ Hypothesis class $\mathcal{H} \subseteq \mathcal{X} \to \mathcal{Y}$ ❷ Loss function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_{0+}$

$\mathcal{H}$ is *PAC-learnable* w.r.t. $\ell$ iff $\exists$ algorithm $A$ s.t. $\forall$

❶ distributions $\mathcal{D}$ over $\mathcal{X} \times \mathcal{Y}$ ❸ failure probabilities $\delta \in (0, 1)$

❷ additive errors $\varepsilon > 0$

$A$ can identify a hypothesis $\hat{h} \in \mathcal{H}$ s.t.

❶ $A$ has $\text{m}(\varepsilon, \delta) < \infty$ sample complexity

❷ with probability at least $1 - \delta$, $\hat{h}$ obeys

$$\mathbb{E}_{(x,y) \sim \mathcal{D}}\big[\ell(y, \hat{h}(x))\big] \leq \underset{h^* \in \mathcal{H}}{\text{argmin}}\; \mathbb{E}_{(x,y) \sim \mathcal{D}}\big[\ell(y, h^*(x))\big] + \varepsilon$$

- Consider *linear classification*: $\mathcal{H}_d \doteq \left\{ \vec{x} \mapsto \mathrm{sgn}(\vec{w} \cdot \vec{x}) \,\middle|\, \vec{w} \in \mathbb{R}^d \right\}$
  - Optimize *risk* $\mathbb{E}_{(x,y)\sim\mathcal{D}}\big[\ell(y, h(x))\big]$, for 0-1 loss $\ell(y, \hat{y}) = 1 - \mathbb{1}_y(\hat{y})$
- Can this class be *efficiently learned*?     What does that even mean?

## Definition (PAC Learning)

Suppose

❶ Hypothesis class $\mathcal{H} \subseteq \mathcal{X} \to \mathcal{Y}$     ❷ Loss function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_{0+}$

$\mathcal{H}$ is *PAC-learnable* w.r.t. $\ell$ iff $\exists$ algorithm $A$ s.t. $\forall$

❶ distributions $\mathcal{D}$ over $\mathcal{X} \times \mathcal{Y}$     ❸ failure probabilities $\delta \in (0, 1)$

❷ additive errors $\varepsilon > 0$

$A$ can identify a hypothesis $\hat{h} \in \mathcal{H}$ s.t.

❶ $A$ has $\mathrm{m}(\varepsilon, \delta) < \infty$ sample complexity

❷ with probability at least $1 - \delta$, $\hat{h}$ obeys

$$\mathbb{E}_{(x,y)\sim\mathcal{D}}\big[\ell(y, \hat{h}(x))\big] \leq \operatorname*{argmin}_{h^* \in \mathcal{H}} \mathbb{E}_{(x,y)\sim\mathcal{D}}\big[\ell(y, h^*(x))\big] + \varepsilon$$

- May also consider *efficient PAC-learnable*: require *poly-time* $A$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# Fairness and Statistical Learning Theory

## Definition (Fair-PAC Learning)

Suppose
1. Hypothesis class $\mathcal{H} \subseteq \mathcal{X} \to \mathcal{Y}$
2. Loss function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_{0+}$

$\mathcal{H}$ is *fair PAC-learnable* w.r.t. $\ell$ iff $\exists$ algorithm $A$ s.t. $\forall$
1. distributions $\mathcal{D}_{1:g}$ over $(\mathcal{X} \times \mathcal{Y})^g$
2. fair malfare functions $\mathrm{M}(\cdot)$
3. additive errors $\varepsilon > 0$
4. failure probabilities $\delta \in (0, 1)$

$A$ can identify a hypothesis $\hat{h} \in \mathcal{H}$ s.t.

1. $A$ has $\mathrm{m}(\varepsilon, \delta, g)$ sample complexity
2. with probability at least $1 - \delta$, $\hat{h}$ obeys

$$\mathrm{M}\left(\mathbb{E}_{(x,y) \sim \mathcal{D}_1}\big[\ell(y, \hat{h}(x))\big], \dots\right) \leq \operatorname*{argmin}_{h^* \in \mathcal{H}} \mathrm{M}\left(\mathbb{E}_{(x,y) \sim \mathcal{D}_1}\big[\ell(y, h^*(x))\big], \dots\right) + \varepsilon$$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

Fairness and Statistical Learning Theory

## Definition (Fair-PAC Learning)

Suppose
- ❶ Hypothesis class $\mathcal{H} \subseteq \mathcal{X} \to \mathcal{Y}$
- ❷ Loss function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_{0+}$

$\mathcal{H}$ is *fair PAC-learnable* w.r.t. $\ell$ iff $\exists$ algorithm $A$ s.t. $\forall$

- ❶ distributions $\mathcal{D}_{1:g}$ over $(\mathcal{X} \times \mathcal{Y})^g$
- ❷ fair malfare functions $\mathbb{M}(\cdot)$
- ❸ additive errors $\varepsilon > 0$
- ❹ failure probabilities $\delta \in (0, 1)$

$A$ can identify a hypothesis $\hat{h} \in \mathcal{H}$ s.t.

- ❶ $A$ has $\mathrm{m}(\varepsilon, \delta, g)$ sample complexity
- ❷ with probability at least $1 - \delta$, $\hat{h}$ obeys

$$\mathbb{M}\left(\mathop{\mathbb{E}}_{(x,y)\sim\mathcal{D}_1}\big[\ell(y, \hat{h}(x))\big], \dots\right) \leq \operatorname*{argmin}_{h^* \in \mathcal{H}} \mathbb{M}\left(\mathop{\mathbb{E}}_{(x,y)\sim\mathcal{D}_1}\big[\ell(y, h^*(x))\big], \dots\right) + \varepsilon$$

- Do we capture a valuable, generic notion of fair learning?
  - Axiomatic social planning problem motivation

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
    Welfare
    Malfare
    Axiomatic
    Characterization

Estimation
and Inference
    Linear Classifiers
    Statistical Estimation

Fair PAC
Learning
    Computational
    Learnability

In Conclusion

## Fairness and Statistical Learning Theory

### Definition (Fair-PAC Learning)

Suppose

❶ Hypothesis class $\mathcal{H} \subseteq \mathcal{X} \to \mathcal{Y}$   ❷ Loss function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_{0+}$

$\mathcal{H}$ is *fair PAC-learnable* w.r.t. $\ell$ iff $\exists$ algorithm $A$ s.t. $\forall$

❶ distributions $\mathcal{D}_{1:g}$ over $(\mathcal{X} \times \mathcal{Y})^g$   ❸ additive errors $\varepsilon > 0$

❷ fair malfare functions $\mathbb{M}(\cdot)$   ❹ failure probabilities $\delta \in (0, 1)$

$A$ can identify a hypothesis $\hat{h} \in \mathcal{H}$ s.t.

❶ $A$ has $\mathrm{m}(\varepsilon, \delta, g)$ sample complexity

❷ with probability at least $1 - \delta$, $\hat{h}$ obeys

$$\mathbb{M}\left( \mathop{\mathbb{E}}_{(x,y) \sim \mathcal{D}_1} \left[ \ell(y, \hat{h}(x)) \right], \dots \right) \leq \underset{h^* \in \mathcal{H}}{\mathrm{argmin}} \, \mathbb{M}\left( \mathop{\mathbb{E}}_{(x,y) \sim \mathcal{D}_1} \left[ \ell(y, h^*(x)) \right], \dots \right) + \varepsilon$$

- Do we capture a valuable, generic notion of fair learning?
  - Axiomatic social planning problem motivation
- Do practical fair PAC-learning algorithms exist?

Axiomatic Fair Learning

Cyrus Cousins

NeurIPS2021

Philosophy, Welfare, and Malfare
  Welfare
  Malfare
  Axiomatic Characterization

Estimation and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC Learning
  Computational Learnability

In Conclusion

## Fairness and Statistical Learning Theory

### Definition (Fair-PAC Learning)

Suppose
1. Hypothesis class $\mathcal{H} \subseteq \mathcal{X} \to \mathcal{Y}$    2. Loss function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_{0+}$

$\mathcal{H}$ is *fair PAC-learnable* w.r.t. $\ell$ iff $\exists$ algorithm $A$ s.t. $\forall$

1. distributions $\mathcal{D}_{1:g}$ over $(\mathcal{X} \times \mathcal{Y})^g$    3. additive errors $\varepsilon > 0$
2. fair malfare functions $\mathbb{M}(\cdot)$    4. failure probabilities $\delta \in (0, 1)$

$A$ can identify a hypothesis $\hat{h} \in \mathcal{H}$ s.t.

1. $A$ has $\mathrm{m}(\varepsilon, \delta, g)$ sample complexity
2. with probability at least $1 - \delta$, $\hat{h}$ obeys

$$\mathbb{M}\left(\mathop{\mathbb{E}}_{(x,y)\sim\mathcal{D}_1}\big[\ell(y, \hat{h}(x))\big], \dots\right) \leq \operatorname*{argmin}_{h^* \in \mathcal{H}} \mathbb{M}\left(\mathop{\mathbb{E}}_{(x,y)\sim\mathcal{D}_1}\big[\ell(y, h^*(x))\big], \dots\right) + \varepsilon$$

- Do we capture a valuable, generic notion of fair learning?
  - Axiomatic social planning problem motivation
- Do practical fair PAC-learning algorithms exist?
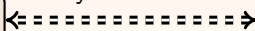- Can we theoretically relate PAC and fair-PAC learning?

Axiomatic Fair
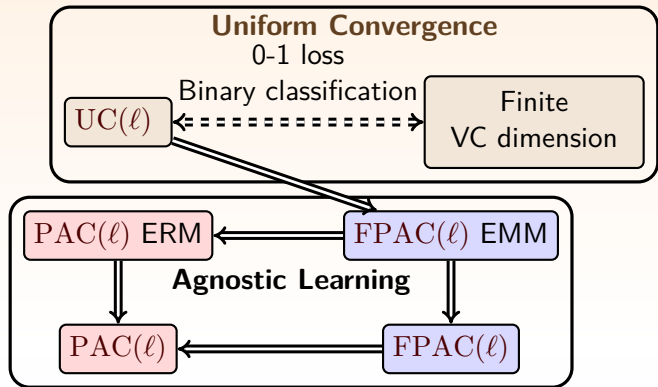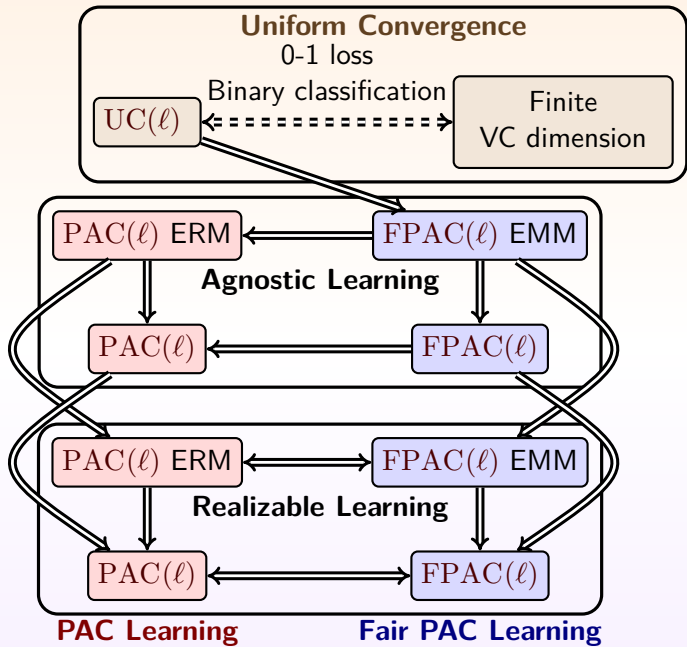Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

Fairness and Statistical Learning Theory

## Definition (Fair-PAC Learning)

Suppose
1. Hypothesis class $\mathcal{H} \subseteq \mathcal{X} \rightarrow \mathcal{Y}$
2. Loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{0+}$

$\mathcal{H}$ is *fair PAC-learnable* w.r.t. $\ell$ iff $\exists$ algorithm $A$ s.t. $\forall$
1. distributions $\mathcal{D}_{1:g}$ over $(\mathcal{X} \times \mathcal{Y})^g$
2. fair malfare functions $\mathbb{M}(\cdot)$
3. additive errors $\varepsilon > 0$
   failure probabilities $\delta \in (0, 1)$

$A$ can identify a hypothesis $\hat{h} \in \mathcal{H}$ s.t.
1. $A$ has $\mathrm{m}(\varepsilon, \delta, g)$ sample complexity
2. with probability

Yes.
Why else would I ask?

$$\mathbb{M}\left(\mathop{\mathbb{E}}_{(x,y)\sim\mathcal{D}_1}\big[\ell(y, \hat{h}(x))\big], \dots\right) \leq \operatorname*{argmin}_{h^*\in\mathcal{H}} \mathbb{M}\left(\mathop{\mathbb{E}}_{(x,y)\sim\mathcal{D}_1}\big[\ell(y, h^*(x))\big], \dots\right) + \varepsilon$$

- Do we capture a valuable, generic notion of fair learning?
  - Axiomatic social planning problem motivation
- Do practical fair PAC-learning algorithms exist?
- Can we theoretically relate PAC and fair-PAC learning?
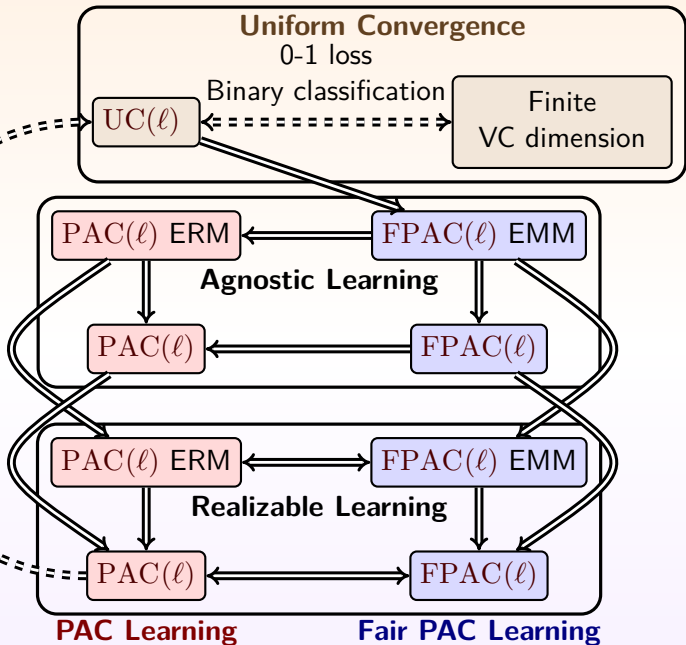
**Uniform Convergence**
0-1 loss
Binary classification

$\text{UC}(\ell)$ ⟸========⟹ Finite VC dimension

- Can we construct *computationally-efficient* FPAC learners from PAC learners?
  - *Efficient* means $\mathrm{Poly}(\frac{1}{\varepsilon}, \frac{1}{\delta}, g)$ sample complexity
  - Realizable case: reduction preserves polynomial-time complexity
  - Agnostic case: Cyrus has no answer (yet)

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
  Welfare
  Malfare
  Axiomatic
  Characterization

Estimation
and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC
Learning
  Computational
  Learnability

In Conclusion

- Can we construct *computationally-efficient* FPAC learners from PAC learners?
  - *Efficient* means $\mathrm{Poly}(\frac{1}{\varepsilon}, \frac{1}{\delta}, g)$ sample complexity
  - Realizable case: reduction preserves polynomial-time complexity
  - Agnostic case: Cyrus has no answer (yet)
- Are conditions for efficient PAC sufficient for efficient FPAC?
  - This area looks promising!

- Can we construct *computationally-efficient* FPAC learners from PAC learners?
  - *Efficient* means $\mathrm{Poly}(\frac{1}{\varepsilon}, \frac{1}{\delta}, g)$ sample complexity
  - Realizable case: reduction preserves polynomial-time complexity
  - Agnostic case: Cyrus has no answer (yet)
- Are conditions for efficient PAC sufficient for efficient FPAC?
  - This area looks promising!
  - Efficiently Coverable Classes:
    - If we can *efficiently approximately enumerate* $\mathcal{H}$
    - And our loss-function is well-behaved
    - Then we can PAC or FPAC-learn in $\mathcal{H}$
    - Think "all separating hyperplanes of bounded dimension"

- Can we construct *computationally-efficient* FPAC learners from PAC learners?
  - *Efficient* means $\mathrm{Poly}(\frac{1}{\varepsilon}, \frac{1}{\delta}, g)$ sample complexity
  - Realizable case: reduction preserves polynomial-time complexity
  - Agnostic case: Cyrus has no answer (yet)
- Are conditions for efficient PAC sufficient for efficient FPAC?
  - This area looks promising!
  - Efficiently Coverable Classes:
    - If we can *efficiently approximately enumerate* $\mathcal{H}$
    - And our loss-function is well-behaved
    - Then we can PAC or FPAC-learn in $\mathcal{H}$
    - Think "all separating hyperplanes of bounded dimension"
  - Convex optimization:
    - Suppose *bounded* parameter space $\Theta$
    - Assume $\ell \circ h_\theta$ is *convex + Lipschitz continuous* in $\theta$
    - Then $\varepsilon$-empirical risk minimization requires *polynomial time*
    - Same for empirical malfare minimization (this work)

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# Convex Optimization

**Strategy**: Assume class $\ell \circ \mathcal{H}$ is:

❶ Uniformly Convergent

- Vapnik-Chervonenkis theory:
  *Uniform* bounds over distribution
  $\mathcal{D}$
- *Estimation error*: $\epsilon(m, \delta)$ s.t.

$$\mathbb{P}\left(\sup_{f \in \ell \circ \mathcal{H}} \left|\mathbb{E}[f] - \widehat{\mathbb{E}}[f]\right| \geq \epsilon(m, \delta)\right) \leq \delta$$

- *Sample complexity*

$$\mathrm{m}(\varepsilon, \delta) \doteq \operatorname{argmin}\left\{m : \epsilon(m, \delta) \leq \varepsilon\right\}$$
$$\in \operatorname{Poly}\left(\tfrac{1}{\varepsilon}, \tfrac{1}{\delta}\right)$$

❷ Bounded

- Bounded parameter space $\Theta \in \mathbb{R}^d$

❸ Lipschitz Continuous

- $\lambda_\ell$-Lipschitz loss $\ell$, $\lambda_{\mathcal{H}}$-Lipschitz $\mathcal{H}$

❹ Convex

- $\ell(\circ h(x; \theta), y)$ is convex in $\theta$ over $\Theta$

## The Algorithm

❶ Draw $\mathrm{m}(\tfrac{\varepsilon}{3}, \tfrac{\delta}{g})$ samples (per group)

❷ Define *empirical malfare* objective

$$f(\theta) \doteq \mathrm{M}_p\big(i \mapsto \hat{\mathrm{R}}(h(\cdot; \theta); \ell, \boldsymbol{x}_i, \boldsymbol{y}_i)\big)$$

❸ Iterations: $n \doteq \left(\tfrac{3 \operatorname{diam}(\Theta) \lambda_\ell \lambda_{\mathcal{H}}}{\varepsilon}\right)^2$

❹ Learning rate $\alpha \doteq \dfrac{\operatorname{diam}(\Theta)}{\lambda_\ell \lambda_{\mathcal{H}} \sqrt{n}} \approx \dfrac{\varepsilon}{3\lambda_\ell^2 \lambda_{\mathcal{H}}^2}$

❺ Shor's *projected subgradient* algorithm

$$\hat{\theta} \leftarrow \mathrm{PSG}(f, \Theta, n, \alpha)$$

❻ Return $h(\cdot; \hat{\theta})$

W.h.p., *estimation + optimization*
error don't exceed $\varepsilon$

Polynomial time + sample complexity

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

Convex Optimization

**Strategy**: Assume class $\ell \circ \mathcal{H}$ is:

❶ Uniformly Convergent
  - Vapnik-Chervonenkis theory:
    *Uniform* bounds over distribution $\mathcal{D}$
  - *Estimation error*: $\epsilon(m, \delta)$ s.t.

  $$\mathbb{P}\left(\sup_{f \in \ell \circ \mathcal{H}} \left|\mathbb{E}[f] - \widehat{\mathbb{E}}[f]\right| \geq \epsilon(m, \delta)\right) <$$

  - *Sample complexity*

  $$\mathrm{m}(\varepsilon, \delta) \doteq \mathrm{argmin}\left\{m : \epsilon(m, \delta) \leq \varepsilon\right\}$$
  $$\in \mathrm{Poly}\left(\frac{1}{\varepsilon}, \frac{1}{\delta}\right)$$

❷ Bounded
  - Bounded parameter space $\Theta \in \mathbb{R}^d$

❸ Lipschitz Continuous
  - $\lambda_\ell$-Lipschitz loss $\ell$, $\lambda_\mathcal{H}$-Lipschitz $\mathcal{H}$

❹ Convex
  - $\ell(\circ h(x; \theta), y)$ is convex in $\theta$ over $\Theta$

**The Algorithm**

❶ Draw $\mathrm{m}(\frac{\varepsilon}{3}, \frac{\delta}{g})$ samples (per group)

❷ ... fin... *pirical malfare* objective

$$\hat{\mathrm{R}}(h(\cdot; \theta); \ell, \boldsymbol{x}_i, \boldsymbol{y}_i))$$

$$\left(\frac{3 \operatorname{diam}(\Theta) \lambda_\ell \lambda_\mathcal{H}}{\varepsilon}\right)^2$$

$$\frac{\operatorname{diam}(\Theta)}{\lambda_\ell \lambda_\mathcal{H} \sqrt{n}} \approx \frac{\varepsilon}{3 \lambda_\ell^2 \lambda_\mathcal{H}^2}$$
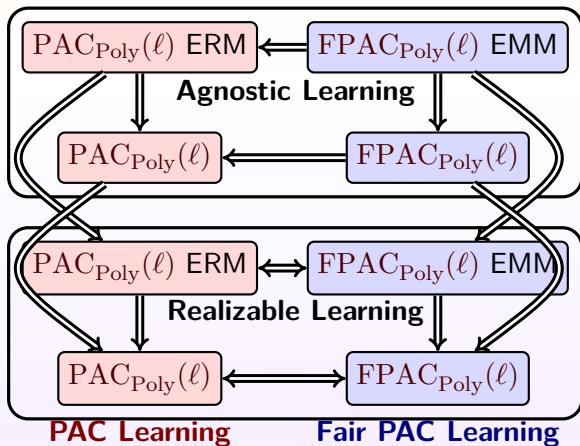
*subgradient* algorithm
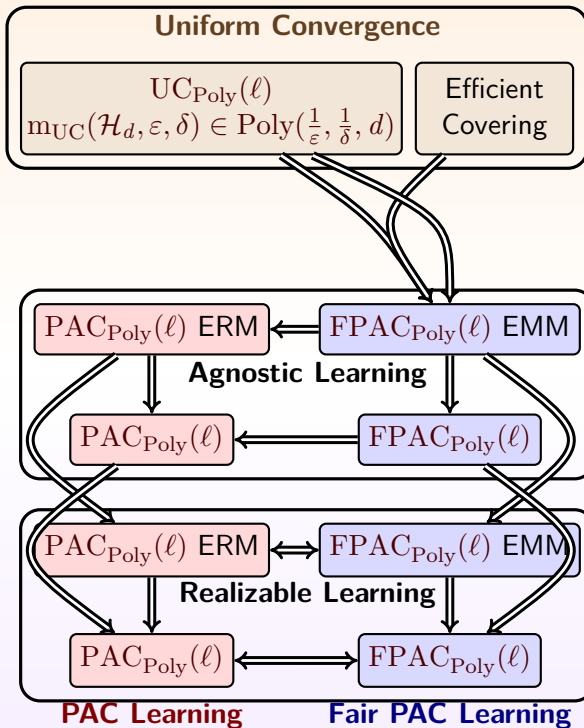
$$\leftarrow \mathrm{PSG}(f, \Theta, n, \alpha)$$

❸ ...eturn... $; \hat{\theta})$

**Read the Paper**
**Today!**
*Now 50% Off!*

W.h.p., *estimation* + *optimization*
error don't exceed $\varepsilon$

Polynomial time + sample complexity

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization

Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

# Recap: Characterizing Fair PAC-Learnability

- Classical method: measure population sentiment with welfare
- This work: welfare and malfare on equal axiomatic footing
  - *Malfare minimization* is *fair extension* of *risk minimization*

- Classical method: measure population sentiment with welfare
- This work: welfare and malfare on equal axiomatic footing
  - *Malfare minimization* is *fair extension* of *risk minimization*
- Under some conditions, $\mathrm{PAC} = \mathrm{FPAC}$ (statistical equivalence)
  - $\mathrm{FPAC} \implies \mathrm{PAC}$ (as a special case)
  - Constructive $\mathrm{PAC} \implies \mathrm{FPAC}$ reduction in realizable case
  - General case is non-constructive, assumes no-free-lunch argument

- Classical method: measure population sentiment with welfare
- This work: welfare and malfare on equal axiomatic footing
  - *Malfare minimization* is *fair extension* of *risk minimization*
- Under some conditions, $\mathrm{PAC} = \mathrm{FPAC}$ (statistical equivalence)
  - $\mathrm{FPAC} \implies \mathrm{PAC}$ (as a special case)
  - Constructive $\mathrm{PAC} \implies \mathrm{FPAC}$ reduction in realizable case
  - General case is non-constructive, assumes no-free-lunch argument
- Open research question: does efficient PAC $\implies$ efficient FPAC?
  - Constructive reduction in realizable case
  - Efficient cover enumerability sufficient for both
  - Standard convex optimization assumptions sufficient for both

- Classical method: measure population sentiment with welfare
- This work: welfare and malfare on equal axiomatic footing
  - *Malfare minimization* is *fair extension* of *risk minimization*
- Under some conditions, $\mathrm{PAC} = \mathrm{FPAC}$ (statistical equivalence)
  - $\mathrm{FPAC} \implies \mathrm{PAC}$ (as a special case)
  - Constructive $\mathrm{PAC} \implies \mathrm{FPAC}$ reduction in realizable case
  - General case is non-constructive, assumes no-free-lunch argument
- Open research question: does efficient PAC $\implies$ efficient FPAC?
  - Constructive reduction in realizable case
  - Efficient cover enumerability sufficient for both
  - Standard convex optimization assumptions sufficient for both
  *Conjecture*: No, $\exists$ PAC-learnable class, where FPAC-learning is $\mathrm{NP}$-hard (and $\mathrm{P} \neq \mathrm{NP}$)

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
  Welfare
  Malfare
  Axiomatic
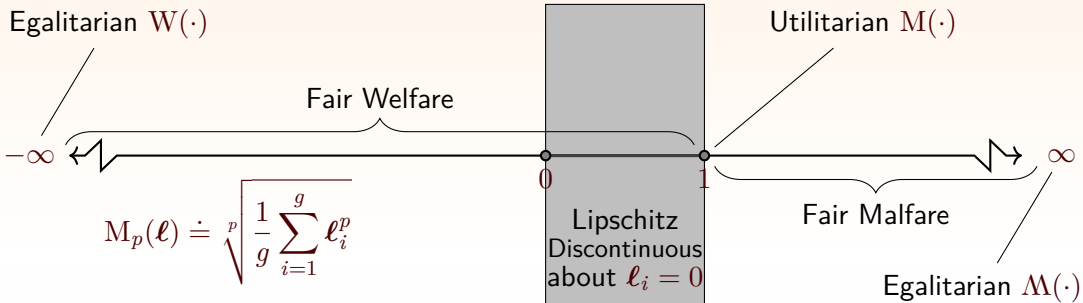  Characterization

Estimation
and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC
Learning
  Computational
  Learnability

In Conclusion

# Recap: Malfare, Welfare, and FPAC Learning



Egalitarian $\mathrm{W}(\cdot)$

Utilitarian $\mathrm{M}(\cdot)$

Fair Welfare

$-\infty$      0      1      $\infty$

Fair Malfare

$$\mathrm{M}_p(\boldsymbol{\ell}) \doteq \sqrt[p]{\frac{1}{g} \sum_{i=1}^{g} \boldsymbol{\ell}_i^p}$$

Lipschitz
Discontinuous
about $\boldsymbol{\ell}_i = 0$

Egalitarian $\mathbb{M}(\cdot)$

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
Welfare
Malfare
Axiomatic
Characterization
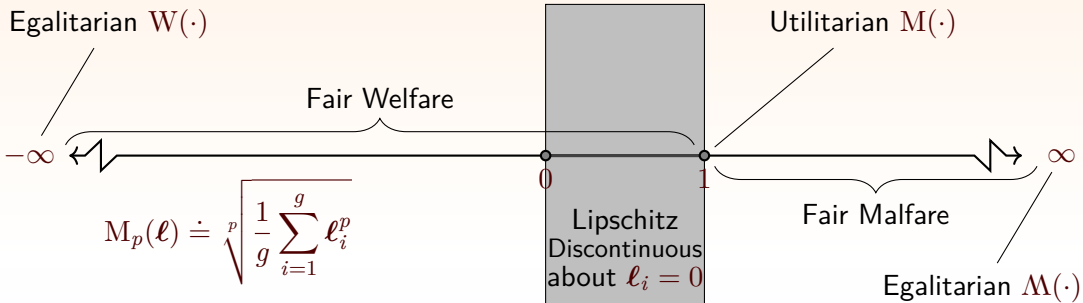
Estimation
and Inference
Linear Classifiers
Statistical Estimation

Fair PAC
Learning
Computational
Learnability

In Conclusion

Recap: Malfare, Welfare, and FPAC Learning

$$\mathrm{M}_p(\boldsymbol{\ell}) \doteq \sqrt[p]{\frac{1}{g} \sum_{i=1}^{g} \boldsymbol{\ell}_i^p}$$

Egalitarian $\mathrm{W}(\cdot)$

Utilitarian $\mathrm{M}(\cdot)$

Fair Welfare

Fair Malfare

$-\infty$  0  1  $\infty$

Lipschitz
Discontinuous
about $\boldsymbol{\ell}_i = 0$

Egalitarian $\mathbb{M}(\cdot)$

Why use malfare instead of welfare?

❶ "Most" machine learning tasks more naturally cast as *loss minimization*
  • Exceptions: reward, profit, accuracy maximization

Axiomatic Fair
Learning

Cyrus Cousins

NeurIPS2021

Philosophy,
Welfare, and
Malfare
  Welfare
  Malfare
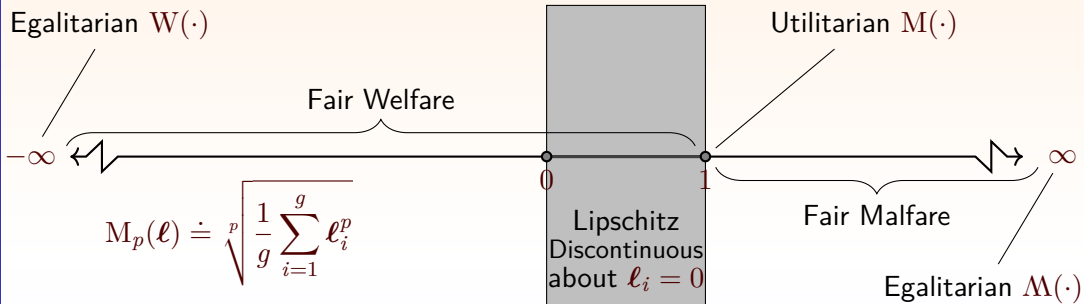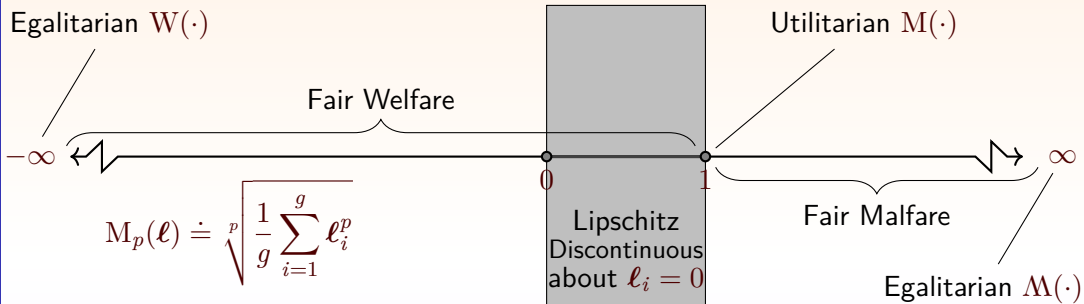  Axiomatic
  Characterization

Estimation
and Inference
  Linear Classifiers
  Statistical Estimation

Fair PAC
Learning
  Computational
  Learnability

In Conclusion

Recap: Malfare, Welfare, and FPAC Learning

Egalitarian $\mathrm{W}(\cdot)$

Utilitarian $\mathrm{M}(\cdot)$

Fair Welfare

$-\infty$ — 0 — 1 — $\infty$

Fair Malfare

$$\mathrm{M}_p(\boldsymbol{\ell}) \doteq \sqrt[p]{\frac{1}{g}\sum_{i=1}^{g}\boldsymbol{\ell}_i^p}$$

Lipschitz
Discontinuous
about $\boldsymbol{\ell}_i = 0$

Egalitarian $\mathbb{M}(\cdot)$

Why use malfare instead of welfare?

❶ "Most" machine learning tasks more naturally cast as *loss minimization*

- Exceptions: reward, profit, accuracy maximization

❷ Fair PAC-Learning *with welfare targets* is tricky

- Inherent statistical instability for $p \in [0, 1)$
- Require additional assumptions, or restricted capabilities

Egalitarian $\mathrm{W}(\cdot)$      Utilitarian $\mathrm{M}(\cdot)$

Fair Welfare

$-\infty$    $0$    $1$    $\infty$

Fair Malfare

$$\mathrm{M}_p(\boldsymbol{\ell}) \doteq \sqrt[p]{\frac{1}{g} \sum_{i=1}^{g} \boldsymbol{\ell}_i^p}$$

Lipschitz
Discontinuous
about $\boldsymbol{\ell}_i = 0$

Egalitarian $\mathbb{M}(\cdot)$

Why use malfare instead of welfare?

❶ "Most" machine learning tasks more naturally cast as *loss minimization*
- Exceptions: reward, profit, accuracy maximization

❷ Fair PAC-Learning *with welfare targets* is tricky
- Inherent statistical instability for $p \in [0, 1)$
- Require additional assumptions, or restricted capabilities

❸ **Why not?**