# Fair Inference Through Semiparametric-Efficient Estimation Over Constraint-Specific Paths

## Nima Hejazi

Group in Biostatistics, and
Center for Computational Biology,
University of California, Berkeley

nimahejazi.org
Twitter: nshejazi
GitHub: nhejazi

slides: bit.ly/jsm_fairtmle_2018

This slide deck is for a short presentation on new work in Targeted Learning, discussing both the construction of constrained ensemble machine learning for the construction of "fair" estimates and a novel algorithm for computing TML estimators with respect to some constraint functional.

Slides: bit.ly/jsm_fairtmle_2018

# Preview: Summary

▶ Recent work suggests that the widespread use of machine learning algorithms has had negative social and policy consequences.

▶ The widespread use of machine learning in policy issues violates human intuitions of bias.

▶ We propose a general algorithm for constructing "fair" optimal ensemble ML estimators via cross-validation.

▶ Constraints may be imposed as functionals defined over the target parameter of interest.

▶ Estimating constrained parameters may be seen as iteratively minimizing a loss function along a *constrained* path in the parameter space $\Psi$.

1

We'll go over this summary again at the end of the talk. Hopefully, it will all make more sense then.

Bernard Parker, left, was rated high risk; Dylan Fugett was rated low risk. (Josh Ritchie for ProPublica)

# Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

# Fairness is machine learning?

Another potential result: a more diverse
workplace. The software relies on data
to surface candidates from a wide variety
of places...free of human biases. But
software is not free of human influence.
Algorithms are written and maintained by
people...As a result...algorithms can
reinforce human prejudices.

-Miller (2015)

Obviously, it's important to explain the motivating example here.

# Addressing bias in a technical manner

- ▶ The careless use of machine learning may induce *unjustified* bias.

- ▶ Problematic discrimination by ML approaches leads to solutions with *practical irrelevance*.

- ▶ Ill-considered discrimination by ML approaches leads to solutions that are *morally problematic*.

- ▶ Two doctrines of discrimination:
  1. Disparate treatment: formal or intentional
  2. Disparate impact: unjustified or avoidable

4

Considering and treating bias using a technical approach is an important way of dealing with the potential negative consequences of machine learning.

# Background, data, notation

- ► An observational unit: $O = (W, X, Y)$, where $W$ is baseline covariates, $X$ a sensitive characteristic, $Y$ an outcome of interest.

- ► Consider $n$ i.i.d. copies $O_1, \ldots, O_n$ of $O \sim P_0 \in \mathcal{M}$.

- ► Here, $\mathcal{M}$ is an infinite-dimensional statistical model (i.e., indexed by an infinite-dimensional vector).

- ► We discuss the estimation of a target parameter $\psi : \mathcal{M} \to \mathbb{R}$, where

$$\Psi(P_0) = \arg\min_{\psi \in \Psi} \mathbb{E}_{P_0} L(\psi)$$

We just need to see this to get a feel for what's going to be happening with the derivation of constraint-specific paths.

# Just a few fairness criteria

▶ Let $C : (X, W) \to Y \in \{0, 1\}$ be a classifier; $X \in \{a, b\}$.

▶ Demographic parity: $\mathbb{P}_{(X=a)}(C = 1) = \mathbb{P}_{(X=b)}(C = 1)$

▶ Accuracy parity: $\mathbb{P}_{(X=a)}(C = Y) = \mathbb{P}_{(X=b)}(C = Y)$

▶ True positive parity:
$\mathbb{P}_{(X=a)}(C = 1 \mid Y = 1) = \mathbb{P}_{(X=b)}(C = 1 \mid Y = 1)$

▶ False positive parity:
$\mathbb{P}_{(X=a)}(C = 1 \mid Y = 0) = \mathbb{P}_{(X=b)}(C = 1 \mid Y = 0)$

▶ Positive predictive value parity:
$\mathbb{P}_{(X=a)}(Y = 1 \mid C = 1) = \mathbb{P}_{(X=b)}(Y = 1 \mid C = 1)$

▶ Negative predictive value parity:
$\mathbb{P}_{(X=a)}(Y = 1 \mid C = 0) = \mathbb{P}_{(X=b)}(Y = 1 \mid C = 0)$

It's a jungle out there

# Wait, where did the fairness go?

▶ Goal: estimate $\Psi(P_0) = \mathbb{E}_{P_0}(Y \mid X, W)$.

▶ Let $Y \in \{0, 1\}$ and use negative log-likelihood loss:

$$L(\psi) = -(Y\log(\mathbb{P}(Y \mid X, W)) + (1-Y)\log(1-\mathbb{P}(Y \mid X, W)))$$

▶ Fairness criterion — *equalized odds*:

$$\Theta_\psi(P_0) = \sum_y \{\mathbb{E}_{P_0}(L(\psi)(O) \mid X = 1, Y = y)$$
$$-\mathbb{E}_{P_0}(L(\psi)(O) \mid X = 0, Y = y)\}^2$$

▶ Let $\Theta_\psi(P_0) : \mathcal{M} \to \mathbb{R}$ be a pathwise differentiable *functional* for each $\psi \in \Psi$.

Equalized odds simultaneously enforces both true positive parity and false positive parity.

# Constrained functional parameters

▶ Estimate target parameter under a constraint:

$$\Psi(P_0) = \underset{\psi \in \Psi, \Theta_\psi(P_0) = 0}{\arg\min} \mathbb{E}_{P_0} L(\psi)$$

▶ Goal: estimate $\Psi^*(P_0)$, the projection of $\Psi(P_0)$ onto the subspace $\Psi^*(P_0) = \{\psi \in \Psi : \Theta_\psi(P_0) = 0\}$:

$$(\Psi^*, \lambda) = (\Psi^*(P_0), \Lambda(P_0)) \equiv \underset{\psi \in \Psi, \lambda}{\arg\min} \mathbb{E}_{P_0} L(\psi) + \lambda \Theta_\psi(P_0).$$

▶ *Lemma*: If $\widetilde{\Psi}(P_0) = (\Psi^*(P_0), \Lambda(P_0))$ is the minimizer of the Lagrange multiplier penalized loss, then

$$\Psi^*(P_0) = \underset{\psi \in \Psi, \Theta_\psi(P_0) = 0}{\arg\min} \mathbb{E}_{P_0} L(\psi).$$

# Learning with constrained parameters

- ▶ Risk function: $R(\widetilde{\psi} \mid P) \equiv P_n L(\psi^*) + \lambda \Theta(\psi^* \mid P)$, where $\widetilde{\psi} = (\psi^*, \lambda)$

- ▶ For $\widetilde{\psi}(P_n) = (\hat{\Psi}^*(P_n), \hat{\lambda}(P_n))$ of $\widetilde{\Psi}(P_0)$, and sample splitting scheme $B_n \in \{0, 1\}^n$:

$$R_0(\widetilde{\psi}, P_n) = \mathbb{E}_{B_n} P_0 L(\hat{\Psi}^*(P^0_{n,B_n})) + \hat{\Lambda}(P_n) \mathbb{E}_{B_n} \Theta(\hat{\Psi}^*(P^0_{n,B_n}) \mid P_0)$$

- • Here $P^0_{n,B_n}$ denotes the empirical distribution of the training sample.

# Learning with constrained parameters

▶ Cross-validated risk:

$$R_{n,CV}(\tilde{\psi}, P_n) = E_{B_n} P^1_{n,B_n} L(\hat{\Psi}^*(P^0_{n,B_n})) \qquad (1)$$
$$+ \hat{\Lambda}(P_n) E_{B_n} \Theta(\hat{\Psi}^*(P^0_{n,B_n}) \mid P^*_{n,B_n}) \qquad (2)$$

▶ Given candidate estimators $\tilde{\psi}_j(P_n) = (\hat{\Psi}^*_j(P_n), \hat{\Lambda}_j(P_n))$, $j = 1, \ldots, J$, the CV selector is given by:
$J_n = \arg\min_j R_{n,CV}(\tilde{\psi}_j, P_n)$.

▶ We may define an optimal estimate of $\tilde{\Psi}$ by
$\tilde{\psi}_n \equiv \tilde{\psi}_{J_n}(P_n) = (\hat{\Psi}_{J_n}(P_n), \hat{\lambda}_{J_n}(P_n))$

# Mappings with constrained learners

A straightforward approach to generating estimators of the constrained parameter would be to simply generate a mapping according to the following simple process:

1. Generate an unconstrained estimate $\psi_n$ of the unconstrained parameter $\psi_0$,

2. Map an estimator $\Theta_{\psi_n,n}$ of the constraint $\Theta_{\psi_n}(P_0)$ into the path $\psi_{n,\lambda}$. The corresponding solution $\psi_n^* = \psi_{n,\lambda_n}$ of $\Theta_{\psi_{n,\lambda_n},n} = 0$ generates an estimator of the constrained parameter.

# Constraint-specific paths

► Consider $\psi_{0,\lambda} = \arg\max_{\psi \in \Psi} \mathbb{E}_{P_0} L(\psi) + \lambda \Theta_0(\psi)$.

► $\{\psi_{0,\lambda} : \lambda\}$ represents a path in the parameter space $\Psi$ through $\psi_0$ at $\lambda = 0$.

► This is a *constraint-specific path*, as it produces an estimate under the desired functional constraint.

► Leverage this construction to map an initial estimator of the unconstrained parameter $\psi_0$ into its corresponding constrained version $\psi_0^*$.

# Future work

- ► Further generalization of constraint-specific paths: the solution path $\{\psi_{0,\lambda} : \lambda\}$ in the parameter space $\Psi$ through $\psi_0$ at $\lambda = 0$.

- ► Further develop relation between constraint-specific paths and universal least favorable submodels.

- ► Integration of the approach of constraint-specific paths with classical classical targeted maximum likelihood estimation — in particular, what, if any, are the implications for inference?

13

# Review: Summary

▶ Recent work suggests that the widespread use of machine learning algorithms has had negative social and policy consequences.

▶ The widespread use of machine learning in policy issues violates human intuitions of bias.

▶ We propose a general algorithm for constructing "fair" optimal ensemble ML estimators via cross-validation.

▶ Constraints may be imposed as functionals defined over the target parameter of interest.

▶ Estimating constrained parameters may be seen as iteratively minimizing a loss function along a *constrained* path in the parameter space $\Psi$.

It's always good to include a summary.

# References I

# Acknowledgments

Thank you.

Slides: bit.ly/jsm_fairtmle_2018

```
nimahejazi.org
```

```
Twitter: nshejazi
```

```
GitHub: nhejazi
```

Here's where you can find me, as well as the slides for this talk.