Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Novel Regularization Approach to Fair Machine Learning

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Bay Area R Users Group GRAIL July 18, 2023

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Plan of this talk:

Here we will introduce a new method for fair machine learning. Package code is available in https://github.com/matloff/EDFfair

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Overview

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Fairness in ML:

Overview

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

A Novel Regularization Approach to Fair Machine Learning

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Fairness in ML:

• The usual ML: predict Y from vector X.

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

A Novel Regularization Approach to Fair Machine Learning

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

- The usual ML: predict Y from vector X.
- But X includes a sensitive variable S (race, gender, age etc.)

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

A Novel Regularization Approach to Fair Machine Learning

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

- The usual ML: predict Y from vector X.
- But X includes a sensitive variable S (race, gender, age etc.)
- Wish to exclude S or at least minimize its impact.

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

ularization Approach to Fair Machine Learning Norman Matloff University of California, Davis Wenxi Zhang

> Columbia University.

A Novel Reg-

- The usual ML: predict Y from vector X.
- But X includes a sensitive variable S (race, gender, age etc.)
- Wish to exclude S or at least minimize its impact.
- But there may be covariates C in X that are proxies for S,

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Approach to Fair Machine Learning Norman Matloff University of California, Davis Wenxi Zhang Columbia

University.

A Novel Regularization

- The usual ML: predict Y from vector X.
- But X includes a sensitive variable S (race, gender, age etc.)
- Wish to exclude S or at least minimize its impact.
- But there may be covariates C in X that are proxies for S, so that you end up "including" S anyway.

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Fair Machine Learning Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Novel Regularization

Approach to

- The usual ML: predict Y from vector X.
- But X includes a sensitive variable S (race, gender, age etc.)
- Wish to exclude S or at least minimize its impact.
- But there may be covariates C in X that are proxies for S, so that you end up "including" S anyway.
- Fairness-Utility Tradeoff: The greater the influence we allow for C, the greater our utility (pred. acc.), but the lesser our fairness.

Example

▲ロト ▲冊ト ▲ヨト ▲ヨト ヨー のくで

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

A Novel Regularization Approach to Fair Machine Learning

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

The "Hello World" of Fair ML: COMPAS algorithm

• Northpointe developed commercial product, COMPAS, to predict recidivism.

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Approach to Fair Machine Learning Norman Matloff University of California, Davis Wenxi Zhang Columbia

University.

A Novel Regularization

- Northpointe developed commercial product, COMPAS, to predict recidivism.
- Used by judges to aid in determining sentence for convicted criminals.

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Novel Regularization

Approach to Fair Machine Learning

- Northpointe developed commercial product, COMPAS, to predict recidivism.
- Used by judges to aid in determining sentence for convicted criminals.
- *ProPublica* expose' claimed COMPAS biased against Black defendants.

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Learning Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Novel Regularization

Approach to Fair Machine

- Northpointe developed commercial product, COMPAS, to predict recidivism.
- Used by judges to aid in determining sentence for convicted criminals.
- *ProPublica* expose' claimed COMPAS biased against Black defendants.
- S = race, C includes # of priors, educ. level etc.

Criteria for "Fairness"

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Many criteria have been proposed.

Criteria for "Fairness"

▲ロト ▲冊ト ▲ヨト ▲ヨト ヨー のくで

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Many criteria have been proposed. (This is a research area, after all. :-))

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Criteria for "Fairness"

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Many criteria have been proposed. (This is a research area, after all. :-)) Why does the criterion one uses matter?

- Criterion itself may be biased. (Northpointe claimed this about *ProPublica*.)
- Many ML unfairness remedies are based on exactly satisfying some chosen fairness citerion.

Example Criteria

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Example Criteria

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Say S is categorical (e.g. race, gender). Set \widehat{Y} = predicted value or class.

Some common criteria:

Example Criteria

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Say S is categorical (e.g. race, gender). Set \widehat{Y} = predicted value or class.

Some common criteria:

• Demographic Parity

Example Criteria

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Say S is categorical (e.g. race, gender). Set \widehat{Y} = predicted value or class.

Some common criteria:

- Demographic Parity
 - \widehat{Y} , *S* independent

Example Criteria

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Say S is categorical (e.g. race, gender). Set $\widehat{Y}=$ predicted value or class.

Some common criteria:

- Demographic Parity
 - \widehat{Y}, S independent
- Equalized Odds

Example Criteria

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Say S is categorical (e.g. race, gender). Set $\widehat{Y}=$ predicted value or class.

Some common criteria:

- Demographic Parity \widehat{Y} . *S* independent
- Equalized Odds

 \$\hat{Y}\$ independent of S, given Y
 Retrospective, e.g. among those who end up *not* recidivating, \$\hat{Y}\$ should not have been affected by S.

There are various others that are popular in the research realm. These can also be phrased in terms of FPR, TPR etc.

Relaxing a Criterion

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Relaxing a Criterion

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

• Allow the criterion to be only approximately met.

Relaxing a Criterion

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

- Allow the criterion to be only approximately met.
- Set a "slider" with which the user can a select a point in the Fairness-Utility spectrum.

Relaxing a Criterion

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

- Allow the criterion to be only approximately met.
- Set a "slider" with which the user can a select a point in the Fairness-Utility spectrum.
- For continuous Y, use correlation.

Relaxing a Criterion

- Norman Matloff University of California, Davis Wenxi Zhang Columbia University.
- Allow the criterion to be only approximately met.
- Set a "slider" with which the user can a select a point in the Fairness-Utility spectrum.
- For continuous Y, use correlation.
- For binary case, we use R(T, W). T and W are \widehat{Y} and S if Y is continuous, $\widehat{P}(Y = 1|X)$ if Y is binary, similarly if S is binary.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Work by Komiyama *et al* and Scutari *et al*

- (Many others, just two examples.)
- Regress Y on S, then regress the residuals U on $\tilde{X} = X$ without the X component.
- User sets an upper bound on $R^2(\widehat{Y}, S)$ to set the level of Fairness-Utility Tradeoff.
- Komiyama use quadratic programming optimization, thus iterative.
- Scutari approach the problem via ridge regression (with λ for the regression on U).
- Scutari is implemented in their fairml package on CRAN.

Scutari Approach

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

• Scutari *et al* (2022), "Achieving Fairness with a Simple Ridge Penalty"

Scutari Approach

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Novel Regularization

Approach to Fair Machine Learning

- Scutari *et al* (2022), "Achieving Fairness with a Simple Ridge Penalty"
- Lots of technical detail. Summary: first linear regress Y on S, then linear ridge-regress residuals on (non-S part of) X.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Our Approach (Linear Case)

• Again use ridge, but differently.

 $\operatorname{argmin}_{b} ||\mathbb{Y} - \mathbb{X}b||^{2} + ||Db||^{2}$ (1)

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Our Approach (Linear Case)

• Again use ridge, but differently.

$$\operatorname{argmin}_{b} ||\mathbb{Y} - \mathbb{X}b||^{2} + ||Db||^{2}$$
(1)

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

 D = diag(d₁,..., d_p) is a hyperparameter, set to desirable point in Fairness-Utility Tradeoff.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Our Approach (Linear Case)

• Again use ridge, but differently.

$$\operatorname{argmin}_{b} ||\mathbb{Y} - \mathbb{X}b||^{2} + ||Db||^{2}$$
(1)

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

- D = diag(d₁,..., d_p) is a hyperparameter, set to desirable point in Fairness-Utility Tradeoff.
- Presumably $d_{i_s} = 0$.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Our Approach (Linear Case)

• Again use ridge, but differently.

$$\operatorname{argmin}_{b} ||\mathbb{Y} - \mathbb{X}b||^{2} + ||Db||^{2}$$
(1)

▲ロ ▶ ▲周 ▶ ▲ ヨ ▶ ▲ ヨ ▶ ● の < ○

- D = diag(d₁,..., d_p) is a hyperparameter, set to desirable point in Fairness-Utility Tradeoff.
- Presumably $d_{i_s} = 0$.
- The positive d_j are for the proxies, i.e. C.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Our Approach (Linear Case)

• Again use ridge, but differently.

$$\operatorname{argmin}_{b} ||\mathbb{Y} - \mathbb{X}b||^{2} + ||Db||^{2}$$
(1)

▲ロ ▶ ▲周 ▶ ▲ ヨ ▶ ▲ ヨ ▶ ● の < ○

- D = diag(d₁,..., d_p) is a hyperparameter, set to desirable point in Fairness-Utility Tradeoff.
- Presumably $d_{i_s} = 0$.
- The positive d_j are for the proxies, i.e. C.
- Advtantage over Scutari *et al*: Not "One size fits all," different *d_i* for different *S_i*.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Our Approach (Linear Case)

• Again use ridge, but differently.

$$\operatorname{argmin}_{b} ||\mathbb{Y} - \mathbb{X}b||^{2} + ||Db||^{2}$$
(1)

▲ロ ▶ ▲周 ▶ ▲ ヨ ▶ ▲ ヨ ▶ ● の < ○

- D = diag(d₁,..., d_p) is a hyperparameter, set to desirable point in Fairness-Utility Tradeoff.
- Presumably $d_{i_s} = 0$.
- The positive d_j are for the proxies, i.e. C.
- Advtantage over Scutari *et al*: Not "One size fits all," different *d_i* for different *S_i*.
- Closed-form solution for b:

$$b = \left[\mathbb{X}'\mathbb{X} + D^2\right]^{-1}\mathbb{X}'\mathbb{Y}$$

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Computational Trick

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

A well-known trick for ridge-regression generalizes to our setting.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Computational Trick

A well-known trick for ridge-regression generalizes to our setting.

Set

$$A = \begin{pmatrix} \mathbb{X} \\ D \end{pmatrix}$$
(2)
$$B = \begin{pmatrix} \mathbb{Y} \\ 0 \end{pmatrix}$$
(3)

▲ロト ▲冊ト ▲ヨト ▲ヨト ヨー のくで

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Computational Trick

A well-known trick for ridge-regression generalizes to our setting.

Set

$$A = \begin{pmatrix} \mathbb{X} \\ D \end{pmatrix}$$
(2)
$$B = \begin{pmatrix} \mathbb{Y} \\ 0 \end{pmatrix}$$
(3)

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

 Then run Im() as usual, using A and B as the design matrix and response variable data, instead of X and Y.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Computational Trick

A well-known trick for ridge-regression generalizes to our setting.

Set

$$A = \begin{pmatrix} \mathbb{X} \\ D \end{pmatrix}$$
(2)
$$B = \begin{pmatrix} \mathbb{Y} \\ 0 \end{pmatrix}$$
(3)

- Then run Im() as usual, using A and B as the design matrix and response variable data, instead of X and Y.
- This gives us our desired ridge estimator.

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Novel Regularization

Approach to Fair Machine Learning

> Package in https://github.com/matloff, small issues. > data(compas) > z <- qeFairRidgeLog(compas, 'two year recid', **list** (decile score = 0.8, gender = 0.8, > priors count = 0.8, age = 0.8), > 'race', yesYVal='Yes', holdout=NULL) > # try a prediction, like row 1 but age 33 not 6 > newx <- compas[1, -9] > newx['age'] <- 33 > predict(z, newx) 1 0.2854387 *# 28.5%* chance to recidivate

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Extension to Other ML Algorithms

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Extension to Other ML Algorithms

• Random forests: Set node-split probability lower for features in C than in the rest of X. (The **ranger** package could be used.)

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Extension to Other ML Algorithms

• Random forests: Set node-split probability lower for features in C than in the rest of X. (The **ranger** package could be used.)

▲ロ ▶ ▲周 ▶ ▲ ヨ ▶ ▲ ヨ ▶ ● の < ○

 k-nearest neighbors (k-NN): In defining the distance metric, place smaller weight on the coordinates corresponding to C. (Could use **qeKNN** in my forthcoming **qeML** package.)

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

Extension to Other ML Algorithms

- Random forests: Set node-split probability lower for features in C than in the rest of X. (The **ranger** package could be used.)
- k-nearest neighbors (k-NN): In defining the distance metric, place smaller weight on the coordinates corresponding to C. (Could use **qeKNN** in my forthcoming **qeML** package.)
- Support vector machines: Apply an ℓ₂ constraint on the portion of the vector w of hyperplane coefficients corresponding to C.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Related Package

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Related Package

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

- dsld, Data Science Looks at Discrimination
- Race, gender, age etc.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Related Package

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

- dsld, Data Science Looks at Discrimination
- Race, gender, age etc.
- "Statistical discrimination analysis in a box"

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Related Package

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

- dsld, Data Science Looks at Discrimination
- Race, gender, age etc.
- "Statistical discrimination analysis in a box"
- Will include paired Quarto book teaching the stat concepts needed for investigation of discrimination.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Related Package

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

- dsld, Data Science Looks at Discrimination
- Race, gender, age etc.
- "Statistical discrimination analysis in a box"
- Will include paired Quarto book teaching the stat concepts needed for investigation of discrimination.
- Part I: General discrimination analysis—effect of S. Part II: Fair ML—predict while avoiding use of S.

Norman Matloff University of California, Davis Wenxi Zhang Columbia University.

A Related Package

- dsld, Data Science Looks at Discrimination
- Race, gender, age etc.
- "Statistical discrimination analysis in a box"
- Will include paired Quarto book teaching the stat concepts needed for investigation of discrimination.
- Part I: General discrimination analysis—effect of S. Part II: Fair ML—predict while avoiding use of S.
- Anticipated usage includes:
 - Teaching and research in the social sciences/economics.
 - Litigation support.
 - Government agencies.
 - Corporate HR analysis.
 - Consulting.