

# Fairness in Criminal Justice Risk Assessment

Goal: Define different kinds of (statistical) fairness at the Tradeoffs

## Terminology

### Confusion Table:

	Failure Predicted (positive)	Success Predicted (negative)	Conditional Procedure Error
Failure - A positive	$tp$ True Positive	$fn$ False Negative	$\frac{fn}{tp+fn}$ = False negative rate
Success - A negative	$fp$ False Positive	$tn$ True Negative	$\frac{fp}{tn+fp}$ = False positive rate
Condition Use Error	$\frac{fp}{tp+fp}$ = Failure Predicted Error	$\frac{fn}{tn+fn}$ = Success Predicted Error	$\frac{fp+fn}{fp+fn+tn+tp}$ = Overall procedure error

Margins of the table → estimates of corresponding probabilities in a population

Sample size: total number of observations  $tp + fn + fp + tn = N$

Base rate: proportion of actual failures:  $\frac{fp+fn}{N}$

or      "      successes:  $\frac{tn+fp}{N}$

Prediction: proportion predicted to fail:  $\frac{tp+fp}{N}$

Distribution:      "      success:  $\frac{tn+fn}{N}$

Overall procedure error: proportion of cases misclassified

Conditional procedure error: proportion of cases ~~incorrectly classified~~ conditional on ~~actual~~ outcome

"      use      "      "      "      predicted      "      "      predicted

Cost ratio: ratio of false negatives to false positives:  $\frac{fn}{fp}$

or      "      "      "      "      "      "      positives      "      negatives:  $\frac{fp}{fn}$

→ Example parole release

## Statistical Framework

Goal: draw inference beyond data being analyzed

→ Idea: each observation is randomly realized from a single joint probability distribution  $P(Y, L, S)$  = "target population"

- $Y$  = outcome of interest that is supposed to be inferred, e.g. arrest while on parole
- $L$  = legitimate predictors, e.g. prior convictions
- $S$  = protected predictors, e.g. race, gender, ...

⇒ population = limitless number of IID observations on  $P(Y, L, S)$

In the population there is some function  $f(L, S)$  linking predictors  $L, S$  to the expectation of  $Y$ :  $E(Y|L, S)$ . If  $Y$  is categorical: cat. expectations = cat. probabilities

A fitting procedure  $h(L, S)$  is applied to the data

- Structure: e.g. linear regression
- Optimization algorithm: minimizing sum of least squares

⇒ ~~fitting~~ algorithm fitting  $h(L, S)$  produces hypothesis  $\hat{f}(L, S)$  which is the source of fitted values  $\hat{Y}$

NOTE: allow  $S$  to participate in fitting procedure because it is associated with  $Y$

→  $\hat{f}(L, S)$  will almost certainly be a biased estimator of the true response surface

- important legitimate predictors not available
- measurement error
- functional form arrived at for  $\hat{f}(L, S)$  might be incorrect

Joint probability distribution is essentially abstraction of a high dimensional histogram for a finite population → expanded to limitless observations

Assumption: Data is realized independently from the same joint prob. dist =  $\approx$  IID

→ whether this assumption makes sense for real data depends

## Defining algorithmic fairness

Proceed as if  $\hat{f}(L, S)$  provides an estimate that are the same as the corresponding population values  $\Rightarrow$  do not conflict discussion about accuracy fairness

1) Overall accuracy equality:  $\frac{tp + tn}{tn + fp + tn + tp}$  is the same across all protected groups

$\rightarrow tp$  and  $tn$  are equally desirable  $\rightarrow$  not true for many cost-weighted settings

Not really used because there is no distinction between successful failure accuracy

2) Statistical parity:  $\frac{tp + fp}{N}$  or  $\frac{tn + fn}{N}$  (marginal dist. of predictions) is the same across all protected groups

$\hookrightarrow$  different from one another, but the same for all protected groups

$\rightarrow$  „demographic parity“: can lead to undesirable results, e.g. ~~higher~~ incarceration for women that are no public safety risk to maintain proportions

3) Conditional procedure accuracy equality:  $\frac{tp}{fp + fn}$  or  $\frac{tn}{tn + fn}$  is the same across all protected groups

$\Rightarrow$  „equalized odds“: the accuracy of a classification, given the true outcome class  
e.g. accuracy of being classified as being arrested while on parole, given the person was actually arrested on parole

4) Conditional use accuracy equality:  $\frac{tp}{tp + fp}$  or  $\frac{tn}{tn + fn}$  is the same across all protected groups

$\Rightarrow$  the accuracy of being part of a particular class, given the predicted classification  
e.g. accuracy of being arrested on parole, given the algorithm predicted no arrest on parole

5) Treatment equality:  $\frac{f_n}{f_p}$  or  $\frac{f_p}{f_n}$  is the same across all protected groups

→ treatment "implies": policy lever to achieve other kinds of fairness

e.g. if  $f_n$  are made to be more costly for men  $\Rightarrow$  conditional procedure accuracy equality

BUT: men and women are not treated equally, us incorrectly classifying a female on parole is more costly for men

6) Total fairness: all previously defined notions of fairness are achieved simultaneously

Other notions of fairness are not discussed, because they cannot be operationalized.

## Estimator Accuracy

Previously: applying  $h(L, S)$  to data  $\Rightarrow \hat{f}(L, S)$  that estimates true response surface

But: no credible claim, that response surface is being estimated in unbiased manner

Still: With more samples  $\rightarrow$  smaller estimator error

$$\boxed{\text{bias.} + \text{variance}} = \text{no conventional confidence limit}$$

$\Rightarrow$  bias cannot be removed w/o comparison to truth, which is unknown

Alternative:  $\hat{f}(L, S)$  estimates a response surface  $\hat{f}_p$ ,  $\hat{f}_p$  is an approximation of the true

in the population  $f(L, S)$  has the same form as the approximation  $\rightarrow$  probability estimator

of cross-tabulation ~~on~~ data = prob. est. of a  $Y$ -by- $Y$  cross-tabulation of applying  $h(L, S)$

$\Rightarrow$  since data is assumed to be IID  $\Rightarrow$  estimates are asymptotically unbiased

Estimator accuracy is then measured by out-of-sample performance

## Tradeoffs

Fundamentally:  $h(L, S)$  capitalizes on non-redundant associations that  $L$  and  $S$  have with  $Y$   
⇒ excluding or discarding  $S$  will reduce accuracy

But fairness is important as well ⇒ examples of tradeoffs

## Proven Impossibility Theorem 1:

„When the base rates differ by protected groups and when there is not separation, one cannot have both conditional use accuracy and equality in false negative/positive rate. (f1 = f2)“  
(reverse cond. procedure acc.)

## Definitions:

Separation: observations are separable if perfectly accurate classification is possible  
⇒ for each possible configuration of predictor values there is some  $h(L, S)$  for which probability of membership in a given outcome class is always either 1.0 or 0.0

Calibration: proportion of people experiencing an outcome (base rate)  
= proportion of people predicted to experience an outcome

• indicator for algorithm performance

• enables of fairness if difference in quality of calibration across groups

⇒ If there is variation in base rates and no separation  
→ goal of complete racial gender neutrality is unachievable  
⇒ ~~Tradeoffs~~

⇒ example parole release

# Potential Solutions

1) Preprocessing: eliminate sources of unfairness in data before  $h(L, S)$

- Remove linear dependencies between  $L$  and  $S$

• regress in turn each predictor in  $L$  on predictors in  $S$

$\hookrightarrow$  work with residuals  $\Rightarrow$  residualized ~~predictor~~ transformations of predictors

Problem: interaction effects are not removed unless specifically included

• transform predictors, so that fair predictions can be obtained while "preserving as much information as possible"  $\Rightarrow$  Euclidean distances between original and transformed predictors

$\rightarrow$  residualize in turn using results from previous residualizations and indicators from protected classes

- Rebalancing base rates:

• direct rebalancing: apply weights for each group separately  $\rightarrow$  base rates same across categories

• randomly rebalance same response values

Problems: 1) accuracy loss

2) differing false positive/negative rates across groups

- Direct vs. Indirect discrimination: features of protected class vs. selected feature of prot. class

• perturb class membership:

• perturbing outcome table ( $\hookrightarrow$  also changes base rate!)

- randomly transform all predictors except for indicators of protected class membership

$\Rightarrow$  joint distribution less ~~dependent~~ on prot. class membership

Constraints: 1) joint distro. of transformed variables is very close to joint distro. of original pred.

2) no individual case is substantially disturbed

Problems: - effect on different kinds of fairness unclear

• accuracy price

2) In-Processing: make fairness adjustments as part of the process of computing  $h(L, S)$

- alter results of ~~inference~~ to more fair results, if ~~prob.~~ itself is very uncertain

→ reduction of cut-off-sample accuracy might be small

Problem: might have unacceptable consequence for FPs/FN rates

- add a penalty term to penalize unfairness

• for example: penalty for violations of conditional procedure acc. equality

Problems:

- loss function typically not convex

- undesirable implications for other kinds of fairness

3) Post-processing: adjust results after applying  $h(L, S)$  to make it more fair

cond. prob. acc. eq./

- "equivalent odds": probabilistic switching of assigned labels

- binary case:  $f(L, S)$  assigns  $\in \{0 \text{ or } 1\}$

- Construct "inner" ~~process~~ model: assign each label a prob  $(0, 1)$  to switchable

$\Rightarrow$  minimize the classification error s.t. the one of the two fairness constraints by changing the values of the probabilities.

Problems:

- implications for other kinds of fairness unclear

- conditional use acc. eq. can suffer

- lower accuracy ( $\Rightarrow$  target: overall class error as sum of worst groups)

- reassigned prob. larger when base rates are more desperate

Making fairness operational: consider question of benchmark:

of equity is achieved, with regard to what?

e.g. fairness in terms of prison sentence length between black and white

offenders does not make statement about length of sentences except for the very  
equally long

Algorithmic

Fair Use:

- Combination of technical & policy challenges

- Benchmark: ~~of~~ current practice → even small steps can lead to meaningful impact

But: used data reflects past practices → current decisions affect future training data

↳ algorithmic solutions/results need to be regularly updated

Technical questions: - only long?

- discarding of historic data ~~than~~ much?

- more weight for recent training data?

⇒ example of In-Process Fairness

Conclusion:

fairness is a subtle problem ⇒ tradeoffs inevitable

⇒ Core Problem: differing base rates across protected groups

Paths forward:

1) Statistical procedures can improve transparency, accuracy and fairness up to point

2) Tradeoffs need to be explicit and available as hyperparameters

Measures of fairness need to be formalized

3) Decisions about tradeoffs need to be made by stakeholders

→ in matters of law, values and the political process

4) There will be no quick solutions