

# CounterFair: Counterfactual Burden Optimization for Bias Detection and Actionability-oriented Fairness

Dataset	Binary / Categorical	Ordinal / Continuous
Adult 488842 ins. 14 feat.	Sex ( $\times$ , -)	EducationLevel ( $\checkmark$ , $\uparrow$ )
	NativeCountry ( $\times$ , -)	AgeGroup ( $\times$ , -)
	Race ( $\times$ , -)	EducationNumber ( $\checkmark$ , $\uparrow$ )
	WorkClass ( $\checkmark$ , $\uparrow$ )	CapitalGain ( $\checkmark$ , $\uparrow$ )
	MaritalStatus ( $\checkmark$ , $\uparrow$ )	CapitalLoss ( $\checkmark$ , $\uparrow$ )
Athlete 1000 ins. 6 feat.	Occupation ( $\checkmark$ , $\uparrow$ )	HoursPerWeek ( $\checkmark$ , $\uparrow$ )
	Relationship ( $\checkmark$ , $\uparrow$ )	
Compas 7214 ins. 52 feat.	Sex ( $\times$ , -)	Age ( $\times$ , -)
	ChargeDegree ( $\checkmark$ , $\uparrow$ )	SleepHours ( $\checkmark$ , $\uparrow$ )
Dutch 60420 ins. 12 feat.	Race ( $\times$ , -)	PriorsCount ( $\checkmark$ , $\uparrow$ )
	HouseholdPosition ( $\checkmark$ , $\uparrow$ )	AgeGroup ( $\checkmark$ , $\uparrow$ )
German 1000 ins. 20 feat.	HouseholdSize ( $\checkmark$ , $\uparrow$ )	
	Country ( $\times$ , -)	EducationLevel ( $\checkmark$ , $\uparrow$ )
Student 395 ins. 33 feat.	EconomicStatus ( $\checkmark$ , $\uparrow$ )	Age ( $\checkmark$ , $\uparrow$ )
	CurEcoActivity ( $\checkmark$ , $\uparrow$ )	
Student 395 ins. 33 feat.	MaritalStatus ( $\checkmark$ , $\uparrow$ )	
	Sex ( $\times$ , -)	
Student 395 ins. 33 feat.	Single ( $\checkmark$ , $\uparrow$ )	Age ( $\checkmark$ , $\uparrow$ )
	Unemployed ( $\checkmark$ , $\uparrow$ )	Credit ( $\checkmark$ , $\uparrow$ )
Student 395 ins. 33 feat.	PurposeOfLoan ( $\checkmark$ , $\uparrow$ )	LoanDuration ( $\checkmark$ , $\uparrow$ )
	InstallmentRate ( $\checkmark$ , $\uparrow$ )	
Student 395 ins. 33 feat.	Housing ( $\checkmark$ , $\uparrow$ )	
	Sex ( $\times$ , -)	
Student 395 ins. 33 feat.	School ( $\checkmark$ , $\uparrow$ )	MotherEducation ( $\checkmark$ , $\uparrow$ )
	AgeGroup ( $\times$ , -)	FatherEducation ( $\checkmark$ , $\uparrow$ )
Student 395 ins. 33 feat.	FamilySize ( $\checkmark$ , $\uparrow$ )	TravelTime ( $\checkmark$ , $\uparrow$ )
	ParentStatus ( $\checkmark$ , $\uparrow$ )	ClassFailures ( $\checkmark$ , $\uparrow$ )
Student 395 ins. 33 feat.	SchoolSupport ( $\checkmark$ , $\uparrow$ )	GoOut ( $\checkmark$ , $\uparrow$ )
	FamilySupport ( $\checkmark$ , $\uparrow$ )	
Student 395 ins. 33 feat.	ExtraPaid ( $\checkmark$ , $\uparrow$ )	
	ExtraActivities ( $\checkmark$ , $\uparrow$ )	
Student 395 ins. 33 feat.	Nursery ( $\checkmark$ , $\uparrow$ )	
	HigherEdu ( $\checkmark$ , $\uparrow$ )	
Student 395 ins. 33 feat.	Internet ( $\checkmark$ , $\uparrow$ )	
	Romantic ( $\checkmark$ , $\uparrow$ )	
Student 395 ins. 33 feat.	MotherJob ( $\checkmark$ , $\uparrow$ )	
	FatherJob ( $\checkmark$ , $\uparrow$ )	
Student 395 ins. 33 feat.	SchoolReason ( $\checkmark$ , $\uparrow$ )	

TABLE I: Mutability and directionality for each feature.  $\checkmark$  means the feature is mutable (may be modified), while  $\times$  means the feature is immutable.  $\uparrow$  means the feature may increase or decrease, while  $\downarrow$  means the feature may rise only.

## APPENDIX

### A. Datasets

The properties of mutability and directionality for each dataset are shown in Table I.

### B. Subgroups identified

The following figures detail the different subgroups identified for the Adult, Athlete and Student datasets.

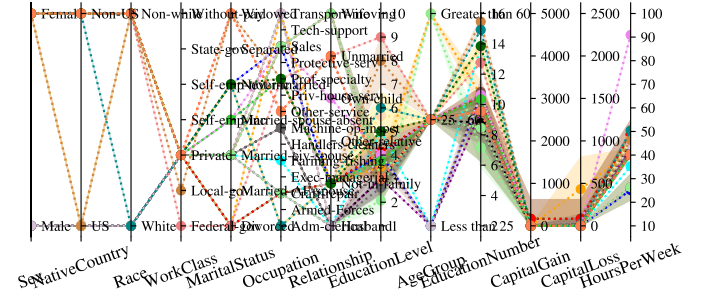


Fig. 1: Adult dataset subgroup details with  $\alpha = 0.1$

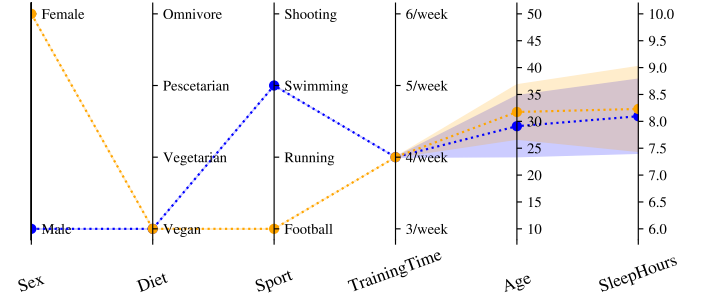


Fig. 2: Athlete dataset subgroup details with  $\alpha = 0.1$

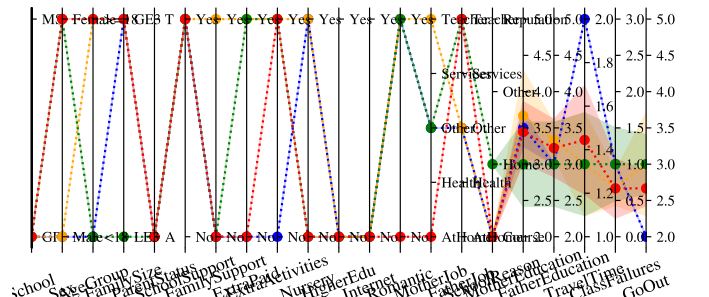


Fig. 3: Student dataset subgroup details with  $\alpha = 0.1$