

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

## APPENDIX

### A. Datasets

The task in the synthetic Athlete dataset is to predict whether an athlete will earn a medal in the olympics. It consists of the following attributes:

- Age: continuous, immutable. Decreases chances linearly from age 25 up to 50 (-1% decrease/year) and remains at 5% at 50 and on.
- Gender: binary, immutable.
- Sport: categorical, mutable, any value. (Football: 0%, Running: +5%, Swimming: +10%, Shooting: +20%).
- Training per week: ordinal, mutable, any direction: (3: -15%, 4: -10%, 5: +15%, 6: +30%).
- daily sleep hours: continuous, mutable, any direction. (+1% increase/hour of daily sleep).
- Diet: Categorical, mutable, any value: (Vegan: +20%, Vegetarian: +15%, Pescetarian: -5%, Omnivore: -10%).

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, Dutch, German and Student datasets.

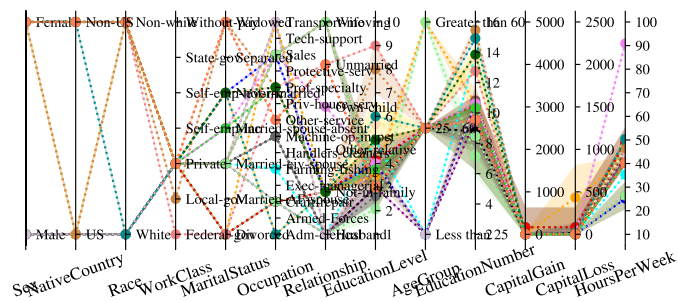


Fig. 1: Adult dataset subgroup details with  $\alpha = 0.1$

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 ( $\checkmark$ , $\uparrow$ )
	Occupation ( $\checkmark$ , $\uparrow$ )	HoursPerWeek ( $\checkmark$ , $\uparrow$ )
Athlete 1000 ins. 6 feat.	Relationship ( $\checkmark$ , $\uparrow$ )	
	Sex ( $\times$ , -)	
	Diet ( $\checkmark$ , $\uparrow$ )	Age ( $\times$ , -)
	Sport ( $\checkmark$ , $\uparrow$ )	SleepHours ( $\checkmark$ , $\uparrow$ )
Compas 7214 ins. 52 feat.	TrainingTime ( $\checkmark$ , $\uparrow$ )	
	Sex ( $\times$ , -)	
	ChargeDegree ( $\checkmark$ , $\uparrow$ )	PriorsCount ( $\checkmark$ , $\uparrow$ )
	Race ( $\times$ , -)	AgeGroup ( $\checkmark$ , $\uparrow$ )
Dutch 60420 ins. 12 feat.	HouseholdPosition ( $\checkmark$ , $\uparrow$ )	
	HouseholdSize ( $\checkmark$ , $\uparrow$ )	
	Country ( $\times$ , -)	EducationLevel ( $\checkmark$ , $\uparrow$ )
	EconomicStatus ( $\checkmark$ , $\uparrow$ )	Age ( $\checkmark$ , $\uparrow$ )
	CurEcoActivity ( $\checkmark$ , $\uparrow$ )	
	MaritalStatus ( $\checkmark$ , $\uparrow$ )	
German 1000 ins. 20 feat.	Sex ( $\times$ , -)	
	Single ( $\checkmark$ , $\uparrow$ )	Age ( $\checkmark$ , $\uparrow$ )
	Unemployed ( $\checkmark$ , $\uparrow$ )	Credit ( $\checkmark$ , $\uparrow$ )
	PurposeOfLoan ( $\checkmark$ , $\uparrow$ )	LoanDuration ( $\checkmark$ , $\uparrow$ )
	InstallmentRate ( $\checkmark$ , $\uparrow$ )	
	Housing ( $\checkmark$ , $\uparrow$ )	
Student 395 ins. 33 feat.	Sex ( $\times$ , -)	
	School ( $\checkmark$ , $\uparrow$ )	
	AgeGroup ( $\times$ , -)	
	FamilySize ( $\checkmark$ , $\uparrow$ )	
	ParentStatus ( $\checkmark$ , $\uparrow$ )	
	SchoolSupport ( $\checkmark$ , $\uparrow$ )	MotherEducation ( $\checkmark$ , $\uparrow$ )
	FamilySupport ( $\checkmark$ , $\uparrow$ )	FatherEducation ( $\checkmark$ , $\uparrow$ )
	ExtraPaid ( $\checkmark$ , $\uparrow$ )	TravelTime ( $\checkmark$ , $\uparrow$ )
	ExtraActivities ( $\checkmark$ , $\uparrow$ )	ClassFailures ( $\checkmark$ , $\uparrow$ )
	Nursery ( $\checkmark$ , $\uparrow$ )	GoOut ( $\checkmark$ , $\uparrow$ )
	HigherEdu ( $\checkmark$ , $\uparrow$ )	
	Internet ( $\checkmark$ , $\uparrow$ )	
	Romantic ( $\checkmark$ , $\uparrow$ )	
	MotherJob ( $\checkmark$ , $\uparrow$ )	
	FatherJob ( $\checkmark$ , $\uparrow$ )	
	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  $\uparrow$  means the feature may rise only.

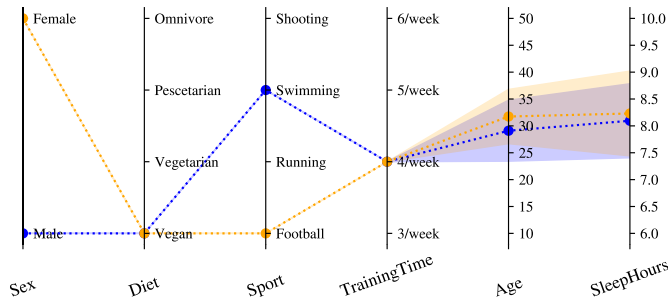


Fig. 2: Athlete dataset subgroup details with  $\alpha = 0.1$

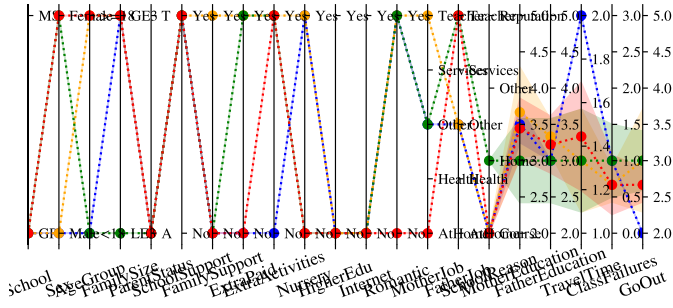


Fig. 3: Student dataset subgroup details with  $\alpha = 0.1$

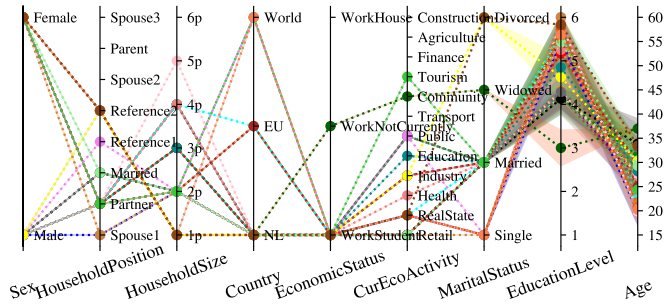


Fig. 4: Dutch dataset subgroup details with  $\alpha = 0.1$

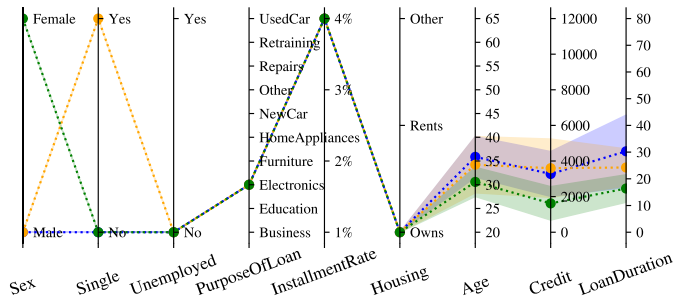


Fig. 5: German dataset subgroup details with  $\alpha = 0.1$