

# **ALGORITHMS FOR SOCIAL GOOD:** FAIRNESS AND BIAS IN DATA-DRIVEN DECISION-MAKING **SYSTEMS**

Elena Beretta

PhD candidate (XXXIII cycle) - Thesis Defense

Nexa Center for Internet & Society, Politecnico di Torino, Italy

Fondazione Bruno Kessler, Trento, Italy

**Supervisors** 

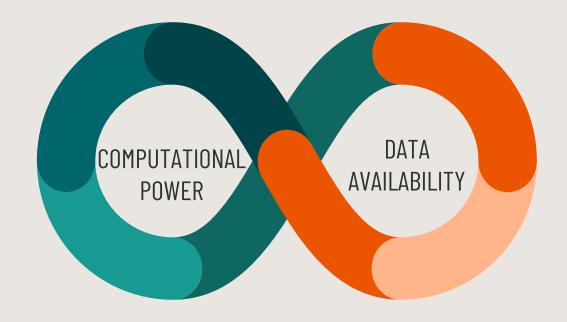
Prof. Juan Carlos De Martin, Politecnico di Torino

Bruno Lepri, FBK

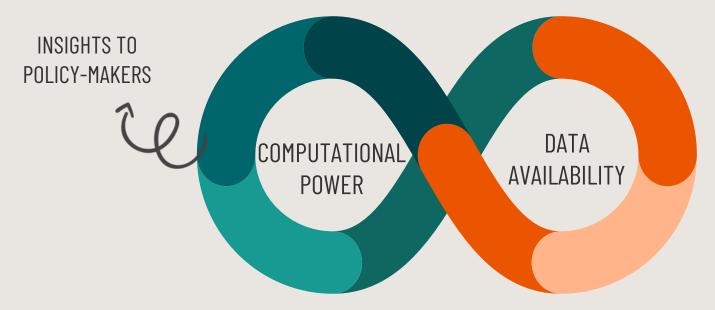
Advisor

Antonio Vetrò, Politecnico di Torino

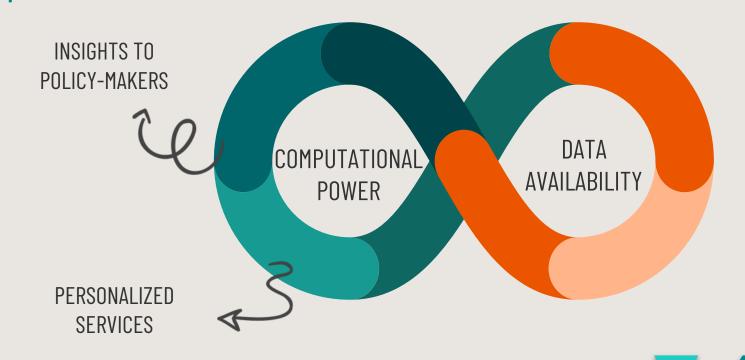
## TABLE OF **CONTENTS** 01 INTRODUCTION & BACKGROUND 02 MOTIVATIONS & GOALS 03 RESEARCH QUESTIONS 04 DATA BIAS AWARENESS 05 FAIRNESS IN RANKING SYSTEMS 06 LONG-TERM FAIRNESS 07 CONCLUSIONS

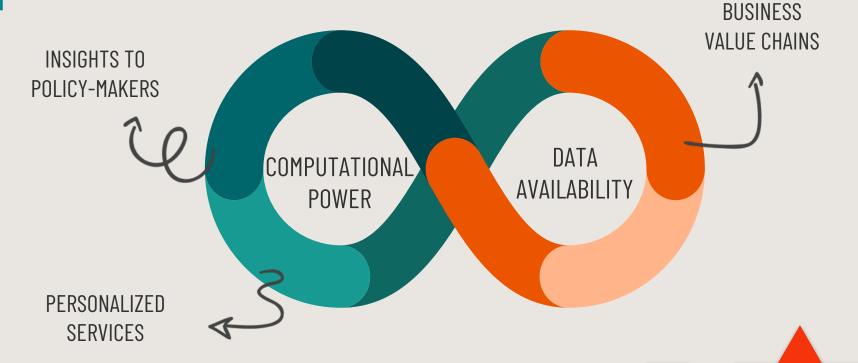


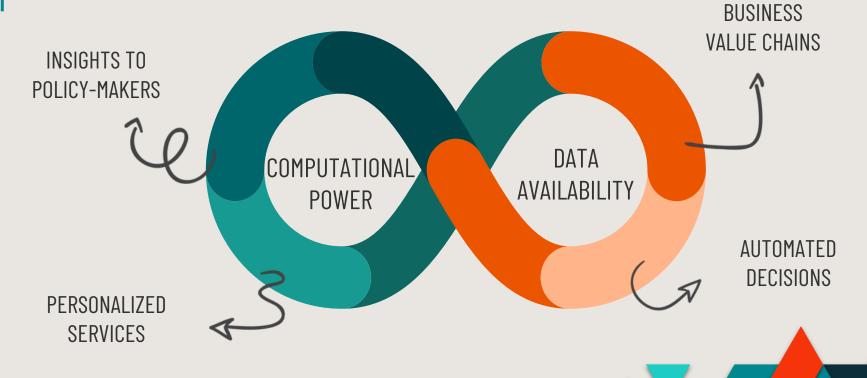












































SEARCH ENGINE

02

INTRODUCTION & BACKGROUND



















SEARCH ENGINE

**PRODUCTS** 



02

















SEARCH ENGINE

**PRODUCTS** 

CULTURE



02











SEARCH ENGINE



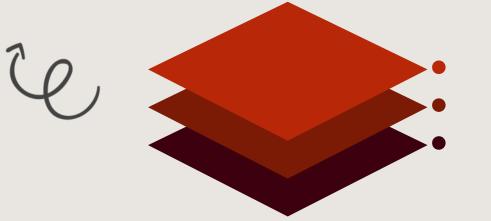
**PRODUCTS** 





# **AUTOMATED DATA-DRIVEN DECISION-MAKING SYSTEMS**

**AUTOMATED DECISION SYSTEM** 



**COMPUTER PROGRAMS** 

**ALGORITHMS** 

DATA

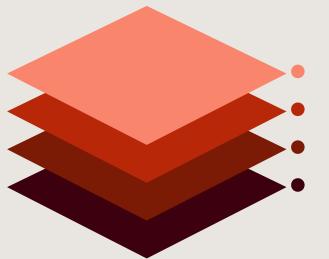


## **AUTOMATED DATA-DRIVEN DECISION-MAKING SYSTEMS**

AUTOMATED DATA-DRIVEN

DECISION-MAKING SYSTEM (ADMs)

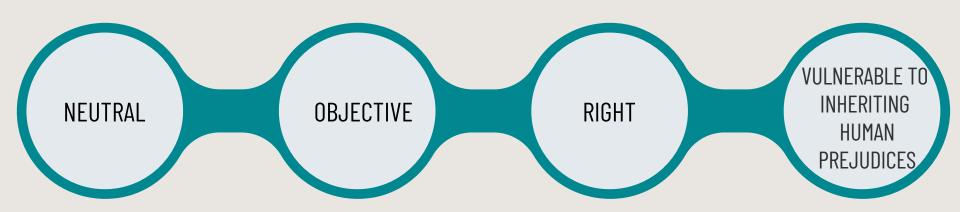




MACHINE LEARNING SYSTEMS
COMPUTER PROGRAMS
ALGORITHMS
DATA

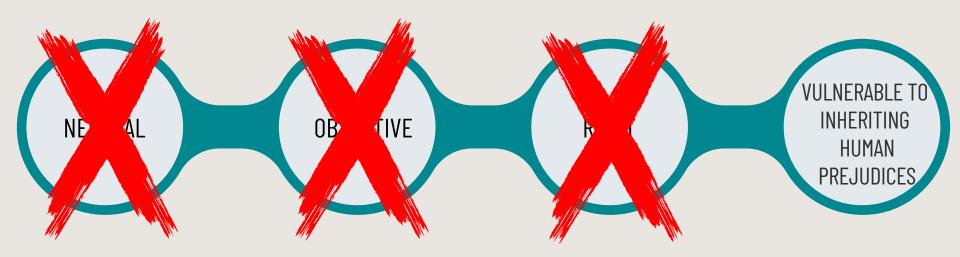


## **DISCRIMINATION IN ADMs**





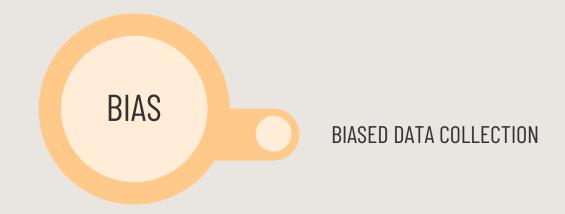
## **DISCRIMINATION IN ADMs**



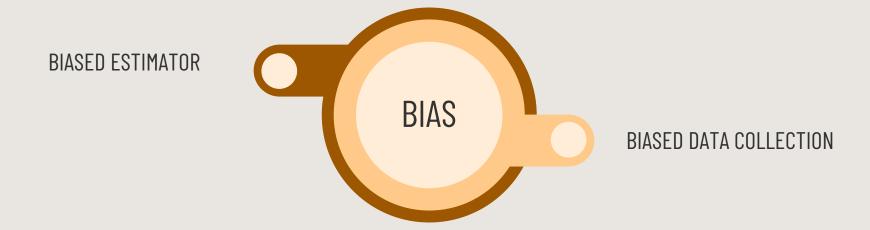


BIAS

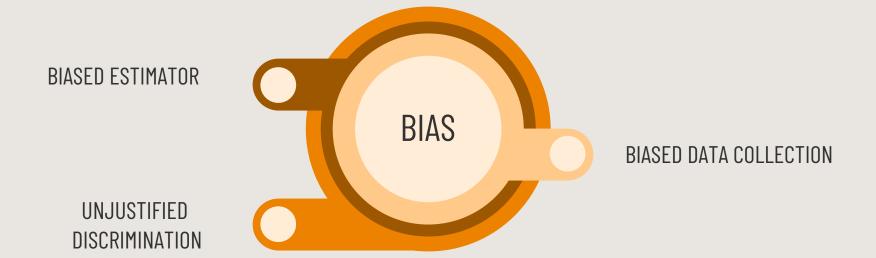




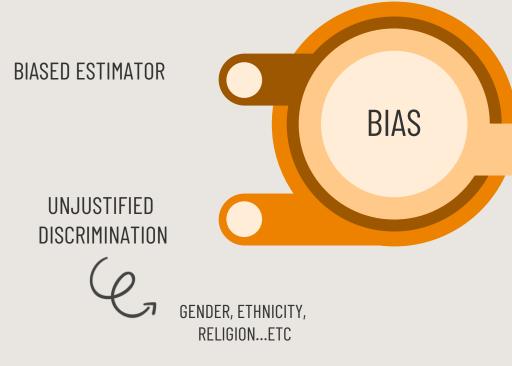








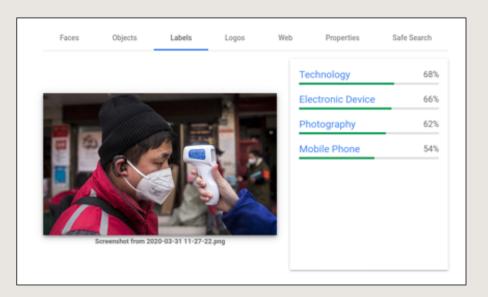


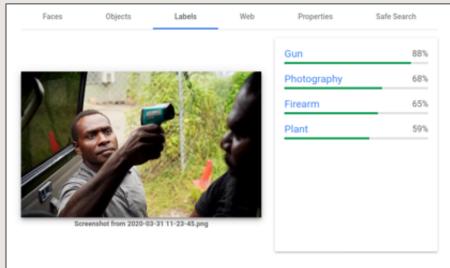


BIASED DATA COLLECTION



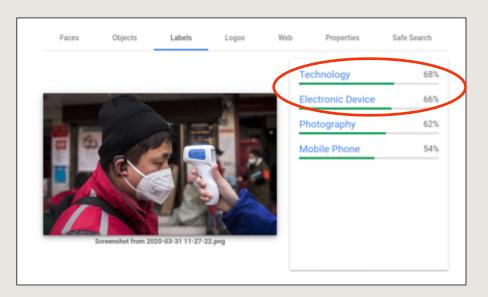
## **UNJUSTIFIED DISCRIMINATION: GOOGLE VISION AI**

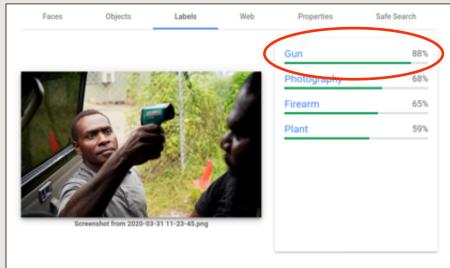






## **UNJUSTIFIED DISCRIMINATION: GOOGLE VISION AI**







## **OVERALL GOALS**



ADDRESSING PROBLEMS OF UNJUSTIFIED DISCRIMINATION IN ADMs



LEVERAGING MULTIPLE
DISCIPLINES TO MITIGATE
THESE UNDESIRABLE
EFFECTS AND TO PROVIDE
CROSS-DISCIPLINARY
PERSPECTIVES OF ANALYSIS

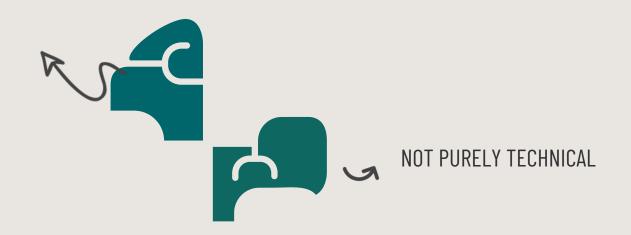


INTERSECTION OF SCIENCE, TECHNOLOGY AND SOCIETY



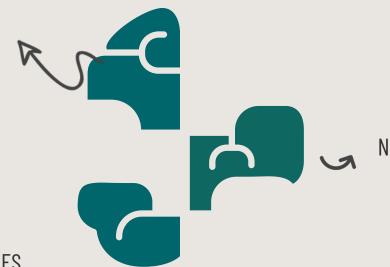


INTERSECTION OF SCIENCE, TECHNOLOGY AND SOCIETY





INTERSECTION OF SCIENCE, TECHNOLOGY AND SOCIETY



NOT PURELY TECHNICAL

SOCIAL SCIENCES AND HUMANITIES

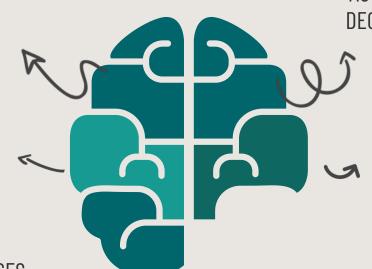


INTERSECTION OF SCIENCE, **TECHNOLOGY** AND SOCIETY **INTERDISCIPLINARY** NOT PURELY TECHNICAL **PERSPECTIVE** SOCIAL SCIENCES AND HUMANITIES

INTERSECTION OF SCIENCE, TECHNOLOGY AND SOCIETY

INTERDISCIPLINARY PERSPECTIVE

SOCIAL SCIENCES AND HUMANITIES



FAIRNESS AND BIAS IN AUTOMATED DATA-DRIVEN DECISION-MAKING SYSTEMS

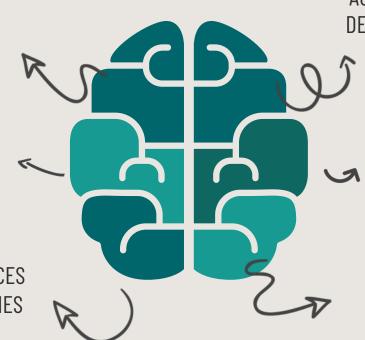
NOT PURELY TECHNICAL



INTERSECTION OF SCIENCE, TECHNOLOGY AND SOCIETY

INTERDISCIPLINARY PERSPECTIVE

SOCIAL SCIENCES AND HUMANITIES



FAIRNESS AND BIAS IN AUTOMATED DATA-DRIVEN DECISION-MAKING SYSTEMS

NOT PURELY TECHNICAL

THEME CHAPTERS

#### THESIS CONTRIBUTION



DATA ANNOTATION
SYSTEM
PREDICTING FUTURE
DISCRIMINATORY RISK
BASED ON BIASED
DATA



RANKING SYSTEM

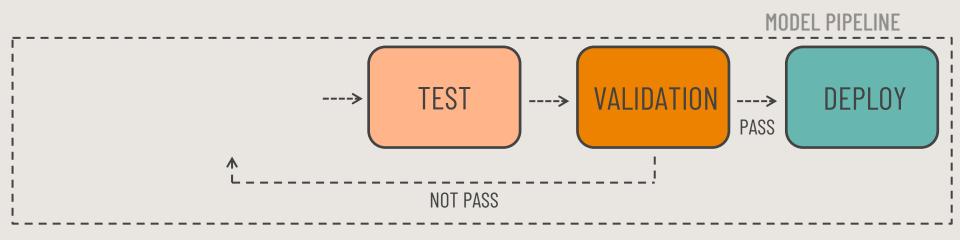
PROPOSING A FAIR-DISTRIBUTIVE RANKING SYSTEM



SYSTEM
MODELING INDIVIDUAL
DYNAMICS FOR
LONG-TERM FAIRNESS

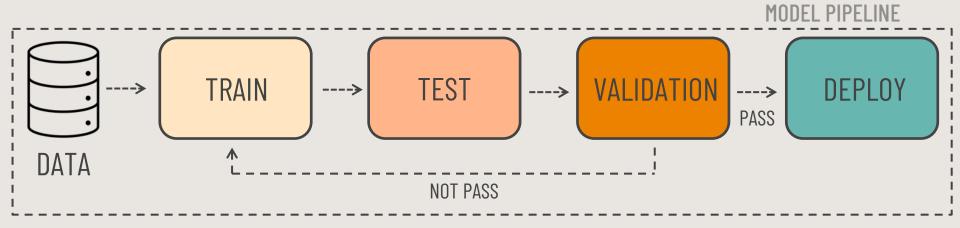


















IS IT POSSIBLE TO ESTABLISH THE A PRIORI PROBABILITY OF THE TRAINING DATA DISTRIBUTION FROM THE AVAILABLE DATASET?

HOW TRAINING DATA COULD INFORM ABOUT THE RISK OF FUTURE DISCRIMINATION?



#### THESIS CONTRIBUTION



DATA ANNOTATION
SYSTEM
PREDICTING FUTURE
DISCRIMINATORY RISK
BASED ON DATA BIAS



RANKING SYSTEM

PROPOSING A FAIR-DISTRIBUTIVE RANKING SYSTEM



DECISION SUPPORT
SYSTEM
MODELING INDIVIDUAL
DYNAMICS FOR
LONG-TERM FAIRNESS





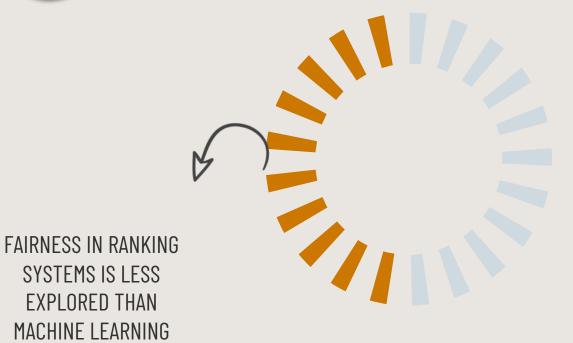
#### **RANKING SYSTEM**







#### **RANKING SYSTEM**



13



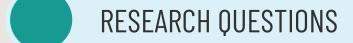
#### RANKING SYSTEM

THE MAJORITY OF THE STUDIES
PROVIDES A DEFINITION OF
EQUITY RATHER THAN GIVING
A SOLUTION TO INEQUALITY

FAIRNESS IN RANKING SYSTEMS IS LESS EXPLORED THAN MACHINE LEARNING







ARE RANKING SYSTEMS BASED ON A
DISTRIBUTIVE FAIRNESS CONSTRAINT ABLE TO
PRESERVE THE ACCURACY OF THE RANKING AND
THE MODEL'S OVERALL UTILITY BY PROVIDING A
RANKING OF THE BEST CANDIDATES?

WHAT ARE THE FACTORS AFFECTING THE FAIRNESS UTILITY TRADE-OFF IN A FAIRNESS CONSTRAINED RANKING SYSTEM?



#### THESIS CONTRIBUTION



DATA ANNOTATION
SYSTEM
PREDICTING FUTURE
DISCRIMINATORY RISK
BASED ON DATA BIAS



RANKING SYSTEM

PROPOSING A FAIR-DISTRIBUTIVE RANKING SYSTEM



DECISION SUPPORT
SYSTEM
MODELING INDIVIDUAL
DYNAMICS FOR
LONG-TERM FAIRNESS

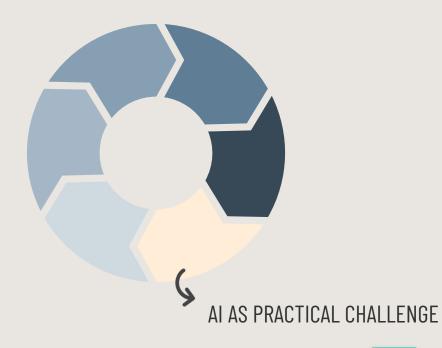




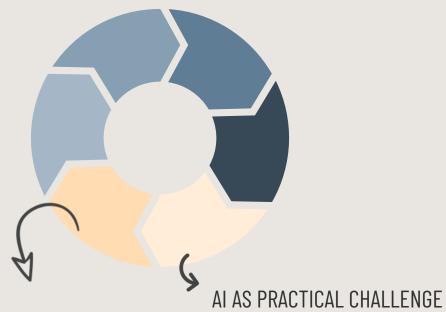








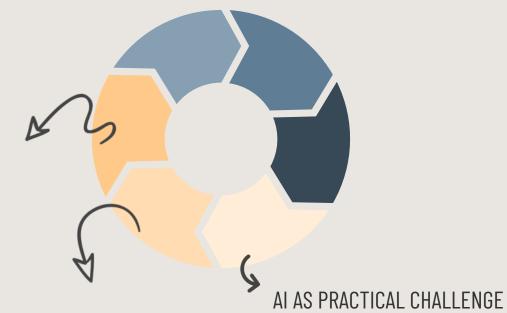






ML STATIC OBJECTIVES

# DECISION SUPPORT SYSTEM







ML STATIC OBJECTIVES





CONSEQUENTIAL DECISIONS RESHAPE THE POPULATION

ML STATIC OBJECTIVES

IMPACT ON SOCIETY

LONG-TERM FAIRNESS
STILL UNEXPLORED





CONSEQUENTIAL DECISIONS RESHAPE THE POPULATION

ONS ION

LONG-TERM FAIRNESS
STILL UNEXPLORED

AI AS PRACTICAL CHALLENGE

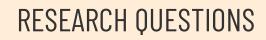


INDIVIDUAL DECISION-MAKING
DYNAMICS CAN AFFECT THE
EFFECTIVENESS OF A POLICY

ML STATIC OBJECTIVES









DO THE FAIRNESS CONSTRAINTS KEEP THEIR VALIDITY FOR AS LONG AS THEY ACT?



# DETECTING DISCRIMINATORY RISK THROUGH DATA ANNOTATION BASED ON BAYESIAN INFERENCES



#### PRINCIPAL TECHNIQUES/CONCEPTS

- BAYESIAN INFERENCE
- BIAS AND FAIRNESS IN DATA AND ML

#### **RESEARCH QUESTIONS**

IS IT POSSIBLE TO ESTABLISH THE A PRIORI PROBABILITY OF THE TRAINING DATA DISTRIBUTION FROM THE AVAILABLE DATASET?

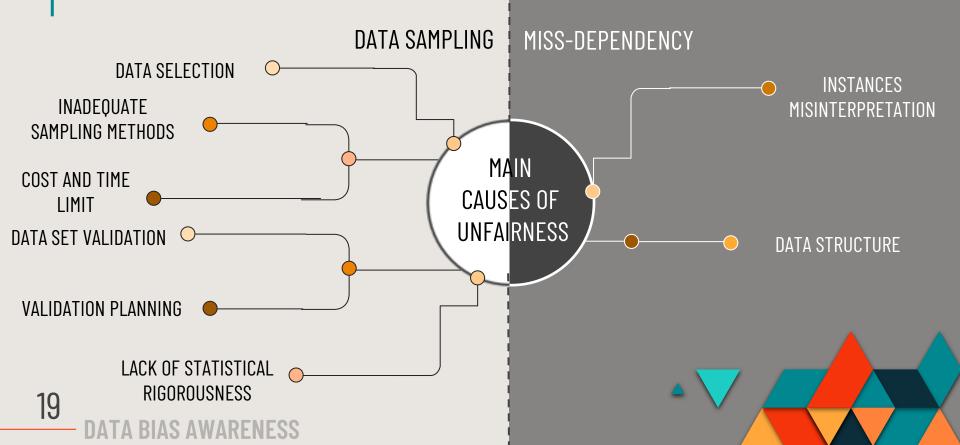
HOW TRAINING DATA COULD INFORM ABOUT THE RISK OF FUTURE DISCRIMINATION?

#### CONTENT

A METHOD OF DATA ANNOTATION BASED ON BAYESIAN STATISTICAL INFERENCE THAT AIMS TO WARN ABOUT THE RISK OF DISCRIMINATORY RESULTS OF A GIVEN DATA SET



### **BACKGROUND**



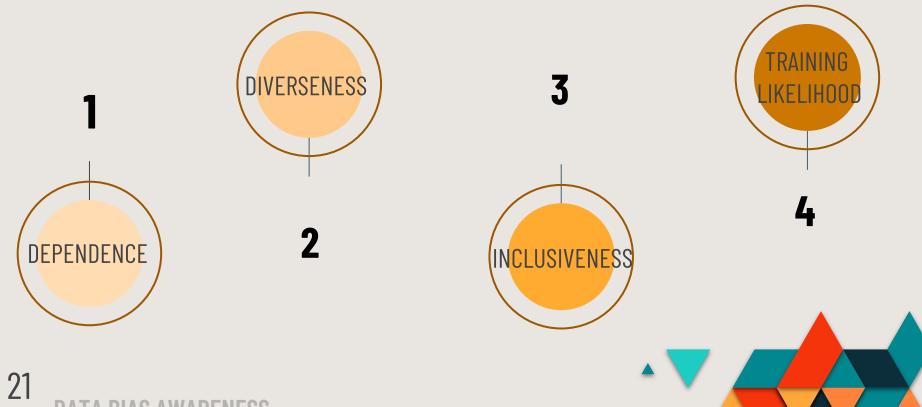
### RESEARCH GOAL



A DIAGNOSTIC FRAMEWORK TO WARN ABOUT THE RISK OF DISCRIMINATORY RESULTS



### THEORETICAL MODEL



#### DEGREE OF CONNECTION





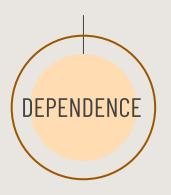
SIMPLIFIED INTERPRETATION OF THE DEPENDENCY

DEGREE OF CONNECTION





DEGREE OF CONNECTION



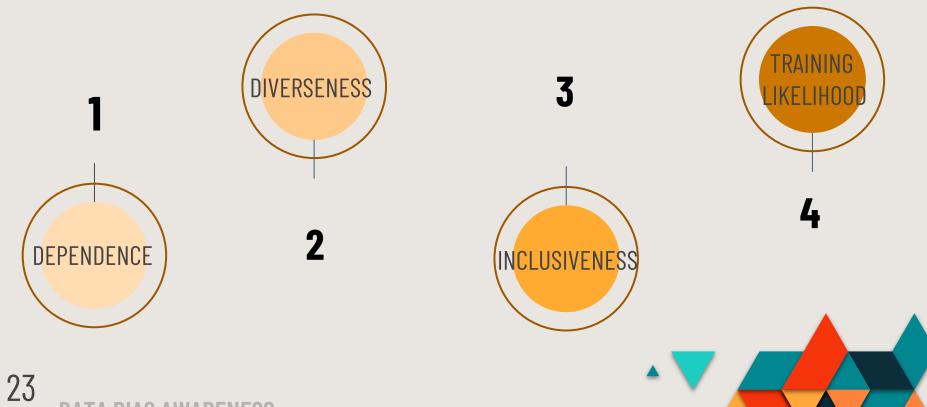
SIMPLIFIED INTERPRETATION OF THE DEPENDENCY



#### EFFECT SIZE INDEX w

MAGNITUDE	VALUE
SMALL	w = 0.1
MEDIUM	w = 0.3
LARGE	w = 0.5







TRAINING DIVERSIFICATION PROBABILITY

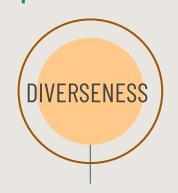




HOW LIKELY A PROPERTY WILL OCCUR

TRAINING DIVERSIFICATION PROBABILITY





TRAINING DIVERSIFICATION PROBABILITY

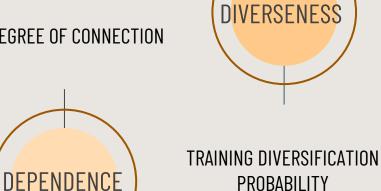
HOW LIKELY A PROPERTY WILL OCCUR



FORMULA	PROBABILITY
P(Y=0) $P(Y=1)$	P = 0.48 P = 0.52
P(A = white) P(A = black) P(A = Asian)	P = 0.6 P = 0.35 P = 0.15



DEGREE OF CONNECTION



PROBABILITY THAT TWO PROPERTIES ARE SIMULTANEOUSLY INCLUDED IN THE TRAINING SET



INCLUSIVENESS

OCCURRENCE LIKELIHOOD OF THE PROTECTED ATTRIBUTE LEVELS GIVEN THE TARGET VARIABLE LEVELS BEFORE THE TRAINING SET IS SAMPLED

PROBABILITY THAT TWO
PROPERTIES
ARE SIMULTANEOUSLY
INCLUDED IN
THE TRAINING SET





PROBABILITY THAT TWO
PROPERTIES
ARE SIMULTANEOUSLY
INCLUDED IN
THE TRAINING SET



THE PROBABILITY THAT THE
TRAINING SET SIMULTANEOUSLY
SHOWS THE PROPERTY
Y = y AND A = a



PROBABILITY THAT TWO
PROPERTIES
ARE SIMULTANEOUSLY
INCLUDED IN
THE TRAINING SET



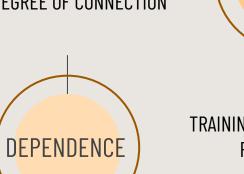
THE PROBABILITY THAT THE
TRAINING SET SIMULTANEOUSLY
SHOWS THE PROPERTY
Y = y AND A = a



FORMULA	PROBABILITY
$P(Y=0 \cap A=white)$ $P(Y=0 \cap A=black)$	P=0.42 P=0.07
$P(Y=0 \cap A=Asian)$	P=0.09
$P(Y=1 \cap A = white)$ $P(Y=1 \cap A = black)$	P=0.18 P=0.28
$P(Y=1 \cap A = Asian)$	P=0.06



**DEGREE OF CONNECTION** 



TRAINING DIVERSIFICATION PROBABILITY



PROBABILITY THAT TWO PROPERTIES
ARE SIMULTANEOUSLY INCLUDED IN
THE TRAINING SET





OCCURRENCE LIKELIHOOD OF THE PROTECTED ATTRIBUTE LEVELS GIVEN THE TARGET VARIABLE LEVELS BEFORE THE TRAINING SET IS SAMPLED





OCCURRENCE LIKELIHOOD OF THE PROTECTED ATTRIBUTE LEVELS GIVEN THE TARGET VARIABLE LEVELS BEFORE THE TRAINING SET IS SAMPLED





OCCURRENCE LIKELIHOOD OF THE PROTECTED ATTRIBUTE LEVELS GIVEN THE TARGET VARIABLE LEVELS BEFORE THE TRAINING SET IS SAMPLED

- i) WHAT IS THE PROBABILITY OF BELONGING TO AN ETHNIC GROUP WITH RESPECT TO THE OUTCOME VARIABLE?
- ii) WHAT IS THE PROBABILITY OF OBTAINING A CERTAIN OUTCOME WITH RESPECT TO THE ETHNIC GROUP?





OCCURRENCE LIKELIHOOD OF THE PROTECTED ATTRIBUTE LEVELS GIVEN THE TARGET VARIABLE LEVELS BEFORE THE TRAINING SET IS SAMPLED

- i) WHAT IS THE PROBABILITY OF BELONGING TO AN ETHNIC GROUP WITH RESPECT TO THE OUTCOME VARIABLE?
- ii) WHAT IS THE PROBABILITY OF
  OBTAINING A CERTAIN
  OUTCOME WITH RESPECT TO
  THE ETHNIC GROUP?

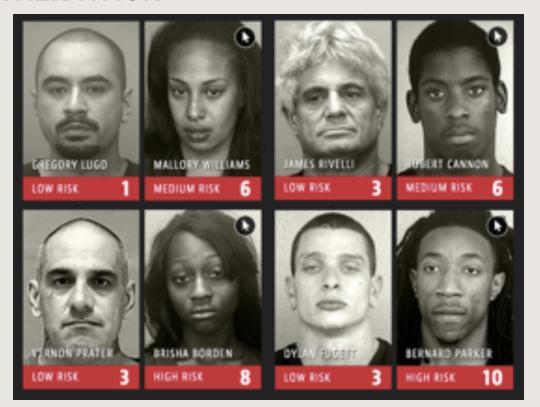
FORMULA	PROBABILITY
P(Y=0 A=white)	P=0.7
P(Y=0 A=black)	P=0.2
P(Y=0 A=Asian)	P=0.6
P(Y=1 A=white)	P=0.3
$P\left(Y=1 A=black\right)$	P=0.8
P(Y=1 A=Asian)	P=0.4
P(A=white Y=1)	P=0.34
P(A=white Y=0)	P=0.87
P(A=black Y=1)	P=0.53
P(A=black Y=0)	P=0.15
P(A=Asian Y=1)	P=0.11
P(A=Asian Y=0)	P=0.18



### **DATASETS PROMINENT PROPERTIES**

	COMPAS	DRUG CONSUMPTION	ADULT CENSUS DATASET
SIZE	6172x9	1885x31	48842x15
TARGET	$0 \rightarrow N0$	0 → NON USER	0 → > 50K
VARIABLE	$1 \rightarrow YES$	1 → USER	1 → ≤ 50K
LEVELS OF	ASIAN	ASIAN BLACK	AMERICAN-INDIAN/ESKIMO
ETHNICITY ATTRIBUTE	BLACK	BLACK/ASIAN	ASIAN-PAC-ISLANDER
ATTRIBUTE	CAUCASIAN	CAUCASIAN	BLACK
	HISPANIC	WHITE/ASIAN	CAUCASIAN
	NATIVE AMERICAN	WHITE/BLACK	OTHER
	OTHER	OTHER	

#### **VALIDATION**



#### **COMPAS**

CORRECTIONAL OFFENDER MANAGEMENT PROFLING FOR ALTERNATIVE SANCTIONS

1 = RECIDIVE 0 = NOT RECIDIVE

BLACK PEOPLE = MORE LIKELY TO BE LABELED AS RECIDIVISTS

WHITE PEOPLE = UNDERESTIMATED RISK OF RECIDIVISM

#### **VALIDATION**



#### **COMPAS**

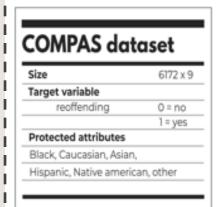
CORRECTIONAL OFFENDER MANAGEMENT PROFLING FOR ALTERNATIVE SANCTIONS

1 = RECIDIVE 0 = NOT RECIDIVE

BLACK PEOPLE = MORE LIKELY TO BE LABELED AS RECIDIVISTS

WHITE PEOPLE = UNDERESTIMATED RISK OF RECIDIVISM

#### **RESULTS**



Dependence	SMALL
	Range [0,1]
Contingency coefficient	0.1413
Effect size	0.1427

#### GRAPHICAL VISUALIZATION OF THE DATA ANNOTATION SYSTEM

			Probability
Target variable			Range [0,1]
0			0.545
1			0.455
Protected a	attribute		
Asian			0.005
Black			0.514
Caucasian			0.341
Hispanic			0.082
Native am.			0.002
other		=	0.056

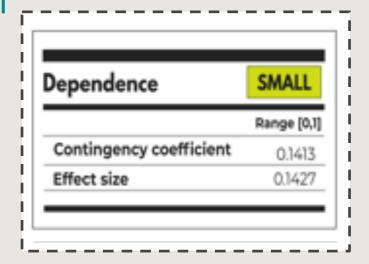
	Probability
	Range [0,1]
P(Asian ∩ 0)	0.0023
P(Asian ∩ 1)	0.0008
P(Black n 0)	0.1514
P(Black n 1)	0.1661
P(Caucasian ∩ 0)	0.1281
P(Caucasian N 1)	0.0822
P(Hispanic ∩ 0)	0.0320
P(Hispanic ∩ 1)	0.0189
P(Native american ∩ 0)	0.0006
P(Native american ∩ 0)	0.0005
P(other ∩ 0)	0.0219
P(other ∩ 1)	0.0124

Training likelihood		
	Probability	
	Range [0,1]	
P(Caucasian   1)	0.293	
P(Caucasian   0 )	0.381	
P(0 Caucasian)	0.609	
P[1  Caucasian)	0.391	
P(Black   1)	0.591	
P(Black   0)	0.450	
P(0 Black)	0.477	
P(1   Black)	0.523	



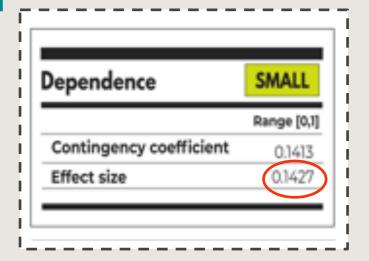
DATASET PROMINENT PROPERTIES





DEPENDENCY RELATIONSHIPS BETWEEN RECIDIVISM AND ETHNIC PROPERTY

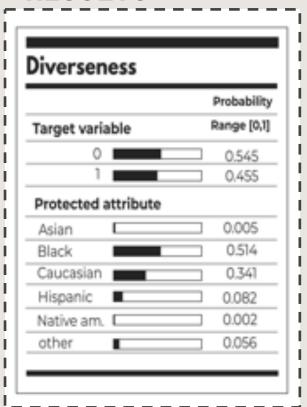




DEPENDENCY RELATIONSHIPS BETWEEN RECIDIVISM AND ETHNIC PROPERTY

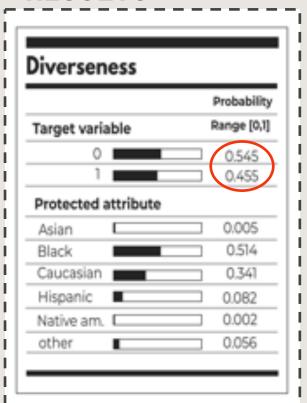






THE PROBABILITY THAT THE TRAINING SET WILL BE EQUALLY COMPOSED BY ETHNIC MINORITIES AND ETHNIC MAJORITIES

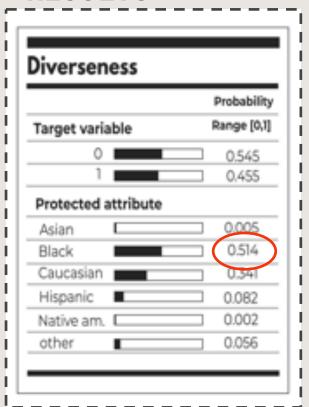




THE PROBABILITY THAT THE TRAINING SET WILL BE EQUALLY COMPOSED BY ETHNIC MINORITIES AND ETHNIC MAJORITIES







THE PROBABILITY THAT THE TRAINING SET WILL BE EQUALLY COMPOSED BY ETHNIC MINORITIES AND ETHNIC MAJORITIES



TARGET VARIABLE EQUALLY DISTRIBUTED



PREVALENCE OF BLACK PEOPLE



	Probability Range [0,1]
P(Asian ∩ 0)	0.0023
P(Asian n 1)	0.0008
P(Black n 0)	0.1514
P(Black n 1)	0.1661
P(Caucasian ∩ 0)	0.1281
P(Caucasian ∩ 1)	0.0822
P(Hispanic ∩ 0)	0.0320
P(Hispanic N 1)	0.0189
P(Native american ∩ 0)	0.0006
P(Native american ∩ 0)	0.0005
P(other n 0)	0.0219
P(other ∩ 1)	0.0124

THE PROBABILITY THAT THE TRAINING SET WILL BE EQUALLY COMPOSED BY ETHNIC MINORITIES AND ETHNIC MAJORITIES SHOWING SIMILAR TARGET LEVELS



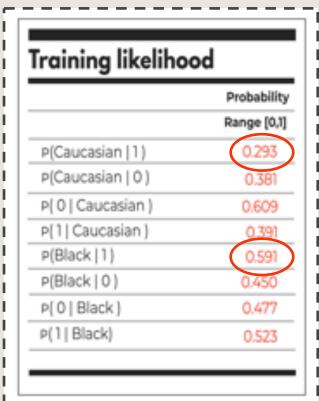
HIGHER LIKELIHOOD THAT CAUCASIANS WILL NOT RECIDIVATE THAN THAT THEY WILL RECIDIVATE



	Probability
	Range [0,1]
P(Caucasian   1)	0.293
P(Caucasian   0 )	0.381
P(0 Caucasian)	0.609
P[1  Caucasian )	0.391
P(Black   1)	0.591
P(Black   0 )	0.450
P(0 Black)	0.477
P(1   Black)	0.523

THE PROBABILITY THAT THE OCCURRENCE OF REOFFENDING IS GIVEN BY THE PROPERTIES OF THE PROTECTED ATTRIBUTE ETHNICITY





# THE PROBABILITY THAT THE OCCURRENCE OF REOFFENDING IS GIVEN BY THE PROPERTIES OF THE PROTECTED ATTRIBUTE ETHNICITY



GIVEN AS VERIFIED THE RECIDIVIST PROPERTY, THE PROBABILITY THAT THE INDIVIDUAL IS BLACK IS MUCH HIGHER THAN THE PROBABILITY THAT THE INDIVIDUAL IS WHITE



	Probability
	Range [0,1]
P(Caucasian   1)	0.293
P(Caucasian   0 )	0.381
P(0 Caucasian)	0.609
P(1   Caucasian )	0.391
P(Black   1)	0.591
P(Black   0 )	0.450
P(0 Black)	0.477
P(1   Black)	0.523

# THE PROBABILITY THAT THE OCCURRENCE OF REOFFENDING IS GIVEN BY THE PROPERTIES OF THE PROTECTED ATTRIBUTE ETHNICITY



GIVEN AS VERIFIED THE RECIDIVIST PROPERTY, THE PROBABILITY THAT THE INDIVIDUAL IS BLACK IS MUCH HIGHER THAN THE PROBABILITY THAT THE INDIVIDUAL IS WHITE



GIVEN AS VERIFIED THE NOT RECIDIVIST PROPERTY, THE PROBABILITY THAT THE INDIVIDUAL IS BLACK IS LOWER THAN THE PROBABILITY THAT THE INDIVIDUAL IS WHITE



- > CLASSICAL SAMPLING VS MACHINE LEARNING PRACTISES
- > REAL POPULATION VS AVAILABLE DATA
- THE STRUCTURE OF THE DATA AFFECTS THE PROBABILITY OF PROPERTIES DISTRIBUTION



# AFters: An automated fair-distributive ranking system for social justice in ai



#### PRINCIPAL TECHNIQUES/CONCEPTS

- RANKING SYSTEMS
- EQUALITY OF OPPORTUNITY
- DISITRIBUTIVE JUSTICE

#### **RESEARCH QUESTIONS**

ARE RANKING SYSTEMS BASED ON A DISTRIBUTIVE FAIRNESS CONSTRAINT ABLE TO PRESERVE THE ACCURACY OF THE RANKING AND THE MODEL'S OVERALL UTILITY BY PROVIDING A RANKING OF THE BEST CANDIDATES?

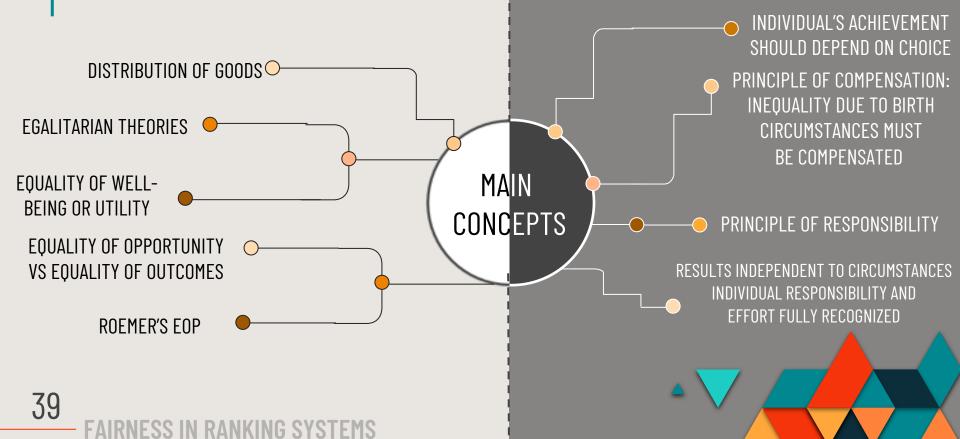
WHAT ARE THE FACTORS AFFECTING THE FAIRNESS UTILITY TRADE-OFF IN A FAIRNESS CONSTRAINED RANKING SYSTEM?

#### CONTENT

AN AUTOMATED FAIR-DISTRIBUTIVE RANKING SYSTEM

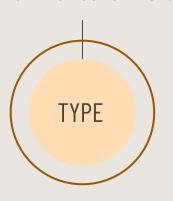


## **BACKGROUND**



# THEORETICAL MODEL

GROUP SAHRING
SAME CIRCUMSTANCES





SAME CHOICES
DIFFERENT CIRCUMSTANCES
DIFFERENT OTUCOME

NEUTRAL TOWARDS INEQUALITIES
WITHIN TYPES



## RESEARCH GOAL



#### AN AUTOMATED FAIR-DISTRIBUTIVE RANKING SYSTEM

- THE BEST TOP-N-RANKING IN A SET OF CANDIDATES.
- MAXIMIZING UTILITY AND SATISFYING FAIRNESS CONSTRAINTS







EQUALITY OF OPPORTUNITY AND DISTRIBUTIVE JUSTICE



#### TYPE

POPULATION PARTITIONED IN A SERIES OF NON-OVERLAPING SETS



VARIABLES DESCRIBING INDIVIDUALS



#### FAIRNESS CONSTRAINTS

A SET OF POLICY: EQUITY, EQUALITY, NEED



#### **EFFORT**

DEGREE OF EFFORT PEOPLE EXERT TO ACCOMPLISH A TASK





TYPES ESTIMATE

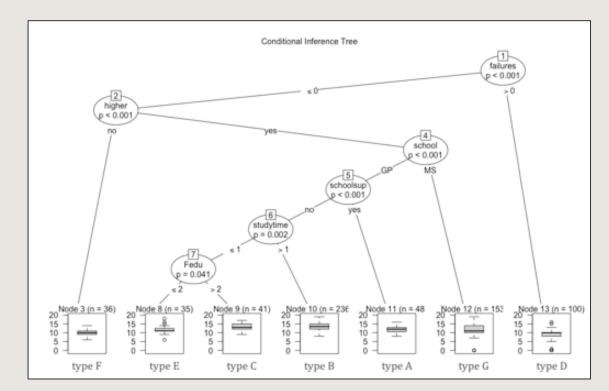
$$T_k = \begin{cases} S_i^X & \text{if } H_0^i : P(Y|X_i) = P(Y) \\ \text{recursion stops} & \text{otherwise} \end{cases}$$

 $S_i^X$  = set of all  $x^i$  possible realizations





TYPES ESTIMATE







#### TYPES ESTIMATE

$$T_k = \begin{cases} S_i^X & \text{if } H_0^i : P(Y|X_i) = P(Y) \\ \text{recursion stops} & \text{otherw} \end{cases}$$

$$if H_0^i: P(Y|X_i) = P(Y)$$
  
otherwise



**EFFORT ESTIMATE** 

# $CDF_{type}(\lambda)$

 $(\lambda)$ = quantile of the Cumulative Distribution **Function** 





#### TYPES ESTIMATE

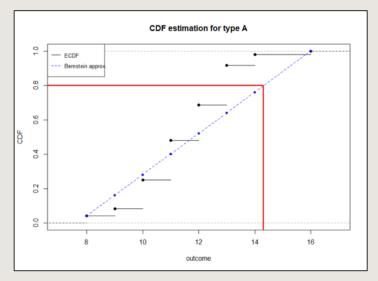
$$T_k = \begin{cases} S_i^X \\ recursion stops \end{cases}$$

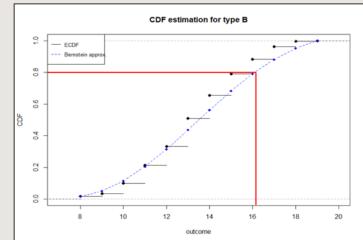
$$if H_0^i: P(Y|X_i) = P(Y)$$
  
otherwise



## EFFORT ESTIMATE

 $CDF_{type}(\lambda)$ 







#### TYPES ESTIMATE

$$T_k = \begin{cases} S_i^X & \text{if } H_0^i : P(Y|X_i) = P(Y) \\ \text{recursion stops} & \text{otherwise} \end{cases}$$



#### **EFFORT ESTIMATE**

$$CDF_{type}(\lambda)$$



UNFAIRNESS DEGREE ESTIMATE

UNFAIRNESS= Gini Index  $\sum CDF_{type}(\lambda)$ 

 $(\lambda)$ = quantile of the Cumulative Distribution Function



#### TYPES ESTIMATE

$$T_k = \begin{cases} S_i^X & \text{if } H_0^i : P(Y|X_i) = P(Y) \\ \text{recursion stops} & \text{otherw} \end{cases}$$

otherwise



#### **EFFORT ESTIMATE**

 $CDF_{type}(\lambda)$ 



#### **UNFAIRNESS DEGREE**

**ESTIMATE** 

UNFAIRNESS= Gini Index  $\sum CDF_{type}(\lambda)$ 

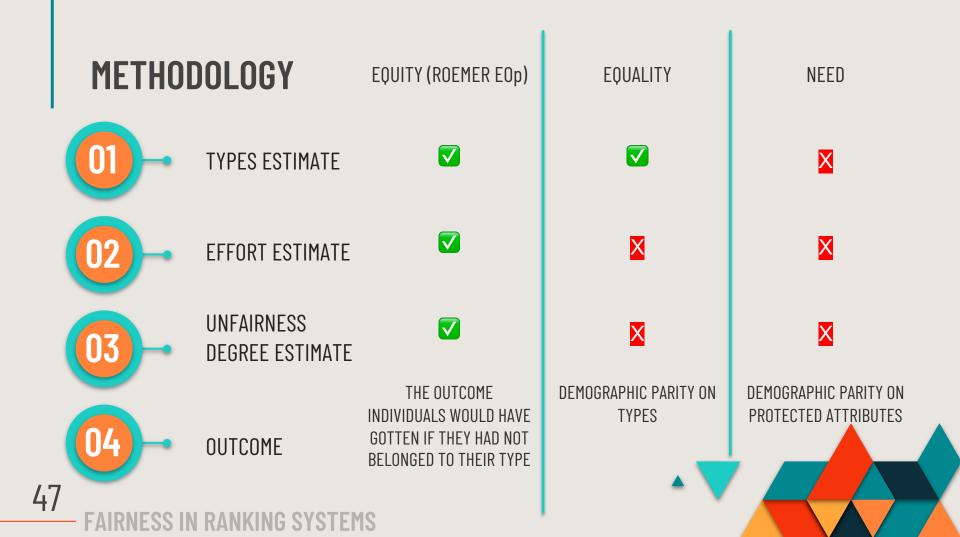


OUTCOME

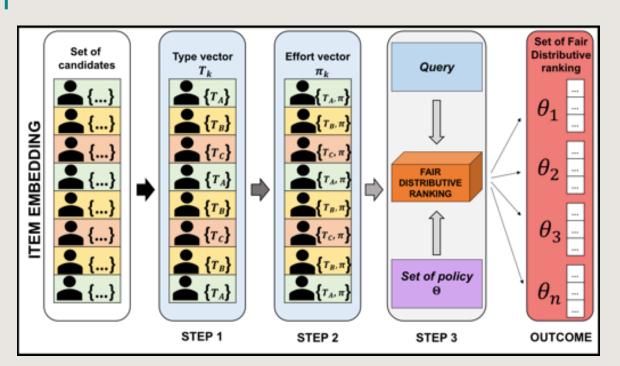
$$y_{\rightarrow} = f(CDF_{type}(\lambda), \sum Gini Index(CDF_{type}(\lambda))$$







#### **APPLICATION SETTING**



DATA: STUDENT PERFORMANCE DATASET

SCENARIO: HYPOTHETICAL SCENARIO OF A UNIVERSITY SELECTION PROCESS IN WHICH THE DECISION-MAKER DETERMINES WHICH STUDENTS ARE SUITABLE ON THE BASIS OF THEIR PERSONAL QUALIFICATIONS AND ACHIEVEMENTS, SO AS TO MAXIMIZE THE INSTITUTION UTILITY

#### SYSTEM'S GOAL:

 $\Gamma = maxmin(utiliy, unfairness)$ 



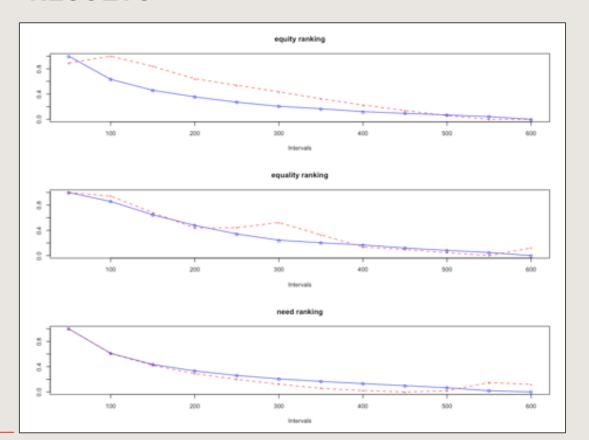
## **VALIDATION**





## **VALIDATION**





UNFAIRNESS-UTILITY
TRADE-OFF FOR RANKINGS
UNDER FAIRNESS
CONSTRAINTS

**RED LINE: UNFAIRNESS** 

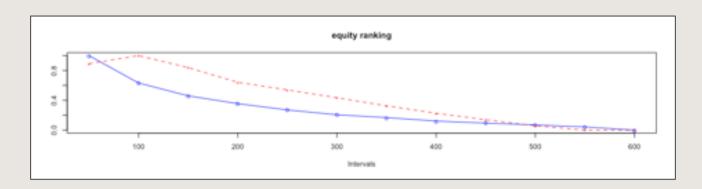
**BLUE LINE: UTILITY** 

Y-AXIS: UTILITY AND INEQUALITY VALUES RANGING FROM 0 TO 1

X-AXIS: NUMEROSITY OF THE RANKINGS



# **EQUITY POLICY: UNFAIRNESS-UTILITY TRADE-OFF**





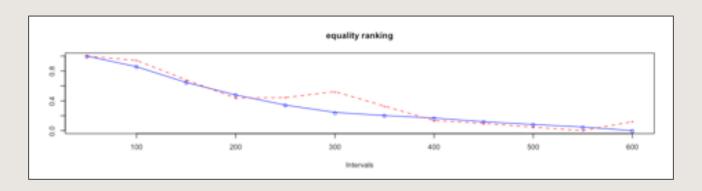
GOOD GENERAL UTILITY IN THE FIRST TOP-N-RANKING



HIGH LEVELS OF INEQUALITY



# **EQUALITY POLICY: UNFAIRNESS-UTILITY TRADE-OFF**





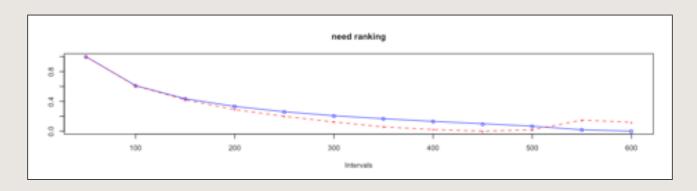
LOW LEVELS OF INEQUALITY



LOW GENERAL UTILITY IN THE FIRST TOP-N-RANKING



# **NEED POLICY: UNFAIRNESS-UTILITY TRADE-OFF**





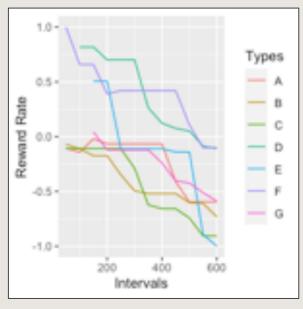
UNIFORM LEVELES OF INEQUALITY AND UTILITY

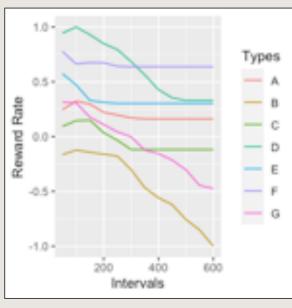


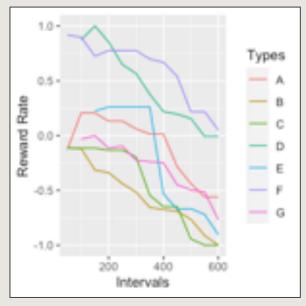
OVERALL WORST PERFORMANCES



## **REWARD RATE**







**EQUITY POLICY** 

**EQUALITY POLICY** 



- > GROUPS VS INDIVIDUAL
- EQUITY AND UTILITY RESULTS ARE POLICY DEPENDENT
- > EQUITY POLICY: BEST CHOICE FOR NUMEROUS RANKING
- > EQUALITY POLICY: BEST CHOICE FOR LESS DENSE RANKING
- > THE MORAL GROUND IS CONTEXT DEPENDENT





## A DECISION SUPPORT SYSTEM FOR LONG-TERM FAIRNESS



#### PRINCIPAL TECHNIQUES/CONCEPTS

- DECISION THEORY
- LONG-TERM FAIRNESS FOR CLASSIFICATION SYSTEMS

#### CONTENT

A DECISION SUPPORT SYSTEM TO ENSURE LONG-TERM FAIRNESS IN MACHINE LEARNING SYSTEMS

#### RESEARCH QUESTIONS

HOW TO CHOOSE THE BEST POLICY TO ENSURE LONG-TERM FAIRNESS?

DO THE FAIRNESS CONSTRAINTS KEEP THEIR VALIDITY FOR AS LONG AS THEY ACT?



## RESEARCH GOAL



#### A DECISION SUPPORT SYSTEM TO ENSURE LONG-TERM FAIRNESS

- DECISION THEORY APPLIED TO ALGORTHMIC DECISION-MAKING
- 2. INDIVIDUAL DYNAMICS INTEGRATED IN THE MODEL



# THEORETICAL MODEL



# THEORETICAL MODEL

#### INSTITUTION'S UTILITY

$$U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta)$$

 $g_{A}$ ,  $g_{B}$  = fractions of the population  $\gamma_{1}$ = positive behavior

**BEST POLICY** 

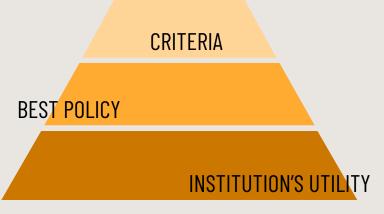
INSTITUTION'S UTILITY

POPULATION IS INDUCED TO PERFORM A POSITIVE BEHAVIOR

$$U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta)$$

 $g_{A'}$   $g_{B}$  = fractions of the population  $\gamma_1$  = positive behavior



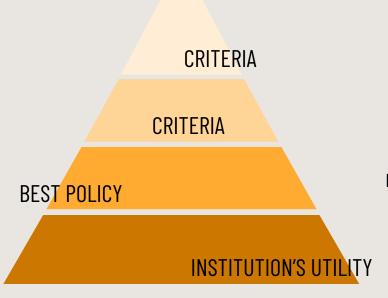


$$max \int_0^\infty \gamma_1(g_A) + \gamma_1(g_B)$$

POPULATION IS INDUCED TO PERFORM A POSITIVE BEHAVIOR

$$U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta)$$

 $g_{A}$ ,  $g_{B}$  = fractions of the population  $\gamma_{1}$ = positive behavior



$$min\left|\int_{0}^{\infty}\gamma_{1}(g_{A})-\gamma_{1}(g_{B})\right|$$

$$max \int_0^\infty \gamma_1(g_A) + \gamma_1(g_B)$$

POPULATION IS INDUCED TO PERFORM A POSITIVE BEHAVIOR

$$U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta)$$

 $g_{A}$ ,  $g_{B}$  = fractions of the population  $\gamma_{1}$ = positive behavior

SOCIETY

 $g_{A'} g_B$ 

MINORITY AND MAJORITY GROUPS

 $\gamma_1$ 

POSITIVE BEHAVIOR

 $\theta \in \Theta$ 

SET OF INDIVIDUALS' ATTRIBUTES

$$U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta) \qquad \text{INSTITUTION'S UTILITY}$$

$$max \int_0^\infty \gamma_1(g_A) + \gamma_1(g_B)$$
 MAXIMIZING DOMINANCE

$$min \left| \int_{\Omega} \gamma_1(g_A) - \gamma_1(g_B) \right|$$

 $\gamma_1(g_A) - \gamma_1(g_B)$  MINIMIZING DOMINANCE AMONG GROUPS

SOCIETY

MINORITY AND MAJORITY GROUPS  $g_{A'} g_B$ 

 $\gamma_1$ 

POSITIVE BEHAVIOR

 $\theta \in \Theta$ 

SET OF INDIVIDUALS' ATTRIBUTES

$$U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta)$$
 INSTITUTION'S UTILITY

$$max \int_0^\infty \gamma_1(g_A) + \gamma_1(g_B)$$
 MAXIMIZING DOMINANCE

$$min \left| \int \gamma_1(g_A) - \gamma_1(g_B) \right|$$
 MINIMIZING DOMINANCE AMONG GROUPS

INDIVIDUAL DYNAMICS

**ALTERNATIVES** 

SOCIETY

MINORITY AND MAJORITY GROUPS

 $\gamma_1$ 

POSITIVE BEHAVIOR

 $\theta \in \Theta$ 

 $g_{A'} g_B$ 

SET OF INDIVIDUALS' ATTRIBUTES

$$U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta)$$
 INSTITUTION'S UTILITY

$$max \int_0^\infty \gamma_1(g_A) + \gamma_1(g_B)$$
 MAXIMIZING DOMINANCE

$$min \left| \int \gamma_1(g_A) - \gamma_1(g_B) \right|$$
 MINIMIZING DOMINANCE AMONG GROUPS

INDIVIDUAL DYNAMICS

**ALTERNATIVES** 

**SCENARIOS** 

Ω



SOCIETY

MINORITY AND MAJORITY GROUPS

 $\gamma_1$ 

POSITIVE BEHAVIOR

 $\theta \in \Theta$ 

 $g_{A'} g_B$ 

SET OF INDIVIDUALS' ATTRIBUTES

$$U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta)$$
 INSTITUTION'S UTILITY

$$max \int_0^\infty \gamma_1(g_A) + \gamma_1(g_B)$$
 MAXIMIZING DOMINANCE

$$min \left| \int \gamma_1(g_A) - \gamma_1(g_B) \right|$$
 MINIMIZING DOMINANCE AMONG GROUPS

INDIVIDUAL DYNAMICS

**ALTERNATIVES** 

**SCENARIOS** 

**IMPACTS** 

Ω

SOCIETY

INDIVIDUAL DYNAMICS

 $g_{A'} g_B$ 

MINORITY AND MAJORITY GROUPS

 $\gamma_1$ 

POSITIVE BEHAVIOR

 $\theta \in \Theta$ 

SET OF INDIVIDUALS' ATTRIBUTES

$$U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta)$$

INSTITUTION'S UTILITY

$$max \int_0^\infty \gamma_1(g_A) + \gamma_1(g_B)$$
 MAXIMIZING DOMINANCE

$$min \left| \int_{-\infty}^{\infty} \gamma_1(g_A) - \gamma_1(g_B) \right|$$

 $\gamma_1(g_A) - \gamma_1(g_B)$  MINIMIZING DOMINANCE AMONG GROUPS

**ALTERNATIVES** 

**SCENARIOS** 

**IMPACTS** 

UTILITY FUNCTION

Ω



SOCIETY

INDIVIDUAL DYNAMICS

 $g_{A'}\,g_B$ 

MINORITY AND MAJORITY GROUPS

 $\gamma_1$ 

POSITIVE BEHAVIOR

 $\theta \in \Theta$ 

SET OF INDIVIDUALS' ATTRIBUTES

$$U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta)$$

INSTITUTION'S UTILITY

$$max \int_0^\infty \gamma_1(g_A) + \gamma_1(g_B)$$

MAXIMIZING DOMINANCE

 $\min\left|\int \gamma_1(g_A) - \gamma_1(g_B)
ight|$  MINIMIZING DOMINANCE AMONG GROUPS

ALTERNATIVES

SCENARIOS  $\Omega$ 

\_

UTILITY FUNCTION

DECIDERS

**IMPACTS** 

D

59

**LONG-TERM FAIRNESS** 

SOCIETY

INDIVIDUAL DYNAMICS

 $g_{A'} g_B$ 

MINORITY AND MAJORITY GROUPS

**ALTERNATIVES** 

Ω

 $\gamma_1$ 

POSITIVE BEHAVIOR

 $\theta \in \Theta$ 

SET OF INDIVIDUALS' ATTRIBUTES

UTILITY FUNCTION

 $U_I^* = \sum \lambda_A g_A(\theta) \cdot \lambda_B g_B(\theta)$ 

 $max \int_0^\infty \gamma_1(g_A) + \gamma_1(g_B)$  MAXIMIZING DOMINANCE

INSTITUTION'S UTILITY

**DECIDERS** 

**SCENARIOS** 

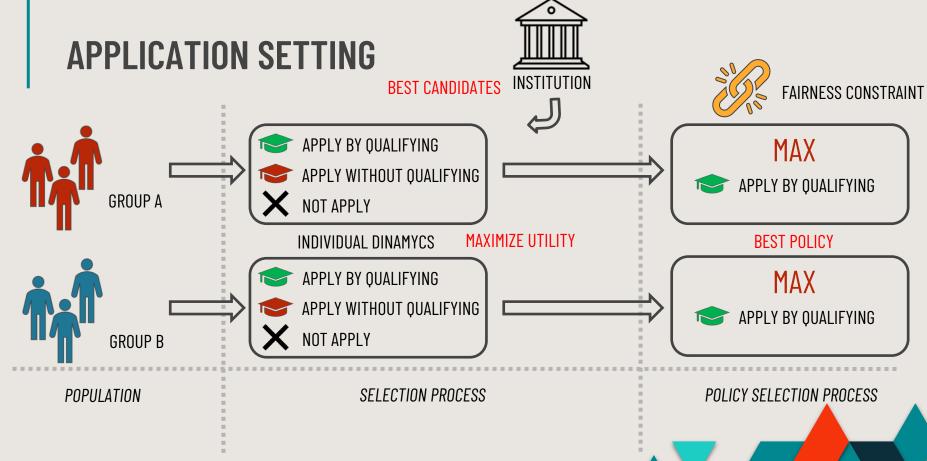
**IMPACTS** 

D

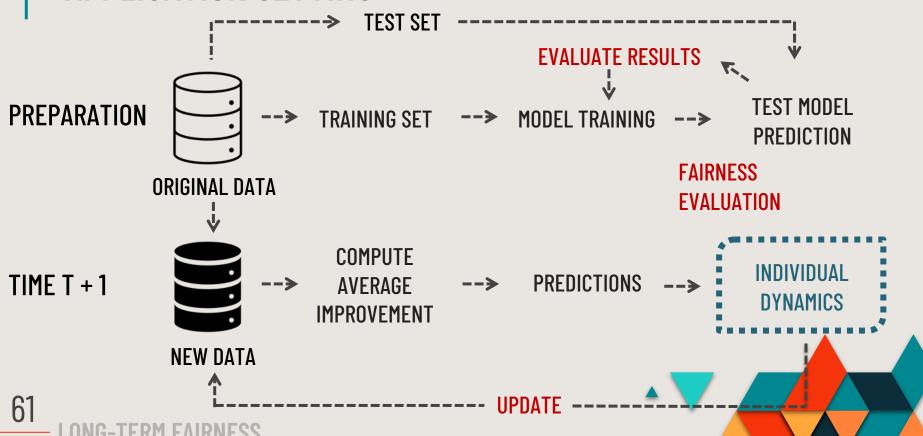
 $\gamma_1(g_A) - \gamma_1(g_B)$  MINIMIZING DOMINANCE AMONG GROUPS

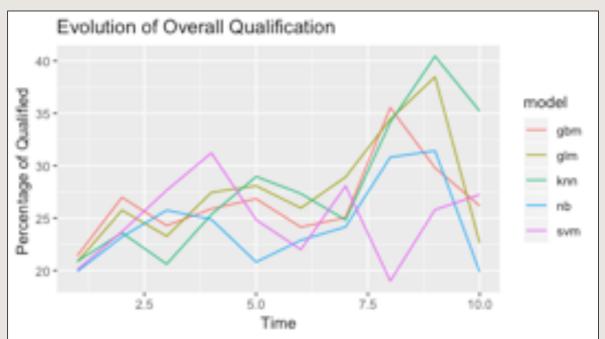
PREFERENCES' FUNCTION

П



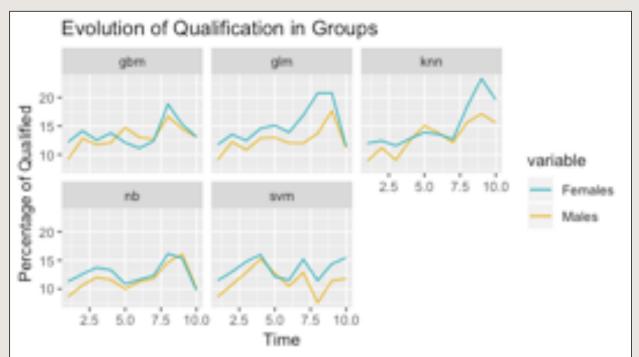
#### **APPLICATION SETTING**

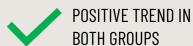














#### MAXIMIZING DOMINANCE

#### MINIMIZING DOMINANCE AMONG GROUPS

$$\int_0^{10} \gamma_1(g_A) + \gamma_1(g_B)$$

$$\left(\left|\int\limits_{0}^{\infty}\gamma_{1}(g_{A})-\gamma_{1}(g_{B})\right|\right)$$

GRADIENT BOOSTING MACHINE

GENERALIZED LINEAR MODEL

K-NEAREST NEIGHBOUR

NAIVE BAYES CLASSIFIER

SUPPORT VECTOR MACHINE

#### MAXIMIZING DOMINANCE

#### MINIMIZING DOMINANCE AMONG GROUPS

$$\int_0^{10} \gamma_1(g_A) + \gamma_1(g_B)$$

$$\left[\left|\int\limits_0^\infty \gamma_1(g_A)-\gamma_1(g_B)
ight|
ight]$$

GRADIENT BOOSTING MACHINE

GENERALIZED LINEAR MODEL

K-NEAREST NEIGHBOUR

NAIVE BAYES CLASSIFIER

SUPPORT VECTOR MACHINE

64

#### MAXIMIZING DOMINANCE

#### MINIMIZING DOMINANCE AMONG GROUPS

$$\int_0^{10} \gamma_1(g_A) + \gamma_1(g_B)$$

$$\left| \int_0^\infty \gamma_1(g_A) - \gamma_1(g_B) \right|$$

-----

GRADIENT BOOSTING MACHINE

242.31

3.48

GENERALIZED LINEAR MODEL

253.33

15.83

K-NEAREST NEIGHBOUR

223.82

9.02

NAIVE BAYES CLASSIFIER

225.95

17.83

SUPPORT VECTOR MACHINE

254.06

25.19

- OUR SYSTEM IS EFFICIENT IN ANALYZING THE LONG-TERM EFFECTS OF POLICIES
- > POLICIES HAVE DIFFERENT INFLUENCES ON GROUPS IN A NON-ONE-STEP MODEL
- FAIRNESS IS NOT CONSISTENT OVER TIME
- FAIRNESS CONSTRAINTS DO NOT NECESSARILY KEEP THEIR VALIDITY FOR AS LONG AS THEY ACT
- > INDIVIDUAL DYNAMICS AFFECT SYSTEM OUTCOMES







#### THESIS CONTRIBUTION



DATA ANNOTATION
SYSTEM
PREDICTING FUTURE
DISCRIMINATORY RISK
BASED ON BIASED
DATA



RANKING SYSTEM

PROPOSING A FAIR-DISTRIBUTIVE RANKING SYSTEM



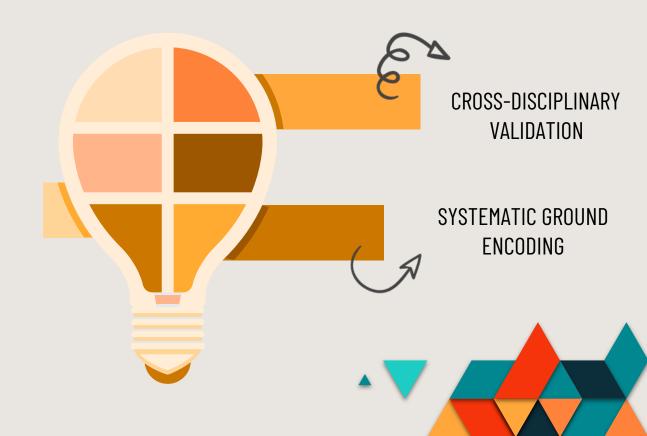
SYSTEM
MODELING INDIVIDUAL
DYNAMICS FOR
LONG-TERM FAIRNESS

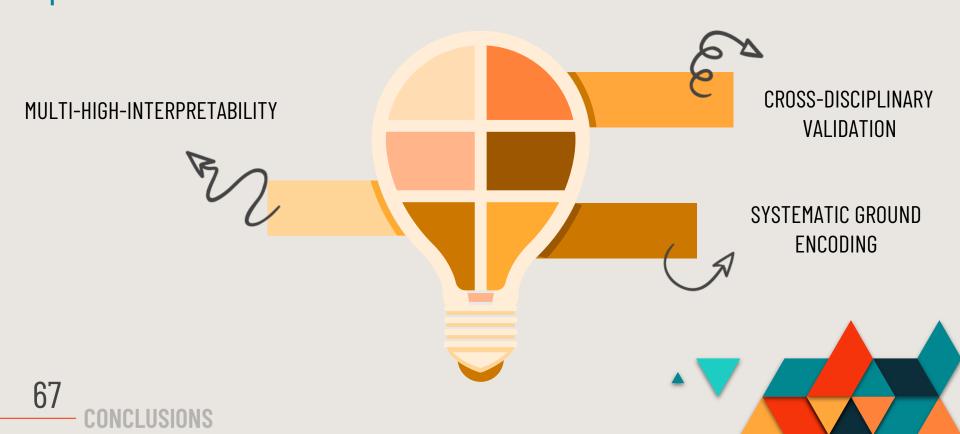












#### **PUBLICATIONS**

#### **Published**

Elena Beretta, Antonio Vetrò, Bruno Lepri, Juan Carlos De Martin. (2021) **Detecting discriminatory risk through data annotation based on Bayesian inferences**. In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21). Association for Computing Machinery, New York, NY, USA, 794–804. DOI:https://doi.org/10.1145/3442188.3445940

Antonio Vetrò, Antonio Santangelo, Elena Beretta, Juan Carlos De Martin. (2019) **Al: from rational agents to socially responsible agents**, *Digital Policy, Regulation and Governance*, https://doi.org/10.1108/DPRG-08-2018-0049

Elena Beretta, Antonio Santangelo, Antonio Vetrò, Bruno Lepri, Juan Carlos De Martin. (2019) **The Invisible Power of Fairness. How Machine Learning Shapes Democracy**. In: Meurs MJ., Rudzicz F. (eds) *Advances in Artificial Intelligence. Canadian Al 2019*. Lecture Notes in Computer Science, vol 11489. Springer, Cham. https://doi.org/10.1007/978-3-030-18305-9\_19

Elena Beretta, Antonio Vetrò, Bruno Lepri, Juan Carlos De Martin. (2018) **Ethical and Socially-Aware Data Labels**, Information Management and Big Data. SIMBig 2018. *Communications in Computer and Information Science*, vol 898, pp. 320-327, Springer, Cham.

#### **Accepted for publication**

Elena Beretta, Antonio Vetrò, Bruno Lepri, Juan Carlos De Martin. (2021) **Equality of opportunity in ranking: a Fair-Distributive Model**, Second International Workshop on Algorithmic Bias in Search and Recommendation (April 2021)

# THANK YOU

Elena Beretta PhD candidate (XXXIII cycle) - Thesis Defense

Nexa Center for Internet & Society, Politecnico di Torino, Italy

Fondazione Bruno Kessler, Trento, Italy

**Supervisors** Prof. Juan Carlos De Martin, Politecnico di Torino

Bruno Lepri, Fondazione Bruno Kessler

Advisor Antonio Vetrò, Politecnico di Torino

# **APPENDIX**



## DATA BIAS AWARENESS

DATA BIAS AND CONDITIONAL PROBABILITIES

QUANTIFYING DEPENDENCE

ESTIMATING DIVERSENESS

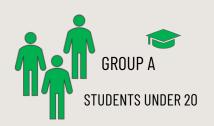
ESTIMATING INCLUSIVENESS

ESTIMATING TRAINING LIKELIHOOD

#### **APPENDIX**



#### DATA BIAS AND CONDITIONAL PROBABILITIES



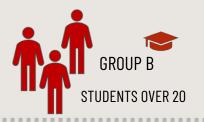
CASE 1

THE ACTUAL PROPERTIES DIFFER ACROSS GROUPS



OBSERVED SPACE = CONSTRUCT SPACE

DECISION SPACE → CORRECT MAPPING



CASE 2

THE ACTUAL PROPERTIES ARE DIFFERENT FROM THOS OBSERVED



OBSERVED SPACE ≠ CONSTRUCT SPACE

DECISION SPACE → ERRONEOUS MAPPING

**POPULATION** 

CASES

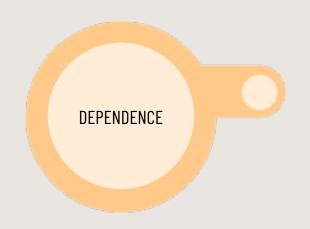


DEPENDENCE
ASSESSES THE DEGREE OF CONNECTION
AMONG THE PROTECTED ATTRIBUTE AND
THE TARGET VARIABLE

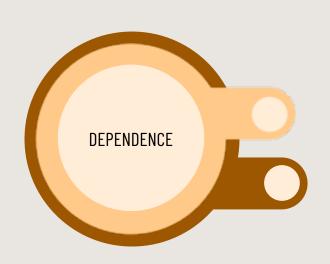
DEPENDENCE



$$C(x_i;y_j) = f(x_i,y_j) - f'(x_i,y_j)$$







$$C(x_i;y_j) = f(x_i,y_j) - f'(x_i,y_j)$$

$$\chi^2 = \sum_{i,j} \frac{C^2(x_i; y_i)}{n_{i,j}} = n \left( \sum_{i,j} \frac{n_{i,j}^2}{n_{i,0} n_{0,j}} \right)$$



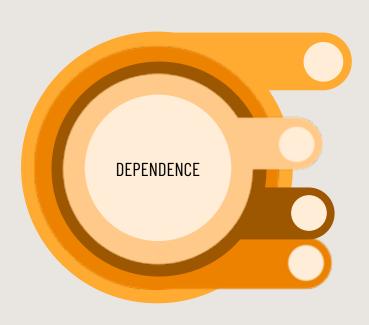


$$C(x_i;y_j) = f(x_i,y_j) - f'(x_i,y_j)$$

$$\chi^2 = \sum_{i,j} \frac{C^2(x_i; y_i)}{n_{i,j}} = n \left( \sum_{i,j} \frac{n_{i,j}^2}{n_{i,0} n_{0,j}} \right)$$

$$C = \sqrt{\frac{\chi^2}{\chi^2 + n}}$$





$$C(x_i;y_j) = f(x_i,y_j) - f'(x_i,y_j)$$

$$\chi^2 = \sum_{i,j} \frac{C^2(x_i; y_i)}{n_{i,j}} = n \left( \sum_{i,j} \frac{n_{i,j}^2}{n_{i,0} n_{0,j}} \right)$$

$$C = \sqrt{\frac{\chi^2}{\chi^2 + n}}$$

$$w = \sqrt{\sum_{i=0}^{\infty} \frac{(P_{1i} - P_{0i})^2}{P_{0i}}}$$
 EFFECT SIZE INDEX w

$$C = \sqrt{\frac{\chi^2}{\chi^2 + n}} = \sqrt{\frac{w^2}{w^2 + 1}} \qquad w = \sqrt{\frac{C^2}{1 - C^2}}$$

MAGNITUDE	VALUE
SMALL	w = 0.1
MEDIUM	w = 0.3
LARGE	w = 0.5



#### **ESTIMATING DIVERSENESS**

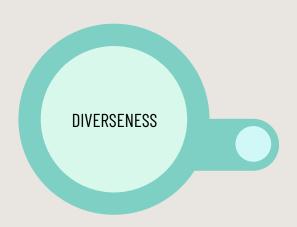
DIVERSENESS
PROVIDES THE TRAINING DIVERSIFICATION
PROBABILITY IN RESPECT TO EACH LEVEL
OF THE PROTECTED ATTRIBUTE AND THE
TARGET VARIABLE

DIVERSENESS



#### **ESTIMATING DIVERSENESS**





$$P = \frac{number\ of\ favorable\ cases}{number\ of\ possibles\ cases}$$



#### **ESTIMATING DIVERSENESS**





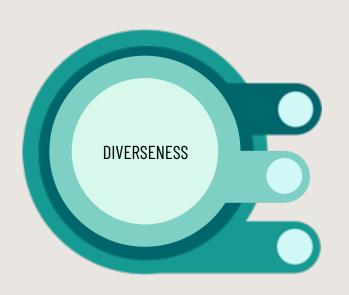
$$P = \frac{number\ of\ favorable\ cases}{number\ of\ possibles\ cases}$$

$$P = \frac{number\ of\ favorable\ properties}{number\ of\ possibles\ properties}$$



#### **ESTIMATING DIVERSENESS**





$$P = \frac{number\ of\ favorable\ cases}{number\ of\ possibles\ cases}$$

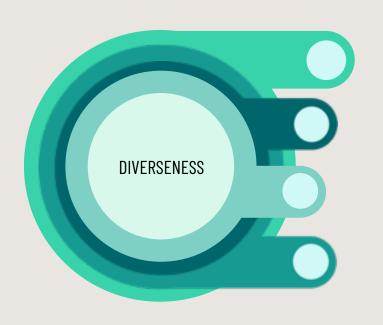
$$P = \frac{number\ of\ favorable\ properties}{number\ of\ possibles\ properties}$$

$$P = (Y = y)$$
  
 $P = (A = a)$  PRIOR PROBABILITIES



#### **ESTIMATING DIVERSENESS**

DIVERSENESS
PROVIDES THE TRAINING DIVERSIFICATION
PROBABILITY IN RESPECT TO EACH LEVEL
OF THE PROTECTED ATTRIBUTE AND THE
TARGET VARIABLE



$$P = \frac{number\ of\ favorable\ cases}{number\ of\ possibles\ cases}$$

$$P = \frac{number\ of\ favorable\ properties}{number\ of\ possibles\ properties}$$

$$P = (Y = y)$$
  
 $P = (A = a)$  PRIOR PROBABILITIES

FORMULA	PROBABILITY
P(Y=0) $P(Y=1)$	P = 0.48 P = 0.52
P(A = white) P(A = black) P(A = Asian)	P = 0.6 P = 0.35 P = 0.15



**DATA BIAS AWARENESS** 

INCLUSIVENESS
PROVIDES THE PROBABILITY THAT TWO
PROPERTIES ARE SIMULTANEOUSLY
INCLUDED IN THE TRAINING SET

INLUSIVENESS

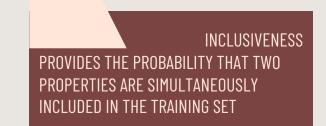


INCLUSIVENESS
PROVIDES THE PROBABILITY THAT TWO
PROPERTIES ARE SIMULTANEOUSLY
INCLUDED IN THE TRAINING SET

P(A|B) POSTERIOR PROBABILITIES





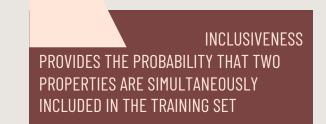




P(A|B) POSTERIOR PROBABILITIES

$$P(A=a\cap Y=y)=P(A=a)P(Y=y|A=a)$$
  
 
$$P(Y=y\cap A=a)=P(Y=y)P(A=a|Y=y)$$





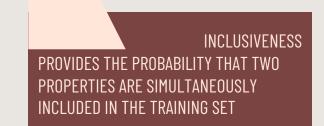


$$P(A|B)$$
 POSTERIOR PROBABILITIES

$$P(A=a\cap Y=y)=P(A=a)P(Y=y|A=a)$$
  
 
$$P(Y=y\cap A=a)=P(Y=y)P(A=a|Y=y)$$

$$P\left(A=a\cap Y=y\right) = P\left(Y=y\cap A=a\right)$$
 COMPOUND PROBABILITY THEOREM







P(A|B) POSTERIOR PROBABILITIES

$$P(A=a\cap Y=y)=P(A=a)P(Y=y|A=a)$$
  
 
$$P(Y=y\cap A=a)=P(Y=y)P(A=a|Y=y)$$

$$P(A = a \cap Y = y) = P(Y = y \cap A = a)$$

COMPOUND PROBABILITY THEOREM

FORMULA	PROBABILITY
$P(Y=0 \cap A=white)$	P=0.42
$P(Y=0 \cap A=black)$	P=0.07
$P(Y=0 \cap A=Asian)$	P=0.09
$P(Y=1 \cap A=white)$	P=0.18
$P(Y=1\cap A=black)$	P=0.28
$P(Y=1 \cap A=Asign)$	P=0.06

**DATA BIAS AWARENESS** 



TRAINING LIKELIHOOD
PROVIDES THE OCCURRENCE LIKELIHOOD
OF THE PROTECTED ATTRIBUTE LEVELS
GIVEN THE TARGET VARIABLE LEVELS AND VICE VERSA - BEFORE THE TRAINING
SET IS SAMPLED

TRAINING LIKELIHOOD





$$P(A = a|Y = y) = \frac{P(A = a)P(Y = y|A = a)}{P(Y = y)}$$

$$P(Y = y|A = a) = \frac{P(Y = y)P(A = a|Y = y)}{P(A = a)}$$





$$P(A = a|Y = y) = \frac{P(A = a)P(Y = y|A = a)}{P(Y = y)}$$

$$P(Y = y|A = a) = \frac{P(Y = y)P(A = a|Y = y)}{P(A = a)}$$

$$\Omega: \bigcup_{i=1}^{N} Y_i = \Omega, \text{ hence } \sum_{i=1}^{N} P(Y_i) = P(\bigcup_{i=1}^{N} Y_i)$$

$$\Omega: \bigcup_{i=1}^{N} A_i = \Omega, \text{ hence } \sum_{i=1}^{N} P(A_i) = P(\bigcup_{i=1}^{N} A_i)$$



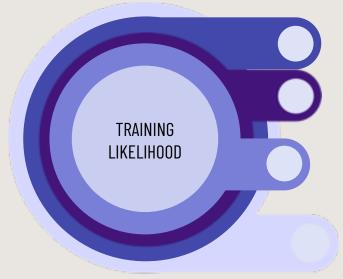
$$P(A = a|Y = y) = \frac{P(A = a)P(Y = y|A = a)}{P(Y = y)}$$

$$P(Y = y|A = a) = \frac{P(Y = y)P(A = a|Y = y)}{P(A = a)}$$

$$\Omega: \bigcup_{i=1}^{N} Y_i = \Omega, \text{ hence } \sum_{i=1}^{N} P(Y_i) = P(\bigcup_{i=1}^{N} Y_i)$$

$$\Omega: \bigcup_{i=1}^{N} A_i = \Omega, \text{ hence } \sum_{i=1}^{N} P(A_i) = P(\bigcup_{i=1}^{N} A_i)$$

$$P(Y = y|A) = \frac{P(Y = y)P(A|Y = y)}{P(A)} = \frac{P(Y = y)P(A|Y = y)}{\sum_{i=1}^{N} P(A|Y_i)P(Y_i)}$$



FORMULA	PROBABILITY
P(Y=0 A=white)	P=0.7
P(Y=0 A=black)	P=0.2
P(Y=0 A=Asian)	P=0.6
P(Y=1 A=white) $P(Y=1 A=black)$ $P(Y=1 A=Asian)$	P=0.3 P=0.8 P=0.4
P(A=white Y=1)	P=0.34
P(A=white Y=0)	P=0.87
P(A=black Y=1)	P=0.53
P(A=black Y=0)	P=0.15
P(A=Asian Y=1)	P=0.11
P(A=Asian Y=0)	P=0.18



# FAIRNESS IN RANKING SYSTEMS

#### **APPENDIX**

MORAL GROUND

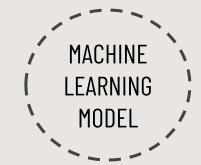
AFteRS: ALGORITHM 1

AFteRS: ALGORITHM 2

TABLE OF METRICS















CASE 1





TARGET: PAST PERFORMANCE

MACHINE GROUND TRUTH
LEARNING
MODEL

TRAINING DATA







TARGET: PAST PERFORMANCE

REPLICATE HUMAN DECISIONS

CASE 1





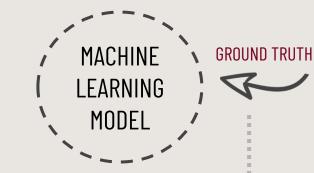


TARGET: PAST PERFORMANCE

REPLICATE HUMAN DECISIONS

CASE 1

CASE 2



TRAINING DATA





TARGET: PAST PERFORMANCE

REPLICATE HUMAN **DECISIONS** 

CASE 1

**MACHINE GROUND TRUTH LEARNING** MODEL

TARGET: ANNUAL REVIEWS

CASE 2



**TRAINING** 

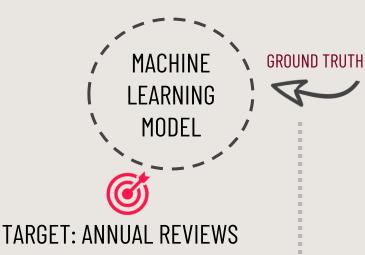
DATA



TARGET: PAST PERFORMANCE

REPLICATE HUMAN DECISIONS

CASE 1



REPLICATE HUMAN MANAGER DECISIONS

CASE 2



**TRAINING** 

DATA

FAIRNESS IN RANKING SYSTEMS



TARGET: PAST PERFORMANCE

REPLICATE HUMAN **DECISIONS** 

CASE 1



TARGET: ANNUAL REVIEWS

REPLICATE HUMAN MANAGER DECISIONS

CASE 2



**TRAINING** 

DATA

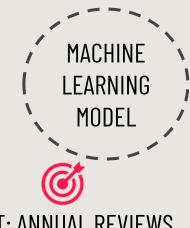
**GROUND TRUTH** 



TARGET: PAST PERFORMANCE



CASE 1



TARGET: ANNUAL REVIEWS

REPLICATE HUMAN MANAGER DECISIONS

CASE 2



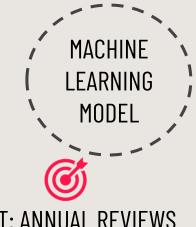


TARGET: SALES FOR THE YEAR





TARGET: PAST PERFORMANCE



TARGET: ANNUAL REVIEWS





REPLICATE HUMAN DECISIONS

CASE 1 CASE 2

REPLICATE HUMAN MANAGER DECISIONS

**REPLICATE CIRCUMSTANCES** 



# AFteRS: Algorithm 1

```
Algorithm 1 Automated Fair Distributive Ranking – Step 1-2 (Figure 3)
The algorithm partitions the population in n types (Section 3.2.1), derives effort (Section 3.2.2), and computes the Standardized Distribution (Equation 5)
```

```
Step 1
```

input: dataset D

output: non-overlapping subsets of D  $\implies$  population partitioned in  $T_k$  types

- 1: for all  $X_i \in D$  do
- Test the null hypothesis of independence between Y and all X<sub>i</sub>
- if  $H^0_{nortial}: P(Y|X) = P(Y)$  couldn't be rejected then
- Stop
- i: else
- select X<sub>i</sub> with the strongest association to Y (smallest adj p-value)
   find the splitting point C\* for X<sub>i</sub> such that
  - S: S<sub>n</sub><sup>x<sub>i</sub></sup> ⊂ χ<sub>i</sub> are all the possible disjoint sets of the sample space χ<sub>i</sub>
- 9: end if
- 10: end for
- 11: return  $T_k$  vectors  $\subset D$

#### Step 2

input:  $T_k$  vectors  $\subset$  D output: Standardized Outcome  $\tilde{y}_i^t(\lambda)$ 

- 1: partition each  $T_k$  in 10 sets  $\Psi_n$ , such that  $\Psi_{k,n} \subset T_k$
- 2: training set  $1^{st} 9^{th} T_k$  sets
- 3: test set 10<sup>th</sup> T<sub>k</sub> set
- 4: for all  $\Psi_{k,n} \subset T_k$  do
  - 1. perform the Bernstein polynomials log-likelihood on the training set to estimate
  - the best type-distribution approximation  $LL_B(p_m = \sum_{i=1}^n \log f_B(x_j, p_m))$ 2. predict the CDF of  $T_k$  on the test set
- 8: end for
- estimate the Standardize Distribution = y<sup>t</sup><sub>i</sub>(λ) μ/μ<sup>λ</sup>
   return ŷ<sup>t</sup><sub>i</sub>(λ)

# AFteRS: Algorithm 2

```
Algorithm 2 Automated Fair Distributive Ranking – Step 3 (Figure 3)
The algorithm computes the \Gamma ranking based on policies Equity, Equality and Need
Step 3
input: Standardized Outcome \tilde{y}_{i}^{t}(\lambda)
output: ranking \Gamma constrained by a policy \theta \in \Theta
 1: if \theta = equity then
 2: for all Y_{t,\lambda} \in D do
        compute the counterfactual outcome from stnd. outcome and decomposed Gini
        \Gamma \leftarrow ranking ordered by counterfactual outcome
 5: end for
 7: if \theta = equality then
 8: for all T_k \in D do
        sorted_{T_b} \leftarrow type-ranking ordered by decreasing stnd. outcome
10: end for
11: for all (j) \in sorted_{T_k} do
        row_n \leftarrow j element of sorted_{T_k}
        array[i] \leftarrow row_n ordered by decreasing std. outcome
14: end for
15: Γ ← merge all j array
17: if \theta = need then
18: G_k \leftarrow n subsets \in D grouped by protected attribute A
19: for all (z) ∈ sorted<sub>Gk</sub> do
        row_n \leftarrow z element of sorted_{G_k}
        \operatorname{array}[\mathbf{z}] \leftarrow row_n ordered by decreasing std. outcome
22: end for
23: \Gamma \leftarrow merge all z array
24:
```

25: return Ranking Γ

#### TABLE OF METRICS

Metric	Formula	Input
Expected ranking	$r = argmaxU(ranking_n q)$	Score distribution
Exposure	$\frac{1}{\log(1+j)}$	Original Distribution
Relevance	$\beta(Rel(item_n user_n, q))$	ScDistr., adj ScoreDistr
Expected ranking-policy	$\Gamma = argmax_{\theta \in \Theta}u^{t}(q e_{i}(\lambda), \theta)$	Adj ScoreDistr
Exposure-policy	$\max_{\theta \in \Theta} \int_{0}^{1} min_{t}exp^{t}(\lambda, \theta) d\lambda$	Adj ScoreDistr
Gini Index	$1 - \frac{1}{a} \int_0^\infty (1 - F(y))^2 dy$	All distributions
Decomposed Gini	$Gini_{\lambda}^{t}$	Stand. score distribution
Richness	$n^t$	Types diversity
Margalef	$\frac{T-1}{\ln N}$	Types diversity
Shannon-Wiener Index	$H = \sum_{i=1}^{R} p_i \ln p_i$	Types diversity
Simpson	$1 - \sum \frac{n^t(n^t-1)}{N(N-1)}$	Types diversity
Theil Index	$\frac{1}{N}\sum_{i=1}^{N} ln(\frac{\mu}{u_i})$	All Distributions
Opportunity-Types Profile	$min/max(y^t - \mu(y))$	Score distribution
Opportunity-Types Rate	$y^t - \mu(y)$	ScDistr
Opportunity-L/G Profile	$min/max(y_{\lambda}^{t} - \mu(y_{\lambda}))$	Stnd. distribution
Opportunity-L/G Rate	$y_{\lambda}^{t} - \mu(y_{\lambda})$	StndDistr
Unexplained Inequality Rate	$\frac{1}{N}\sum y_i - \tilde{y}_i$	ScDistr, stndDistr
Reward Profile	$min/max(j(y_{\lambda}^{t}) - j(adj(\tilde{y}_{\lambda}^{t})))$	ScDistr, adj score distr.
Reward Rate	$j(y_{\lambda}^{t}) - j(adj(\tilde{y}_{\lambda}^{t}))$	ScDistr, adj score distr.

Table 3: Summary of metrics employed. Notation: F(y)= cumulative distribution function of the score,  $\mu$  = mean score; R = number of types,  $p_i$  = frequency of types;  $y_{\lambda}^t$  = score distribution aggregated by type and quantile;  $\tilde{y}_i$ = standardized score;  $adj(\tilde{y}_{\lambda}^t)$ = adjusted mean-type score at each effort degree (after policy); j = ranking position



## LONG-TERM FAIRNESS

**APPENDIX** 

APPLICATION SETTING INDIVIDUAL DYNAMICS



#### **APPLICATION SETTING**

SCENARIO	UNIVERSITY SELECTION PROCESS	MAXIMIZE LONG-TERM SELECTION	BEST CANDIDATES SELECTION	FAIRNESS AS POSITIVE BEHAVIOR	QUALIFICATION AT TIME t+1	
DATA	SYNTHETHIC DATA	RETAKING SAT STATISTICS	GPA, SAT AND GRE SCORE, AGE, SEX	ONLY IMPROVE SAT SCORE	ESTABLISH Maximum Score	
POLICY	GRADIENT BOOSTING MACHINE	GENERALIZED LINEAR MODEL	 K-NEAREST NEIGHBOUR	NAIVE Bayes Classifier	SUPPORT Vector Machine	-
TIME	NON EVOLVING PREFERENCES	DATA AT TIME t+1	KNOWLEDGE OF THE STATE OF THE NATURE	NOT A DETERMINISTIC MODEL	POLICY SELECTION 10 YEARS	



**ALTERNATIVES** 

**SCENARIOS** 

IMPACTS AND UTILITY FUNCTION

**PREFERENCES** 

 $x_1$ : APPLYING WITH QUALIFICATION

 $x_2$ : APPLYING WITHOUT

QUALIFICATION

 $x_3$ : NOT APPLYING

 $x_1: 0$ 

such that X = [0, N0, 0]

 $x_2$ : NQ

 $x_3$ : N



**ALTERNATIVES** 



IMPACTS AND UTILITY FUNCTION

**PREFERENCES** 

$\pi(\omega x)$	Q	NQ	N
OPTIMISTIC	0.75   0.86	0.75   0.76	0.75   0.76
PESSIMISTIC	0.85   0.86	0.85   0.76	0.85   0.76
AGNOSTIC	0.8   0.86	0.8   0.76	0.8   0.76



**ALTERNATIVES** 

**SCENARIOS** 

IMPACTS AND UTILITY FUNCTION

**PREFERENCES** 

f(ω x)	Q	NQ	N
OPTIMISTIC	0.9	1	-1
PESSIMISTIC	0.9	0	0
AGNOSTIC	0.9	0	0



**ALTERNATIVES** 

**SCENARIOS** 

IMPACTS AND UTILITY FUNCTION



f(ω x)	Q	NQ	N
OPTIMISTIC	0.9	1	-1
PESSIMISTIC	0.9	0	0
AGNOSTIC	0.9	0	0

$$\max_{x \in X} Laplace(x) = \max_{x \in X} \frac{\sum_{\omega \in \Omega} f(x, \omega)}{|\Omega|}$$



**ALTERNATIVES** 

**SCENARIOS** 

IMPACTS AND UTILITY FUNCTION

PREFERENCES	
7	

$f(\omega x)$	Q	NQ	N
OPTIMISTIC	0.9	1	-1
PESSIMISTIC	0.9	0	0
AGNOSTIC	0.9	0	0

Laplace(x)	0.9	0.33	- 0.33

