

D6.3

Social, Ethical and Legal Al-Fairness Methodologies

This content is based on research conducted for the AEQUITAS report "Preliminary Social, Ethical and Legal AI-Fairness Methodologies"

www.aequitas-project.eu

01.

Social, Ethical and Legal Al-Fairness Methodologies

Our Goal To develop a methodology

for the creation of the social, ethical and legal requirements for the 3 AEQUITAS Engines:

> Reparation and Mitigation Engine

Awareness and Diagnostics Engine ADE

Fairness-by-Design Engine







What do we need to address?

- What legal elements should be considered.
- What ethical elements should be considered, based on the Ethics Guidelines of Trustworthy AI (EGTAI).
- What social elements should be considered, based on existing literature, methodologies, and processes.

The **preliminary table** shown below aims to **highlight the requirements** that must be considered for each of the **AEQUITAS Engines** at the **data, model and interpretation levels**.

	Awareness and Diagnostics Engine ADE	Reparation and Mitigation Engine RME	Fairness-by-Design Engine FDE
Data	Legal : Article 10 of the Al Act, specifically, those under art. 10.2, art. 10.3, and art. 10.4.	Technical reparation Stakeholder engagement + Re-diagnosis (e.g. re-evaluating art. 10.3)	Fair data collection methodology
Model	Ethical : ALTAI/EGTAI Requirement 4 - Transparency	Technical reparation Stakeholder engagement + Re-diagnosis (e.g. re-evaluating art. 10.3)	Al Act Risk Classification
Interpretation	Legal : Article 14 of the Al Act	Technical reparation Stakeholder engagement + Re-diagnosis (e.g. re-evaluating art. 10.3)	Stakeholder Engagement

User Stories

We aim to develop the requirements needed for the detection of bias in ethical, social, legal, economic and cultural contexts. More specifically, we aim to develop user stories that can diagnose and assess fairness in Al-systems at the data, model, and interpretation level.

Awareness and Diagnostics Engine ADE				
TITLE Data compliance	PRIORITY XXXXXXXXXX	ID XXXXXXXXXX	PROPOSER XXXXXXXXXX	
** As an AI developer (researcher, engineer, data scientists) I want to be able to check if the dataset is fair as regards the persons or groups on which the system is intended to be used (art. 10.3 & 10.4 AIA). So that I can evaluate the level of compliance with the AIA and act accordingly.				
ACCEPTANCE CRITERIA Notes Elements to be identified/visualized for the dataset as regards the person or groups on which the system is intended to be used:				
Relevancy	Representativeness	Free of errors (as far as possible)	Completeness	
Appropriateness of Statistical properties	Characteristics of Geographical setting	Characteristics of Behavioural setting	Characteristics of Functional setting	

		aration and ation Engine RME	
TITLE Data reparation	PRIORITY XXXXXXXXXX	ID XXXXXXXXXX	PROPOSER XXXXXXXXXX
		gineer, data scientists ata as expressed in us	
ACCEPTANCE CRITER	combination wi	method should include tee th the re-involvement of r ial and legal aspects are s	

03.

	En	-by-Design gine DE		
identify requirements which we refer to a ligned to the AI	other engines, the Fair ents that can be trans as 'building blocks'. In ifecycle so AI develop rness-by-design met	s lated into too n turn, these bu pers or other us	Is and meth ailding block sers of the A	odologies s can be
TITLE Data reparation	PRIORITY XXXXXXXXXX	ID XXXXXXXXXXX		OPOSER XXXXXXXXX
As an	l want		So that	
Al developer	a visual rep of the Al life	ecycle.		e a smooth n with the nodoloay
	Analysis, Developr (Testing, Deploym Decommissioning ded to be implemente In engine :	ent), Monitoring ir		
Fairness-by-Desig				
The need to have a trustworthy Al elaboration process	The need to use EGTAI/ALTAI	The need to identify and involve multidiscipl expertise	l id in	ne need to entify and volve akeholders



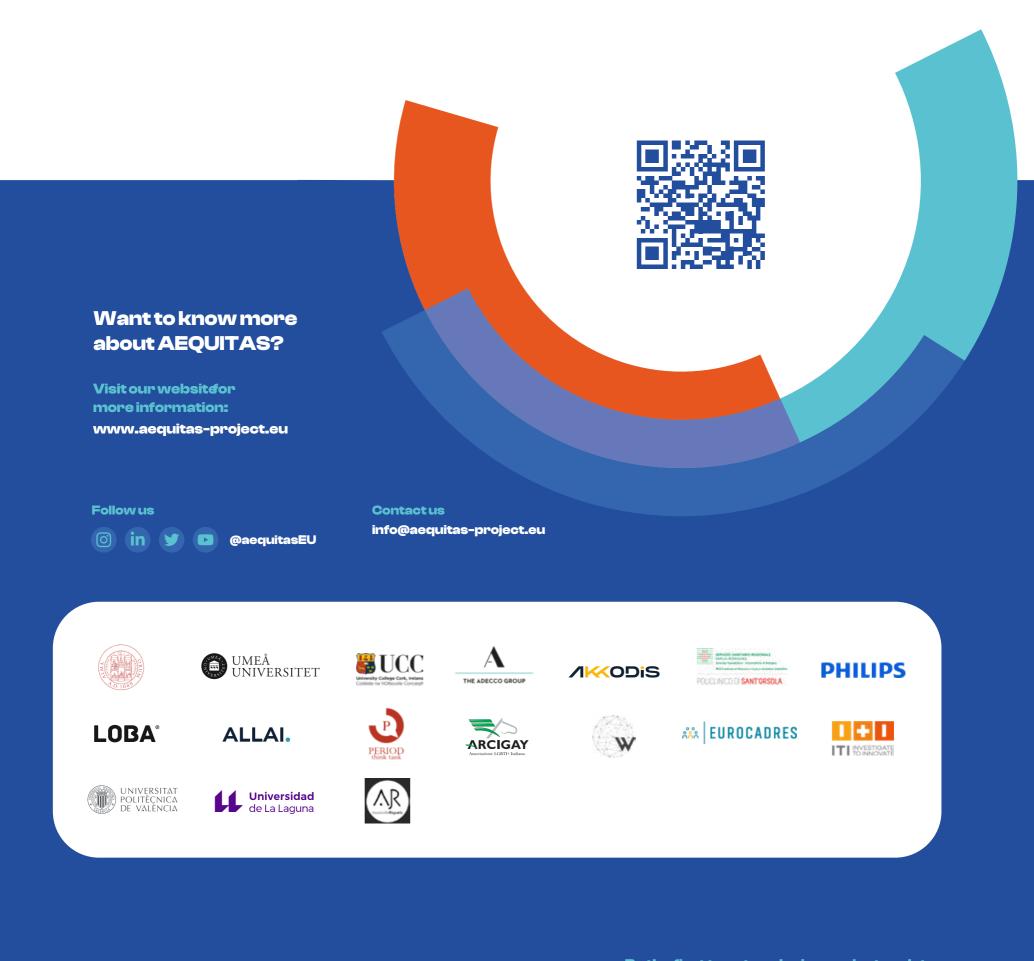
Al fairness Impact Assessment/Al fairness Readiness Assessment Ethics Guidelines for Trustworthy AI and Assessment List for Trustworthy AI (HLEG AI) Socio-technical Matrix

Stakeholder Identification Methodology (developed as Stakeholder Engagement Methodology Trustworthy AI Deliberation (based on the 7-step

Methodology (developed as part of Deliverable 6.1)	Methodology	(based on the 7-step exercise for Trustworthy Al, developed for the Trustworthy Al Project [Erasmus+])
Al Fairness Regulatory Landscape Identification and Assessment	AI Act Risk Classification	Al Act High-Risk Requirements Guidance for Fairness
Fundamental Rights Impact Assessment	Fair software engineering methodologies architectures and methods:	 Fair System Architecture Methodology Fairness Criteria Definition, Assessment and Monitoring Fair Data Collection Methodology
Al fairness Evaluation/Bias Audit	Monitoring in Operation through Critical Control Points	

The requirements identified and observed will be used as building blockfor the 3 AEQUITAS ENGINES.

What's next?



Views and opinions expressed are however those of the authors only and do not necessarily reflect those of the European Union. Neither the European Union nor the granting authority can be held responsible for them. Be the first to get exclusive project updates by subscribing toour newsletter!