

A FROM VALUES TO SPECIFIC CRITERIA MANIFESTATIONS

Value	Criteria	Manifestations
Privacy	(1) Consent for data usage [3, 56, 60]	• Written declaration of consent [56]
	(2) Data protection [3, 60, 61]	• Description of what data is collected [125]
	(3) Control over data / ability to restrict processing [56, 60]	• Description of how data is handled [125]
	(4) Right to rectification [3, 56, 60]	• Purpose statement of data collection [125]
	(5) Right to erase the data [3, 56, 60]	• Statement of how long the data is kept [125]
	(6) Right of access by data subject, data agency [56, 167]	• Form and submission mechanisms to object data collection and to make complaints [27]
Conservation		• Obfuscation of data [3]
		AGAINST INTEGRITY THREATS [183]:
		• Training time [183] Ex.:
		• Data sanitization ¹² [23, 40]
		• Robust learning ¹³ [23, 73]
		• Prediction time [183]
Security		• Model enhancement [23, 74, 122, 141] Ex.:
		• Adversarial Learning ¹⁴
		• Gradient masking ¹⁵
		• Defensive Distillation ¹⁶
	(1) Resilience to attacks : protection of privacy [86, 127, 178], vulnerabilities, fallback plans [3, 60, 75, 133]	AGAINST PRIVACY THREATS [183]:
	(2) Predictability [3, 57, 60]	• Mitigation techniques [136]:
	(3) Robustness / reliability : prevent manipulation [3]	• Restrict prediction vector to top k classes ¹⁷ [161]
		• Coarsen the precision of the prediction vector ¹⁸ [161]
		• Increase entropy of the prediction vector ¹⁹ [161]
		• Use regularization ²⁰ [101, 161]
		• Differential privacy mechanisms [136]:
		• Differential privacy ²¹ [53, 187]. Ex.:
		• Adversarial regularization ²² [136]
		• MemGuard ²³ [97]

¹²It ensures data soundness by identifying abnormal input samples and by removing them [183].¹³It ensures that algorithms are trained on statistically robust datasets, with little sensitivity to outliers [183].¹⁴Adversarial samples are introduced to the training set [183].¹⁵Input gradients are modified to enhance model robustness [183].¹⁶The dimensionality of the network is reduced [183].¹⁷Applicable when the number of classes is very large. Even if the model only outputs the most likely k classes, it will still be useful [161].¹⁸It consists in rounding the classification probabilities down [161].¹⁹Modification of the softmax layer (in neural networks) to increase its normalizing temperature [161].²⁰Technique to avoid overfitting in ML that penalizes large parameters by adding a regularization factor λ to the loss function [161].²¹It prevents any adversary from distinguishing the predictions of a model when its training dataset is used compared to when other dataset is used [187]²²Membership privacy is modeled as a min-max optimization problem, where a model is trained to achieve minimum loss of accuracy and maximum robustness against the strongest inference attack [136].²³Noise is added to the confidence vector of the attacker so as to mislead the attacker's classifier [97]

	Value	Criteria	Manifestations
Conservation	Performance	(1) Correctness of predictions [26, 57, 60] (2) Memory efficiency [3, 26] (3) Training efficiency [26] (4) Energy efficiency [3, 26] (5) Data efficiency [26]	<ul style="list-style-type: none"> • Accuracy (for classification, sum of true positive and true negative rates) [130, 180] • False Positive and False Negative rates [130, 180] • False Discovery and Omission Rate [130] • Mean and median error [180] • R2 score [25] • Precision and recall rates [180] • Area under ROC curve (AUC) [25] • Estimation of energy consumption through [68]: <ul style="list-style-type: none"> • performance counters • simulation • instruction- or architecture-level estimations • real-time estimation • Estimation of GPU memory consumption [67, 123] • Wall-clock training time [14, 41]
		(1) Desirability of technology [1, 34, 104] (2) Benefit to society [60–62, 133] (3) Environmental impact [3, 21]	<ul style="list-style-type: none"> • Diverse and inclusive forum for discussion [60, 129] • Measure of social and environmental impact [21, 133, 147]
Universalism	Fairness	(1) Individual fairness ²⁴ [18, 52, 110, 126] (2) Demographic parity ²⁵ [18, 52, 80, 86, 102, 110, 126, 163, 177] (3) Conditional Statistical parity ²⁶ [126, 177] (4) Equality of opportunity ²⁷ [79, 126, 175] (5) Equalized odds ²⁸ [126] (6) Treatment equality ²⁹ [22, 126] (7) Test fairness ³⁰ [37, 126, 177] (8) Procedural fairness ³¹ [77, 110, 126]	<ul style="list-style-type: none"> • Accuracy across groups (for classification, sum of true positive and true negative rates) [37, 80, 105, 133] • False positive and negative rates across groups [37, 105, 126, 151, 179] • False discovery and omission rates across groups [130, 151] • Pinned AUC [48, 130] • Debiasing algorithms [19] • Election of protected classes based on user considerations [77]
	Non-discrimination	(1) Quality and integrity of data [60, 70, 86, 133, 144] (2) Inclusiveness in design [57, 60, 133] (3) Accessibility [3, 26, 60, 133]	<ul style="list-style-type: none"> • Inclusive data generation process [3, 34, 70, 133] • Analysis of data for potential biases, data quality assessment [3, 60, 69, 86, 126] • Diversity of participant in development process [3, 60, 114, 189] • Access to code and technology to all [3, 26, 60, 133]

²⁴Similar individuals should be treated in a similar way. Diverging definitions state that: two individuals that are similar with respect to a common metric should receive the same outcome (*fairness through awareness*); or any protected attribute should not be used when making a decision (*fairness through unawareness*); or the outcome obtained by an individual should be the same if this individual belonged to a counterfactual world or group (*counterfactual fairness*) [126].

²⁵The probability of getting a positive outcome should be the same whether the individual belongs to a protected group or not [126].

²⁶Given a set of factors L, individuals belonging to the protected or unprotected group should have the same probability of getting a positive outcome [126].

²⁷The probability for a person from class A (positive class) of getting a positive outcome, which should be the same regardless of the group (protected group or not) that the individual belongs to [126].

²⁸The probability for a person from class A (positive class) of getting a positive outcome and the probability for a person from class B (negative class) of getting a negative outcome should be the same [126].

²⁹The ratio of false positives and negatives has to be the same for both groups [126].

³⁰For any probability score S, the probability of correctly belonging to the positive class should be the same for both the protected and unprotected group [126].

³¹It deals with the fairness of the decision-making process that leads to the outcome in question [77].

Value	Criteria	Manifestations
Openness	Transparency	<ul style="list-style-type: none"> • Description of data generation process [3, 20, 34, 69, 70, 133] • Disclosure of origin and properties of models and data [3, 130, 168] • Open access to data and algorithm [3, 26, 60, 168] • Notification of usage/interaction [60] • Regular reporting [60]
	Explainability	<ul style="list-style-type: none"> • Interpretability by design [18] • Post-hoc explanations [18]
Individual empowerment	Contestability	<ul style="list-style-type: none"> • Information of who determines and what constitutes a contestable decision and who is accountable [121] • Determination of who can contest the decision (subject or representative) [121] • Indication of type of review in place [121] • Information regarding the contestability workflow [121] • Mechanisms for users to ask questions and record disagreements with system behavior [87, 131]
	Human Control	<ul style="list-style-type: none"> • Continuous monitoring of system to intervene [57, 60, 166] • Establishment levels of human discretion during the use of the system [57, 127] • Ability to override the decision made by a system [57]
	Human agency	<ul style="list-style-type: none"> • Give knowledge and tools to comprehend and interact with AI system [57] • Opportunity to self-assess the system [57]

Table 3. Summary of the specific criteria that relate to each value considered in our ML assessment framework. These criteria are then translated into specific manifestations in the form of signifiers (orange), process-oriented practices (olive) or quantifiable indicators (magenta).

B MAPPING STAKEHOLDERS

Stakeholder	Mapping [164]	Nature of knowledge	Purpose of insight
Development team	ML, Formal + Instrumental + Personal	<ul style="list-style-type: none"> • “Knowledge of the math behind the architecture” [164] • “Stakeholder involved in an ex-ante impact assessment of the automatic decision system” [84] 	<ul style="list-style-type: none"> • Ensure/improve product efficiency and debug [18] • Research new functionalities [18]
Auditing team	Milieu, Formal + Instrumental	<ul style="list-style-type: none"> • “Familiarity with broader ML-enabled systems” [164] • “Experts who intervene wither upstream or downstream” [84] 	<ul style="list-style-type: none"> • Verify model compliance with legislation [18]
Data domain experts	Data domain, Formal + Instrumental	<ul style="list-style-type: none"> • “Theories relevant to the data domain” [164] • “Professional involved in the operational phase of the automatic decision system” [84] 	<ul style="list-style-type: none"> • Gain scientific or domain-specific knowledge [18, 164] • Trust the model [18, 164] • Act based on the output [164]
Decision subjects	Data domain + Milieu, Personal	<ul style="list-style-type: none"> • “Lived experience and cultural knowledge” [164] • “Layperson affected by the outcomes of the automatic decision system” [84] 	<ul style="list-style-type: none"> • Understand their situation [18] • Verify fair decision [18] • Contest decision [164] • Understand how one’s data is being used [164]

Table 4. Description of potential stakeholders that can be brought together as part of our value-based framework. These stakeholders have been mapped following the two dimensional criteria (type of knowledge —formal, instrumental or personal— and contexts in which this knowledge manifests —ML, data domain, milieu—) outlined by Suresh et al. [164]. The nature of their knowledge and the purpose of gaining insight for each of them have also been defined.

A TAILORED COMMUNICATION OF SYSTEM-RELATED INFORMATION

		Development team	Auditing team	Data Domain experts	Decision subjects
Conservation	Privacy	[K]	[K]		[A] [B]
	Security	[K] [W] [AB]	[K] [W]		
	Performance	[F] [G] [H] [Y] [Z] [AE]	[G] [H] [Y] [Z] [AE]	[I] [J]	[J]
Universalism	Respect for public interest	[E] [AE]	[E] [AE]	[E]	[C] [D]
	Fairness	[G] [H] [K] [W] [X] [Y] [Z] [AD]	[G] [H] [K] [W] [X] [Y] [Z] [AD]	[I] [J]	[J]
	Non-discrimination	[H] [K] [X] [Y] [AD]	[H][K] [X] [Y] [AD]	[J] [L]	[J] [L]
Openness	Transparency	[H] [K] [M]	[H][K] [M]	[I] [J] [L] [M]	[B] [J] [L] [M]
	Explainability	[M] [N] [O] [Q] [AC] [AD] [P]	[M] [N] [O] [Q] [AC] [AD] [P]	[J] [M] [N] [O] [Q] [P]	[J] [M] [N] [O] [Q] [R] [S] [P]
Individual empowerment	Contestability	[U]	[U]	[T] [U]	[T] [AF]
	Human Control	[V]	[V]	[T] [V]	[C] [T] [V]
	Human Agency			[T]	[T] [B] [AA]

Table 14. Mapping of available means for transmitting value-specific manifestations to different stakeholders based on the purpose of their insight and the nature of their knowledge. These means have been classified into three main categories: descriptive documents specifying whether/how a value manifestation is fulfilled (red), strategies for fulfilling value manifestations (blue), and complete tools for enabling the fulfillment of value manifestations (green). This table aims at facilitating the navigation of table 15, where each means is documented.

Means	Value	Manifestation(s)	Stakeholder			Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS			
Iconsets for data privacy declarations [A] [56, 89, 125, 149]	Privacy	<ul style="list-style-type: none"> Description of what data is collected Description of how data is handled Purpose statement of data collection Statement of how long the data is kept 				✓	Agnostic	Iconsets	
Privacy dashboards [B] [54, 58, 59, 85, 191]	Privacy	<ul style="list-style-type: none"> Description of what data is collected Description of how data is handled Purpose statement of data collection 				✓	Agnostic	<ul style="list-style-type: none"> Timelines Bar charts Maps Network graphs 	
	Human agency	<ul style="list-style-type: none"> Opportunity to self-assess the system 							
Risk matrix [C] [3, 107]	Transparency	<ul style="list-style-type: none"> Disclosure of origin and properties of data 							
	Respect for public interest	<ul style="list-style-type: none"> Measure of social impact 							<ul style="list-style-type: none"> Two dimensional space (vulnerability vs dependence of the decision)
Moral space [D] [86]	Human Control	<ul style="list-style-type: none"> Ability to override the decision made by a system 				✓	Agnostic		
	Respect for public interest	<ul style="list-style-type: none"> Measure of social impact 				✓	Agnostic	Based on human judgement	<ul style="list-style-type: none"> Three dimensional moral space. Wrongness as a function of intention and harm

Means	Value	Manifestation(s)	Stakeholder			Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS			
[E] Social impact assessment [147]	Respect for public interest	<ul style="list-style-type: none"> Measure of social impact 	✓	✓	✓		Agnostic	Anticipate scenarios	
									<ul style="list-style-type: none"> Summary statistics Confusion matrices Labels chart Precision-recall curves Connector lines to identify similar examples in feature space Highlighted boxes for correlations between features and target classes
[F] Model Tracker interactive visualization [9]	Performance	<ul style="list-style-type: none"> Accuracy False Positive and Negative rates 	✓				Classification tasks		
[G] Model cards for models [130]	Performance	<ul style="list-style-type: none"> Accuracy False Positive and Negative rates False Discovery and omission rates 	✓						<ul style="list-style-type: none"> Confidence bars Bar charts
	Fairness	<ul style="list-style-type: none"> Accuracy across groups False Positive and Negative rates across groups False Discovery and omission rates across groups 					Agnostic		

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
What-if tool ³² [H] [180]	Performance	<ul style="list-style-type: none"> • Accuracy • False Positive and Negative Rates • False Discovery and omission rates 							<ul style="list-style-type: none"> • Confusion matrices • (Two-dimensional) Histograms • Scatterplots • Summary statistics of datasets • Partial dependence plots 	Interactive modules include: list of feature values, inference values, and counterfactual controls
	Fairness	<ul style="list-style-type: none"> • Accuracy across groups • False Positive and Negative Rates across groups • False Discovery and omission rates across groups 	✓				Classification tasks, Regression tasks			
	Transparency	<ul style="list-style-type: none"> • Disclosure of origin and properties of data 								
	Non-discrimination	<ul style="list-style-type: none"> • Analysis of data for potential biases, data quality assessment 								
Interactive transfer learning tools [128] [I]	Performance	<ul style="list-style-type: none"> • Accuracy • False Positive and Negative Rates 							<ul style="list-style-type: none"> • Confusion matrices • Z-scored of each filter • Bar charts • Activation heatmaps • t-SNE clusters 	
	Fairness	<ul style="list-style-type: none"> • Accuracy across groups • False Positive and Negative Rates across groups 			✓		Convolutional Neural Networks			
	Transparency	<ul style="list-style-type: none"> • Disclosure of properties of data 								
	Performance	<ul style="list-style-type: none"> • Accuracy 							<ul style="list-style-type: none"> • Summary statistics (percentage scores) for data explanations and performance metrics • Feature importance • Contrastive explanations 	End users were more interested in the limitation of the model: uncertainty
Question-Driven XAI Design [118] [J]	Fairness	<ul style="list-style-type: none"> • Accuracy across groups 								
	Transparency	<ul style="list-style-type: none"> • Disclosure of origin and properties of data 								
	Non-discrimination	<ul style="list-style-type: none"> • Analysis of data for potential biases, data quality assessment 	✓		✓		Agnostic			

³²<https://github.com/pair-code/what-if-tool>

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
Datasheets for [K] datasets [69]	Explain- ability	• Post-hoc explanations								
	Trans- parency	<ul style="list-style-type: none"> • Description of data generation process • Disclosure of origin properties of models and data 								
	Non-discrimi- nation	<ul style="list-style-type: none"> • Analysis of data for potential biases, data quality assessment 								
	Privacy	<ul style="list-style-type: none"> • Written declaration of consent • Description of what data is collected • Description of how data is handled • Purpose statement of data collection • Statement of how long the data is kept 	✓	✓			Agnostic		<ul style="list-style-type: none"> • Summary statistics • Visual examples of datasets (if images, for instance) 	
	Fairness	<ul style="list-style-type: none"> • Election of protected classes 								
Data centric ex- planations [L] [12]	Security	<ul style="list-style-type: none"> • Membership inference 								
	Trans- parency	<ul style="list-style-type: none"> • Description of data generation process • Disclosure of origin and properties of the models and data 			✓	✓	Agnostic		<ul style="list-style-type: none"> • Interactive list • Q&A format • Pie charts • Bar charts • Process diagrams • timelines • Icons 	
	Non-discrimi- nation	<ul style="list-style-type: none"> • Analysis of data for potential biases, data quality assessment 								

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
Example-based explanations [M]	Transparency	• Disclosure of properties of data	✓	✓	✓	✓	Agnostic	• Similar example • Typical example • Counter-factual example	• Example images from dataset if in the visual domain	Normative vs comparative explanations [32]
	Explainability	• Post-hoc explanations								
Explanation by simplification [N]	Explainability	• Post-hoc explanations	✓	✓	✓	✓	Agnostic	• Decision rule • Decision tree		
Feature relevance explanation [O]	Explainability	• Post-hoc explanations	✓	✓	✓	✓	Agnostic	• Feature attribute • Feature shape • Feature interaction • Sensitivity / perturbation - based • Saliency maps (visual domain)	• Bar charts • Visualization of element importance, saliency (visual domain)	Usability of saliency maps for non-experts [7]. They should be accompanied by global descriptors
	Explainability	• Post-hoc explanations	✓	✓	✓	✓	Agnostic			
Contrastive explanations [P]	Explainability	• Post-hoc explanations	✓	✓	✓	✓	Agnostic	• Example of minimum change that leads to different outcomes		
Text-based explanation [Q]	Explainability	• Post-hoc explanations	✓	✓	✓	✓	Agnostic	• With or without outcome comparison		

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
Interactive demonstrations [R] [120]	Explainability	<ul style="list-style-type: none"> Post-hoc explanations 								
						✓	Agnostic			
Experiential AI [S] [82]	Explainability	<ul style="list-style-type: none"> Post-hoc explanations 								<ul style="list-style-type: none"> Art mediated between computer code and human comprehension
						✓	Agnostic			
Interactive contestations [T] [84, 106]	Contestability	<ul style="list-style-type: none"> Mechanisms for users to ask questions and record disagreements with system behavior 								<ul style="list-style-type: none"> Statements restricted to natural language
	Human Control	<ul style="list-style-type: none"> Ability to override the decision made by the system 						✓		
	Human agency	<ul style="list-style-type: none"> Opportunity to assess the system 								
Challenge justifications provided by operator using the same means [U] [84]	Contestability	<ul style="list-style-type: none"> Mechanisms for users to ask questions and record disagreements with system behavior 								<ul style="list-style-type: none"> Further testing Verification
			✓	✓	✓		Agnostic			

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
Mapping of actors and tasks depending on automation level [33]	Human Control	<ul style="list-style-type: none"> Establishment of levels of human discretion during the use of the system 	✓	✓	✓	✓	Agnostic		<ul style="list-style-type: none"> Relationship diagrams 	
	Failure Modes and Effects Analysis [147]	<ul style="list-style-type: none"> Threats against integrity (adversarial learning) and mitigation techniques Accuracy across groups False positives and negatives across groups 	✓	✓			Agnostic			
Aequitas [X] 33 [151]	Fairness	<ul style="list-style-type: none"> Accuracy across groups False positives and negatives across groups 	✓	✓			Agnostic			
	Fairness	<ul style="list-style-type: none"> Accuracy across groups False Positive and Negative rates across groups False Discovery and Omission rates across groups Counterfactual examples 	✓	✓			Agnostic			
	Non-discrimination	<ul style="list-style-type: none"> Analysis of data for potential biases, data quality assessment 								
AI Fairness [Y] 360 34 [19]	Performance	<ul style="list-style-type: none"> False Positive and Negative rates 					Classifiers: logistic regression, random forest classifier and neural networks		<ul style="list-style-type: none"> Bar charts Confidence bars 	
	Fairness	<ul style="list-style-type: none"> False positive and negative rates across groups Debiasing algorithms 	✓	✓						
	Non-discrimination	<ul style="list-style-type: none"> Analysis of data for potential biases, data quality assessment 								

³³<https://github.com/dssg/aequitas>

³⁴<https://github.com/Trusted-AI/AIF360>

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
Fairlearn [Z] ³⁵ [25]	Performance	<ul style="list-style-type: none"> • Accuracy • False Positive and False Negative rates • Precision and recall rates 			✓		Agnostic		<ul style="list-style-type: none"> • Bar charts • Pie charts 	
	Fairness	<ul style="list-style-type: none"> • Accuracy across groups • False negative and false positive rates across groups • Debiasing algorithms 								
Playbook [AA] ³⁶ AI ³⁶ [90]	Human agency	<ul style="list-style-type: none"> • Give knowledge and tools to comprehend and interact with AI systems • Opportunity to self-assess the system 				✓	NLP	Early AI prototyping	<ul style="list-style-type: none"> • Interactive survey 	
Counterfit [AB] ³⁷	Security	<ul style="list-style-type: none"> • Defence against integrity threats • Defence against privacy threats 			✓		Agnostic			
InterpretML [AC] ³⁸ 39 [134, 137]	Explainability	<ul style="list-style-type: none"> • Interpretability by design • Post-hoc explanations 			✓		Both white-box and blackbox models		<ul style="list-style-type: none"> • Bar charts • Line charts • Decision trees 	

³⁵<https://github.com/fairlearn/fairlearn>³⁶<https://github.com/microsoft/HAXPlaybook>³⁷<https://github.com/Azure/counterfit>³⁸<https://github.com/interpretml/interpret/>³⁹<https://github.com/interpretml/DiCE>

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
Error analysis [AD] dashboard ⁴⁰	Non-discrimination	• Analysis of data for potential biases, data quality assessment								
	Explainability	• Post-hoc explanations	✓				Agnostic		• Decision tree • Error heatmap	
	Fairness	• Accuracy across groups								
Breakend Impact tracker ⁴¹ [AE]	Performance	• Estimation of energy consumption • Estimation of GPU memory consumption								
	Respect for public interest	• Measure of environmental impact	✓				Agnostic		• Dot plots • Bar charts	
	Representative contestations [174]	• Mechanisms for users to ask questions and record disagreement with system behaviour								
Active contestations [AF]	Contestability					✓	Agnostic			

Table 15. Mapping of available means for transmitting value-specific manifestations to different stakeholders based on the purpose of their insight and the nature of their knowledge (DT = Development Team; AT = Auditing Team; DE = Data Domain Experts; DS = Decision Subjects). The identification and color code correspond to those on table 14. Each means is linked to the value and criteria manifestations that they communicate, the stakeholders that the original papers address, model specificity, deployed approach, visual elements and any additional details.

⁴⁰<https://github.com/microsoft/responsible-ai-toolbox/blob/main/docs/erroranalysis-dashboard-README.md>

⁴¹<https://github.com/Breakend/experiment-impact-tracker>