Post-hoc XAI Methods: Counterfactuals and XCBR Applications

Betül Bayrak

19.12.2024



Outline

Motivation

Part-1:

- Terms and concepts
- XAI framework: Factors, Components, Tasks
- Overview of explainers

Part-2:

My Ph.D.:

- RQs
- Papers
- Contributions
- Limitations

Motivation (XAI)



https://tensor-solutions.com/explainable-ai https://www.geeksforgeeks.org/explainable-artificial-intelligencexai/

PART-1



Explainability vs Interpretability

Interpretability refers to the visibility and **understanding of the inner logic** and mechanics of the AI model.

Explainability is the **ability to describe the behavior of a system** in an understandable medium to humans.





XAI Framework Overview



Each component supports the broader goal of explainability, allowing AI systems to generate insights that align with technical demands and user requirements.

XAI framework:

- 1. Human-centered design
- 2. Factors in development
- 3. Components and tasks

Human-centered design

Users are the **recipients** of the system (in most cases)

The human-centered approach focuses on the **what**, **when**, **and how to explain things to human end users** through the explanation process that involves the users in the development process.



Factors in development

I. Domain Analysis and Requirements

- End Users
- Knowledge Sources
- Actionability
- Ethical Considerations
- Inputs and Outputs
- Assessments
- Reusability

II. Multi-modal Interaction and Human-Centered Design

- Coverage
- Personal Preferences
- Usability Testing





Explainers						
Explainers	Ante-hoc	Decision Trees				
	Explainers	Generalized Additive Models (GAMs Bayesian Rule Lists (BRLs)				
		Linear Models & Logistic Regression				
	Post-hoc	Instance-Based Explainers	Counterfactuals Semifactuals Alterfactuals Prototypes & Criticisms			
	Explainers	Attribution-Based Explainers	SHAP Integrated Gradients Grad-CAM Saliency Maps			
		Rule-Based Explainers	Anchors Decision rules			
		Visualization-Based Explainers	PDPs: Marginal effect of features ICE Plots ALE Plots			



Instance-based Explainers





NLN: Nearest Like Neighbour NUN: Nearest Unlike Neighbour

Instance-based Explainers

- NLN: Nearest Like Neighbour
- **NUN:** Nearest Unlike Neighbour
- **Counterfactual Explanations:** Generate hypothetical scenarios to describe how changing the input features can alter the prediction.
- Semifactual Explanations: Suggest that even with changes in certain attributes, the model's outcome would remain the same. "even if... still..." scenarios
- Alterfactual Explanations: Demonstrate the irrelevance of certain features by proposing changes that do not alter the model





NTNU | Norwegian University of Science and Technology

17

Case-Based Reasoning (CBR)

CBR is a problem-solving methodology that uses past experiences to address new problems.

CBR operates through a four-step cycle:

- 1. **Retrieve**: Identify and retrieve the most relevant past cases from the case base.
- 2. Reuse: Adapt the retrieved case(s) to solve the new problem.
- 3. **Revise**: Test the proposed solution and refine it if necessary.
- 4. Retain: Store the new case and solution in the case base for future use.





XCBR

- CBR is a flexible and interpretable methodology.
- It is used to explain AI models



https://gaia.fdi.ucm.es/research/excbr/

Evaluation





NTNU | Norwegian University of Science and Technology
Corce on a

Coroama, Loredana, and Adrian Groza. "Evaluation metrics in explainable artificial intelligence (XAI)." *International conference on advanced research in technologies, information, innovation and sustainability* Cham: Springer Nature Switzerland, 2022.

PART-2



Focus



Focus

Explainers	Ante-hoc	Decision Trees Generalized Additive Models (GAMs Bayesian Rule Lists (BRLs) Linear Models & Logistic			
	Explainers				
		Linear Models & Logistic Regression			
	Post-hoc	Instance-Based Explainers	Counterfactuals Sem if a ctuals Alterfa ctuals Prototy pes & Critic ism s		
	Explainers	Attribution-Based Explainers	SHAP Integrated Gradients Grad-CAM Saliency Maps		
		Rule-Based Explainers	Anchors Decision rules		
		Visualization-Based Explainers	PDPs: Marginal effect of features ICE Plots ALE Plots		



My Ph.D.





Research questions

ResearchRQ1: Domain KnowledgeQuestionsand Design

RQ2: Generating plausible explanations	RQ2.1: High-quality counterfactuals		
•	RQ2.2: Multimodel counterfactuals		
RQ3: Improving comprehensibility	RQ3.1: Instance-based explanations		
. ,	RQ3.2: Evaluating explanations		



Research questions - contributions

Research Questions			P-II	P-III	P-IV	\mathbf{P} -V
RQ 1: Domain Knowledge and Design			*			*
BO 2. Concepting	RQ 2.1: High-quality			*	*	
NG 2: Generating	Counterfactuals					
Plausible Explanations	RQ 2.2: Multimodal					*
	Counterfactuals					
BO 2. Incompanying	RQ 3.1: Instance-	*	*		*	*
Comprehensibility	Based Explanations					
	RQ 3.2: Evaluating				*	*
	Explanations				•	

Paper 1: When to Explain?

A model-agnostic XCBR framework that **selectively** triggers explanations



Paper 1: When to Explain?



" The prediction result is the same with 4 out of 4 closest samples. However, the sample is at risk; in similar cases, when the f1 feature increases by 0.14 and f2 decreases by 0.05, decisions change. "

- 1. Store Data: Maintain a case base with pairs of samples and their counterfactuals.
- 2. Retrieve Pair: Identify the most similar sample-counterfactual pair.
- **3. Show Differences**: Highlight key differences between the sample and counterfactual to provide actionable insights.
- 4. Visualization: Use a combination of a bidirectional bar graph and text annotations to present the results effectively.

Paper 1: Contributions

Research Questions RQ1: Domain Knowledge and Design

RQ2: Generating plausible explanations	RQ2.1: High-quality counterfactuals		
•	RQ2.2: Multimodel counterfactuals		
RQ3: Improving comprehensibility	RQ3.1: Instance-based explanations		
	RQ3.2: Evaluating explanations		



Paper2 - A Twin XCBR System Using Supportive and Contrastive Explanations



Paper2 - A Twin XCBR System Using Supportive and Contrastive Explanations





Paper2 - A Twin XCBR System Using Supportive and Contrastive Explanations



$support = \frac{ S }{ S }$	
$rigidity = \begin{vmatrix} n \\ 1 - \frac{su}{su} \end{vmatrix}$	$\frac{pport}{acc}$

NTNU | Norwegian University of Science and Technology

	accuracy	support	rigidity
Use-case 1	0.24	0.476	0.984
Use-case 2	0.58	0.696	0.199
Use-case 3	0.77	0.9897	0.2197

Use-case 3: Wine Quality



Paper 2: Contributions

Research Questions RQ1: Domain Knowledge and Design

RQ2: Generating plausible explanations	RQ2.1: High-quality counterfactuals		
•	RQ2.2: Multimodel counterfactuals		
RQ3: Improving comprehensibility	RQ3.1: Instance-based explanations		
	RQ3.2: Evaluating explanations		



Paper3 - PertCF: A Perturbation-Based Counterfactual Generation Approach



Norwegian University of

Science and Technology



(d)

NUN X NLN

PertCF: a perturbation based CF generation approach

PertCF: a perturbation based CF generation approach



NUN С Х NLN

PertCF: a perturbation based CF generation approach

- Perturb 'x' to generate candidate c₁
- SHAP (SHapley Additive exPlanations)



Norwegian University of Science and Technology

PertCF: a perturbation based CF generation approach

- Perturb 'x' to generate candidate c₁
- SHAP (SHapley Additive exPlanations)

$$c_i = \langle C_{i1'}, C_{i2'}, \cdots, C_{im} \rangle$$

$$c_{ia} = s_a + shap_target_a * dist(t_a, s_a)$$

s: source t: target



PertCF: a perturbation based CF generation approach





Norwegian University of Science and Technology

PertCF: a perturbation based CF generation approach

Termination criteria:

- 1. Number of iterations
- 2. Distance between the last two candidate



Paper3 - PertCF

	Dissim.	Sparsity	Sparsity Instability	
		South Ger	rman Credit	
DICE	0.055744	0.911053	0.056013	0.073553
CFshap	0.255488	0.584211	0.255525	0.005750
PertCF	0.051702	0.798246	0.051817	0.406937
	Ŭ	ser Knowl	edge Modelin	ıg
DICE	0.172721	0.642276	0.176896	0.117947
CFshap	0.179172	0.029268	0.180618	0.001123
PertCF	0.063630	0.058537	0.066436	0.282734





Paper 3: Contributions

ResearchRQ1: Domain KnowledgeQuestionsand Design

RQ2: Generating plausible explanations	RQ2.1: High-quality counterfactuals		
•	RQ2.2: Multimodel counterfactuals		
RQ3: Improving comprehensibility	RQ3.1: Instance-based explanations		
	RQ3.2: Evaluating explanations		



Paper4 - Evaluation of Instance-based Explanations: An In-depth Analysis of Counterfactual Evaluation Metrics, Challenges, and the CEval Toolkit

- Motivation:
 - Confusion in Literature
 - Benchmarking Needs
 - Lack of Open-Source Tools
- Contributions:
 - **CEval Toolkit:** An open-source platform for evaluating and optimizing instance-based explanations.
 - **Customizable Framework:** Adaptable for different datasets and user needs.
 - Focus on Optimization: Integrates methods to enhance explanation quality.
- Impact:
 - Provides clarity in evaluating instance-based explanations.
 - Bridges technical evaluation with user-centric goals.
 - Serves as a benchmark for future XAI systems.

Metrics

$$validity = \frac{1}{m} \sum_{i=1}^{m} \mathbb{1}_{z_i \neq y} \mathbb{1}_{z_i = f(x'_i)}$$
$$proximity = \frac{1}{m} \sum_{i=1}^{m} dist(x, x'_i)$$
$$sparsity = \frac{1}{m\rho} \sum_{i=1}^{m} \sum_{j=1}^{\rho} \mathbb{1}_{a'_j \neq a_j}$$

 $number_of_explanations = |e|$

 $diversity_dpp = det(K)$ $lcc = \frac{u_z}{l} \quad diversity_{lcc} = lcc * diversity_dpp$ $yNN = \frac{1}{mk} \sum_{i=1}^{m} \sum_{\psi_j \in kNN(x'_i)} \mathbb{1}_{z_i = f(\psi_j)}$ $feasibility = \frac{1}{mk} \sum_{i=1}^{m} \sum_{\psi_j \in kNN(x'_i)}^{k} dist(x'_i, \psi_j)$ $kNLN_distance = \frac{1}{mk} \sum_{i=1}^{m} \sum_{\vartheta_j \in kNLN(x'_i)}^{k} dist(x'_i, \vartheta_j)$

realtive_distance =
$$\frac{1}{m} \sum_{i=1}^{m} \frac{dist(x_i', x)}{dist(NUN(x), x)}$$

NTNU | Norwegian University of Science and Technology

$$\begin{aligned} \text{redundancy} &= \frac{1}{m|S|} \sum_{i=1}^{m} \sum_{s \in S} \mathbb{1}_{f(\tau'(x'_{i}, x, s)) = z_{i}} \\ \text{robustness} &= \frac{1}{m} \sum_{i=1}^{m} \text{dist}(e_{i}, explain(P(x))) \\ \text{plausibility} &= \frac{1}{m} \sum_{i=1}^{m} \frac{\text{dist}(x'_{i}, NLN(x'_{i}))}{\text{dist}(NLN(x'_{i}), NUN(NLN(x'_{i})))} \\ \text{coverage} &= \frac{1}{|X|} \sum_{x \in X} \mathbb{1}_{|e| > 0}, \ e = explain(x) \\ \text{rigidity} &= \left|1 - \frac{\sup p}{acc}\right| \\ \text{discriminative power} \\ \text{vulnerability} \\ \text{computational complexity} \\ \text{constraint violation} \end{aligned}$$

 TABLE 1. Unified Notation for the Rest of the Article

Symbol	Description
x	An instance to explain, $x = \langle a_1,, a_\rho \rangle$.
ρ	The number of attributes.
X	Set of instances to explain, $x \in X$.
f()	Prediction function.
y	Prediction result for $x, y = f(x)$
explain()	Explanation function, returns a set of counterfactuals for
	an instance.
e	Set of counterfactuals for x, $e = explain(x)$ and $e =$
	$\{e_1,, e_m\}.$
m	Number of counterfactuals that are provided for x.
e_i	$e_i = \{x'_i, z_i\}, e_i \in e.$
x'_i	$x'_{i} = \langle a'_{1},, a'_{o} \rangle$
zi	$z_i = f(x_i')$
u_z	The number of unique class labels for the counterfactuals
	generated for x.
1	The number of unique class labels.
R_i	Range of <i>j</i> th attribute.
k	Number of neighbours.
kNN()	kNN(x) finds k Nearest Neighbours to x.
kNLN()	kNLN(x) finds k Nearest Like Neighbours to x.
NUN()	Nearest Unlike Neighbour, $NUN(x)$ returns the nearest
	neighbour which is labelled different than x .
$dist()^*$	Distance function, quantifies the distance between two
	instances.
$P()^*$	Perturbation function, $P(x)$ returns perturbed/corrupted
	version of x, x^p .
$\tau()$	$\tau(x'_i, x, j)$ copies j^{th} attribute of x to j^{th} attribute of x'_i and
	returns x'_i .
A	$A = \{1, 2,, \rho\}$
$[A]^l$	All subsets of A that have l elements. $[A]^2 = \{\{a, b\}:$
	$a, b \in A, a \neq b$
S	All possible subsets (power set) of A without \emptyset and A
	(super set). $S = \bigcup_{l=1}^{p-1} [A]^{l}$
	$ S = 2^{\rho} - 2$
$\tau'()$	$\tau'(x'_i, x, s)$ copies s attributes of x to s attribute of x' and
	returns x_i^{\prime} .

* These functions can be implemented in different ways.

		Applica	ble		Requires		Short Description
	Generated	Existed	Single	Multi	Data	Model	Short Description
Validity ¹	x	x	x	x	-	х	Whether the decision was altered.
Proximity ¹	x	x	x	x	-	-	Mean of feature-wise distance between the instance and its counterfactuals.
Sparsity ¹	х	x	x	x	-	-	Mean number of altered features between the instance and its counterfactuals.
Number of counter- factuals ¹	x	x	-	x	-	-	Number of counterfactuals generated for an instance.
Diversity ¹	x	x	-	x	-	-	Mean proximity between counterfactuals.
Diversity_ lcc ¹	x	x	-	x	x	-	Diversity with class coverage coefficient.
yNN ¹	X	x	x	x	x	х	Amount of support counterfactuals receive from positively clas- sified background data.
Feasibility ¹	X	x	x	x	x	-	Mean proximity of the counterfactuals to their nearest observa- tions in the background data.
kNLN Distance ¹	x	x	x	x	x	-	Mean distance of counterfactuals to their k-NLN.
Relative Distance ¹	X	-	x	x	x	-	The ratio of the mean distance between the instance and coun- terfactuals to the mean distance between the instance and its NUN.
Redundancy	x	x	x	x	-	х	Mean count of unnecessary feature changes.
Robustness	x	-	x	x	x	х	Mean proximity between the explanation of the instance and the explanation of a slightly corrupted version of the instance.
Plausibility	x	-	x	x	x	-	The degree of credibility in the context.
Discriminative Power	x	X	-	x	x	-	The ability to differentiate two distinct classes through a naive approach.
Vulnerability	x	-	x	x	x	-	The extent of susceptibility to manipulations.
Complexity	x	x	x	x	x	х	Cost for the explainer to generate a single explanation.
Constraints	x	x	x	x	-	-	Mean count of violated pre-defined constraints.
Coverage	x	X	x	x	-	-	The ability to generate valid counterfactuals across various types of instances.

¹This metric is implemented in the package.

Paper 4: Contributions

ResearchRQ1: Domain KnowledgeQuestionsand Design

RQ2: Generating plausible explanations	RQ2.1: High-quality counterfactuals			
	RQ2.2: Multimodel counterfactuals			
RQ3: Improving comprehensibility	RQ3.1: Instance-based explanations			
	RQ3.2: Evaluating explanations			

Paper5 - An Empirical Analysis of User Preferences Regarding XAI metrics

• Motivation:

- Understanding user preferences is crucial for designing effective XAI systems.
- A need to bridge the gap between technical metrics and usercentric needs.
- Aim:
 - Conduct an empirical study to evaluate which XAI metrics align with user priorities.
 - Provide actionable insights for improving XAI frameworks.





e-based



This image has been classified as a(n)

OBELISK



Please, according to your own criteria, click on the approach that best explains this classification from the options below.

Occlusion

Key Contributions:

- Empirical Analysis: Quantifies user preferences for common XAI metrics (e.g., interpretability, fidelity, robustness).
- Insights on Trade-offs: Highlights trade-2. offs between technical metrics and user satisfaction.
- 3. Guidelines for XAI Design: Provides recommendations to align technical and user-centric goals.



Norwegian University of Science and Technology

GradCAMPP Sobol HSIC

Paper5



Paper5



- Some of the objective evaluation metrics align better with user preferences.
- Metrics such as insertion and sparsity are good predictors of user satisfaction (for this specific task).
- Novice users prefer instance-based explanations.
- Expert users prefer feature-based explanations.
- Explanation representation MATTERS.
- Objective + Subjective Evaluation

Paper 5: Contributions

Research Questions RQ1: Domain Knowledge and Design

RQ2: Generating plausible explanations	RQ2.1: High-quality counterfactuals		
<u>-</u>	RQ2.2: Multimodel counterfactuals		
RQ3: Improving comprehensibility	RQ3.1: Instance-based explanations		
. ,	RQ3.2: Evaluating explan	ations	



Contributions - overall

Research Questions		P-I	P-II	P-III	P-IV	P-V
RQ 1: Domain Knowledge and Design		*	*			*
PO 2. Concepting	RQ 2.1: High-quality			*	*	
Plausible Explanations	Counterfactuals			-	-	
	RQ 2.2: Multimodal					*
	Counterfactuals					
DO 2. Immenouing	RQ 3.1: Instance-	* *	*		*	*
Comprehensibility	Based Explanations				-	
	RQ 3.2: Evaluating				*	*
	Explanations					



Limitations

- Human-in-the-Loop Challenges:
 - Reliance on human evaluation introduces biases and subjectivity.
 - Balancing objective metrics with user feedback remains complex.
- Multi-Modality Constraints:
 - Current methods focus on single or limited number of data types (e.g., tabular, image).
 - Limited application to real-world multimodal data (e.g., text, audio, video).
- **Scalability Issues:** High computational demands restrict application to large datasets and real-time systems.
- **Objective Mismatch:** Misalignment between stakeholder needs (users, developers, regulators) and explanation goals.
- **Customization Gaps:** Current metrics do not align with user preferences for tailored, flexible explanations.
- Empirical Validation: Limited generalizability to diverse domains and broader user groups.



Future Directions

	Scalability and Efficiency	Develop algorithms for scalable, real-time explanations with minimal computational overhead.	
	Multimodal Explainability	Extend methods to handle diverse data types and cross-modal relationships.	
0	Automated and Adaptive Explanations	Create systems that adapt explanations to user expertise and preferences in real time.	
\$.	Human-Centered XAI	Incorporate participatory design approaches and longitudinal studies to refine explanations.	
	Integration with Emerging Trends	Adapt XAI methods for transformers, LLMs, and conversational AI applications.	
NTNU Norwegian University of Science and Technology			

Conclusion



Conclusion

Advancing Instance-	 Developed innovative methods for generating and evaluating
Based Explanations:	instance-based explanations, particularly counterfactuals.
Bridging Human and	 Focused on human-centered design to align AI explanations
Machine Understanding:	with user needs, fostering trust and usability.
Comprehensive	 Proposed and implemented metrics and toolkits like CEval to
Evaluation Frameworks:	assess and optimize XAI techniques effectively.
Contribution to XAI Literature:	 Enhanced understanding of trade-offs between technical metrics and user preferences, providing actionable insights for future research.
Comprehensive	 Proposed and implemented metrics and toolkits like CEval to
Evaluation Frameworks:	assess and optimize XAI techniques effectively. Enhanced understanding of trade-offs between technical
Contribution to XAI	metrics and user preferences, providing actionable insights for
Literature:	future research.



Thank you for your attention ©

Betül Bayrak

betul.bayrak@ntnu.no

