Questioning the AI: Towards Human Centered Explainable AI (XAI)

Research work 2018-2021

Q. Vera Liao IBM **Research**

Our HCI research: Bridging work

Transfer emerging AI research into tangible tools and guidelines that support product teams to navigate the design space



Explainable AI (XAI): Definition

Narrow (ML) definition:

Techniques and methods that make a ML model's decisions understandable by people **Broader definition:**

Everything that makes Al understandable (e.g., also including data, function, performance, etc.)

XAI is not just ML (also explainable robotics, planning, etc.), but our current work focuses on **explaining supervised ML**

Al is increasingly used in many high-stakes tasks



The quest for explainable AI (XAI)

Companies Grapple With AI's Opaque Decision-Making Process

We Need AI That Is Explainable, Auditable, and Transparent

Why "Explainability" Is A Big Deal In AI

From black box to white box: Reclaiming human power in Al

How Explainable AI Is Helping Algorithms Avoid Bias



The needs for XAI algorithms



XAI "post-hoc" algorithm example: LIME



Neural network, not directly explainable

(a) Original Image



LIME (Ribeiro et al. 2016)

Use a post-hoc XAI technique



(b) Explaining *Electric guitar* (c) Explaining *Acoustic guitar*

(d) Explaining Labrador



We will be teaching a virtual tutorial on this at CHI 2021! https://hcixaitutorial.github.io/



Machine Learning Interpretability: A Survey on **Methods and Metrics**

Diogo V. Carvalho ^{1,2,*}, Eduardo M. Pereira ¹ and

- ¹ Deloitte Portugal, Manuel Bandeira Street, 43, 4150-47
- ² Faculty of Engineering, University of Porto, Dr. Rober
- ³ INESC TEC, Dr. Roberto Frias Street, 4200-465 Porto, 1
- * Correspondence: diocarvalho@deloitte.pt

Received: 21 June 2019; Accepted: 24 July 2019; Published

Abstract: Machine learning systems are becoming in has been expanding, accelerating the shift towar algorithmically informed decisions have greater pe most of these accurate decision support systems rem logic and inner workings are hidden to the user (ratic

Explaining Explanations: An Overview of Interpretability of Machine Learning

Leilani H. Gilpin, David Bau, Ben Z. Yuan, Ayesha Bajwa, Michael Specter and Lalana Kagal Computer Science and Artificial Intelligence Laboratory Massachusetts Institute of Technology Cambridge, MA 02139 {lgilpin, davidbau, bzy, abajwa, specter, lkagal}@ mit.edu

Abstract-There has recently been a surge of work in ex-As a first step towards creating explanation mechanisms planatory artificial intelligence (XAI). This research area tackles there is a new line of research in interpretability, loosel the important problem that complex machines and algorithms defined as the science of comprehending what a model did (c

le models and learning method les include visual cues to fin A Multidisciplinary Survey and Framework for Design and

networks in image recognition

Received August 5, 2018, accepted September 4, 2018, date of publication September 17, 2018, date of current version October 12, 2018. Digital Object Identifier 10.1109/ACCESS.2018.2870052

Peeking Inside the Black-Box: A Survey on **Explainable Artificial Intelligence (XAI)**

AMINA ADADI[©] AND MOHAMMED BERRADA

arv Physics Laboratory, Sidi Mohammed Ben Abdellah University, Fez 30050, Morocco Corresponding author: Amina Adadi (amina.adadi@gmail.com)

ABSTRACT At the dawn of the fourth industrial revolution, we are witnessing a fast and widespread adoption of artificial intelligence (AI) in our daily life, which contributes to accelerating the shift towards a more algorithmic society. However, even with such unprecedented advancements, a key impediment to the use of AI-based systems is that they often lack transparency. Indeed, the black-box nature of these systems allows powerful predictions, but it cannot be directly explained. This issue has triggered a new debate on explainable AI (XAI). A research field holds substantial promise for improving trust and transparency of

SINA M

The Evalua

ques

ERIC D. The need f

intelligenc reasoning to define, o on differen

A growing collection of XAI techniques

netot^{b,e,f}, olina^g.

omies,

challenges for identifying appropriate design and evaluation methodology and consolidating knowledge from across efforts. To this end, this paper presents a survey and framework intended to share knowledge and experiences of XAI design and evaluation methods across multiple disciplines. Aiming to support diverse design goals and evaluation method in XAI research, after a thorough review of XAI related papers in the fields of machine learning, visualization, and human-computer interaction we pre-

A Survey of Methods for Explaining

RICCARDO GUIDOTTI, ANNA MONREALE, SALV/ FRANCO TURINI, KDDLab, University of Pisa, Italy FOSCA GIANNOTTI, KDDLab, ISTI-CNR, Italy DINO PEDRESCHI, KDDLab, University of Pisa, Italy

In recent years, many accurate decision support systems have systems that hide their internal logic to the user. This lack of ex ethical issue. The literature reports many approaches aimed at c at the cost of sacrificing accuracy for interpretability. The appli can be used are various, and each approach is typically develope and, as a consequence, it explicitly or implicitly delineates its ov tion. The aim of this article is to provide a classification of the m respect to the notion of explanation and the type of black box box type, and a desired explanation, this survey should help the

Explanation Methods in Deep Learning: Users, Values, Concerns and Challenges^{*}

Gabriëlle Ras, Marcel van Gerven, Pim Haselager

Radboud University, Donders Institute for Brain, Cognition and Behaviour, Nijmegen, the Netherlands {g.ras, m.vangerven, w.haselager}@donders.ru.nl

Abstract

Issues regarding explainable AI involve four components: users, laws & regulations, explanations and algorithms. Together these components provide a context in which explanation methods can be evaluated regarding their adequacy. The goal of this chapter is to bridge the gap between expert users and lay users. Different kinds of users are identified and their concerns revealed, relevant statements from the General Data Protection Regulation are analyzed in the context of Deep Neural Networks (DNNs), a taxonomy for the classification of existing explanation methods is introduced, and finally, the various classes of explanation methods are analyzed to verify if user concerns are justified. Overall, it is clear that (visual) explanations can be given about various aspects of the influence of the input on the output. However, it is noted that avalanation mathods or interfaces for law users are missing and we encoulate which criteria

[•]ENSTA, Institute Polytechnique Paris and INRIA Flowers Team, Palaiseau, France ^cUniversity of the Basque Country (UPV/EHU), 48013 Bilbao, Spain ^dBasque Center for Applied Mathematics (BCAM), 48009 Bilbao, Bizkaia, Spain ^eSegula Technologies, Parc d'activité de Pissaloup, Trappes, France ^fInstitut des Systèmes Intelligents et de Robotique, Sorbonne Universitè, France

mputational Intelligence, University of Granada, 18071 Granada, Spain nica, 28050 Madrid, Spain

(AI) has achieved a notable momentum that, if harnessed tions over many application sectors across the field. For this ire community stands in front of the barrier of explainability, brought by sub-symbolism (e.g. ensembles or Deep Neural type of AI (namely, expert systems and rule based models). in the so-called *eXplainable* AI (XAI) field, which is widely ctical deployment of AI models. The overview presented in id contributions already done in the field of XAI, including a r this purpose we summarize previous efforts made to define ning a novel definition of explainable Machine Learning that th a major focus on the audience for which the explainability propose and discuss about a taxonomy of recent contributions

XAI in Practice



An abundance of XAI algorithms

IBM Research Trusted AI			Home	Demo	Resources		
AI Explainability 360							
This extensible open source toolkit can help you comprehend how machine learning models predict labels by various means throughout the AI application lifecycle. We invite you to use it and improve it.							
API Docs / Get Code /							
Not sure what to do first? S	Start here!						
Read More	Try a Web Demo	Watch Vide	eos	Re	ad a Paper		
Learn more about explainability concepts, terminology, and tools before you begin.	Step through the process of explaining models to consumers with different personas in an interactive	Watch videos to about AI Explain toolkit.	learn more nability 360	Rea we Exp	ad a paper descr designed AI plainability 360 t		

Toolbox of XAI techniques

From academic research into a practitioners' toolbox



An abundance of XAI algorithms

🔊 ALIBI , 🎯 Captum IBM Research Trusted AI Demo AI Explainability 360 This extensible open source toolkit can help you comprehend how machine learning models predict labels by various means throughout the AI application lifecycle. We invite you to use it and improve it. README.md lack-InterpretML · API Docs Get Code / license MIT python 3.6 | 3.7 maintained yes Not sure what to do first? Start here! In the beginning n aggled in the v Read More Try a Web Demo Watch Videos **Read a Paper** Learn more about Step through the process of Watch videos to learn more Read a paper descr Let there be light. explainability concepts, explaining models to about AI Explainability 360 we designed AI InterpretML is an open-se Explainability 360 t terminology, and tools before consumers with different toolkit. interpretability technique vou begin. personas in an interactive models and explain blackbox systems, interpretivit, neips you understand your through s groupes or Py forch models and can be used behavior, or understand the reasons behind individual predictions. odification to the original neural network.

Toolbox of XAI techniques

XAI in Practice

From academic research into a practitioners' toolbox



Inter-disciplinary perspectives

XAI in Practice

What-If Tool demo - regression model fo Datapoint setter Purturano Fast Visuates © Integrating O Parts	or predicting age - UCI census income dataset	ß	ALIBI	
Control of the second sec	IBM Research Trusted AI AI Explainability 3 This extensible open sour models predict labels by v you to use it and improve	60 ce toolkit can help you comp various means throughout th it.	Home prehend how machine lea he AI application lifecycle	Demo Resources urning . We invite d lack-
2. EEE InterpretMI Elicense MIT python 3.6 maintained yes	API Docs Z Get Code Z	Start here!		,
In the beginning straggled in the Let there be ligi InterpretML is an oper interpretability technic models and explain bil behavior, or understar	g n e v Read More ht. Learn more about explainability concepts, h-si terminology, and tools before you begin. ackoox systems. merpretime neips you un d the reasons behind individual predictio	Try a Web Demo Step through the process of explaining models to consumers with different personas in an interactive interstahu' youn 'hrouch's: goorpes or Py In ns. dification to	Watch Videos Watch videos to learn more about AI Explainability 360 toolkit. oron models and can be used o the original neural network.	Read a Paper Read a paper descr we designed AI Explainability 360 t
	Toolbo	x of XAI	techniq	ues
			> ((\mathbf{x}
) (\$) ار ب
ר r	Fowards r many don	eal-worl nains an	d XAI: s d user g	erving groups



An abundance of XAI algorithms



Inter-disciplinary perspectives

Inter-disciplinary perspectives

- The gaps between XAI output and human explanations: contrastive, selective, socially interactive (Miller 2019; Mittelstadt et al. 2019)
- The plurality of motivation for explanation: diagnosis, predicting the future, sense-making, justification, reconciling dissonance, etc. (Kiel 2006; Lombrozo, 2006)
- Explanatory power is recipient dependent, including the question asked (explanatory relevance) (Hilton, 1990; Walton, 2004)
- More complexities:
 - The plurality of cognitive processes (Petty and Cacioppo, 1986; Kahneman, 2003)
 - Socio-technical systems (Ehsan et al., 2021)

From XAI algorithms to XAI UX

With a toolbox: How to **select?** How to **translate**?

Our paths:

- Develop domain-specific guidelines: HCI research with key XAI use cases
- Tackle the design process: User centered design of XAI

XAI in Practice



XAI use cases in AI lifecycle

Model evaluation and selection (IUI2021)

XAI consumer: Data scientist

Model development

Training

Evaluation

Debugging

Decision aid

Decision aid

Delation

Automation

Model

Data

Construction

Task

definition

Model auditing

Explainable active learning (CSCW 2020) XAI consumer: Annotator Trust calibration and decision support (FAT* 2020, CHI 2021 🞖) XAI consumer: Decision-maker

> Delegation support (ongoing) XAI consumer: Domain expert

Fairness assessment (IUI 2019 8) XAI consumer: Regulator, impacted groups

XAI for model evaluation and selection

	f1	accuracy	roc_auc	precision	recall	neg_log_loss
LGBM_2	0.922	0.923	0.923	0.926	0.918	-2.66
LogisticRegression_2	0.699	0.712	0.712	0.725	0.675	-9.95
DecisionTree_2	0.694	0.707	0.706	0.719	0.67	-10.1
RandomForest_2	0.752	0.755	0.755	0.756	0.747	-8.46

(a) Screenshot of the Metrics Table showing metrics for four selected models.







How does each model make predictions? Why are these instances predicted differently by these models? Why is this model making a wrong prediction?

Narkar et al. Model LineUpper: Supporting Interactive Model Comparison at Multiple Levels for AutoML. IUI 2021

XAI for fairness assessment





Is the way the model makes risk predictions fair? Why is this person predicted of high risk? Is he/she treated fairly?

Lessons learned: From XAI algorithms to XAI UX

- XAI UX does not end with a single XAI algorithm
 - Often need multiple types of explanations
 - Able to anticipate when and where users want what explanations
- XAI algorithmic output is not necessarily neutral
 - People have different ability and motivation to process it deliberatively
 - Invoke different reactions
 - Can unequalize even marginalize certain groups
- "Translation" design is often necessary
 - Algorithmic output needs communication, elaboration, constraints, integration, etc.
 - Sometimes requires adapting the algorithm
- Break-downs are more often, translation are more imperative, on the "model usage" side

XAI in Practice



Why break-downs in model usage?



Why break-downs in model usage?



From XAI algorithms to XAI UX

With a toolbox: How to select? How to translate? How to expand?

Our paths:

- Develop domain-specific guidelines: HCI research with key XAI use cases
- Tackle the design process: User centered design of XAI
- Socially situated explainability by making visible the AI contexts

XAI in Practice



Towards "social transparency" in AI systems

Custo Reco Justi [0] Q	omer: Scout Inc. ommendation: Sell at ification: the AI syste Duota goals	Product: Access I \$100 per account per m m considered the follow [o] <i>Comparative</i>	Management (SaaS) nonth ving components <i>pricing: what similar custo</i>	Product ID (PID) mers pay): 43523X [O] Cost: \$55 /account/month
	For this customer, However, 1 out 3 h	3 members of your tean have sold at the recomm	m received pricing recomm nended price. Click to see m	endations in past nore details.	⁺ sales.
	(Nadia M. ■ Sales Assoc. (AB34)	Action: Reject Recommendat Comment: Long-term profitat selling at cost price to maintat © Oct 2, 2019	tion ↔ ble customer; main r ain relationship	Outcome: No Sale revenue from a different vertical ; 3
	(Eric C. Sales Manager (XZ89)	Action: Accept Recommendation Comment: Recommended private fair Was fair Dec 14, 2019	ation ↔ ice aligned with prof	Outcome: Sale it margins; customer felt the price
4W	What • Who • Why • When •	Jess W. ■ Sales Director (RE43)	 Action: Reject Recommendation Comment: Covid-19 pandemic offered 10% below cost price May 6, 2020 	tion ↔ ic mode; cannot lose	Outcome: Sale e long-term profitable customer; 5

Ehsan et al. Expanding Explainability: Towards Social Transparency in AI systems. To appear in CHI 2021 8

Why break-downs in model usage?



From XAI algorithms to XAI UX

With a toolbox: How to select? How to translate? How to expand?

Our paths:

- Develop domain-specific guidelines: HCI research with key XAI use cases
- Tackle the design process: User centered design of XAI
- Socially situated explainability by making visible the social contexts

XAI in Practice



Where we started: Research into XAI Design Practices

Why AI design practitioners?

 Bridging roles connecting user needs and XAI techniques

Research questions:

- What is the design space of XAI UX?
- What are the design challenges?





MDP

Machine Learning Interpretability: A Survey on Methods and Metrics

Diogo V. Carvalho ^{1,2,*}, Eduardo M. Pereira ¹ and

- ¹ Deloitte Portugal, Manuel Bandeira Street, 43, 4150-47
- ² Faculty of Engineering, University of Porto, Dr. Rober
- ³ INESC TEC, Dr. Roberto Frias Street, 4200-465 Porto, 1
- * Correspondence: diocarvalho@deloitte.pt

Received: 21 June 2019; Accepted: 24 July 2019; Published

Abstract: Machine learning systems are becoming in has been expanding, accelerating the shift towar algorithmically informed decisions have greater por most of these accurate decision support systems rem logic and inner workings are hidden to the user a ratic

Explaining Explanations: An Overview of Interpretability of Machine Learning

Leilani H. Gilpin, David Bau, Ben Z. Yuan, Ayesha Bajwa, Michael Specter and Lalana Kagal Computer Science and Artificial Intelligence Laboratory Massachusetts Institute of Technology Cambridge, MA 02139 {lgilpin, davidbau, bzy, abajwa, specter, lkagal}@ mit.edu

Abstract—There has recently been a surge of work in explanatory artificial intelligence (XAI). This research area tackles the important problem that complex machines and algorithms defined as the science of comprehending what a model did (c le models and learning method

mad
quesA Multidisciplinary Survey and Framework for Design and
TheTheEvaluat

Received August 5, 2018, accepted September 4, 2018, date of publication September 17, 2018, date of current version October 12, 2018. Digital Object Identifier 10.1109/ACCESS.2018.2870052

Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI)

AMINA ADADI[®] AND MOHAMMED BERRADA

Computer and Interdisciplinary Physics Laboratory, Sidi Mohammed Ben Abdellah University, Fez 30050, Morocco Corresponding author: Amina Adadi (amina.adadi@gmail.com)

ABSTRACT At the dawn of the fourth industrial revolution, we are witnessing a fast and widespread adoption of artificial intelligence (AI) in our daily life, which contributes to accelerating the shift towards a more algorithmic society. However, even with such unprecedented advancements, a key impediment to the use of AI-based systems is that they often lack transparency. Indeed, the black-box nature of these systems allows powerful predictions, but it cannot be directly explained. This issue has triggered a new debate on explainable AI (XAI). A research field holds substantial promise for improving trust and transparency of

SINA MOH

The need for intelligence a reasoning be to define, des on different c challenges fo across efforts experiences design goals

A technical space people are not quite in there yet... how to talk about it?

les include visual cues to fin

),e,f

es,

fields of machine learning, visualization, and human-computer interaction, we present a categorization of the cost of the second s

A Survey of Methods for Explaining

RICCARDO GUIDOTTI, ANNA MONREALE, SALV/ FRANCO TURINI, KDDLab, University of Pisa, Italy FOSCA GIANNOTTI, KDDLab, ISTI-CNR, Italy DINO PEDRESCHI, KDDLab, University of Pisa, Italy

In recent years, many accurate decision support systems have systems that hide their internal logic to the user. This lack of ex ethical issue. The literature reports many approaches aimed at o at the cost of sacrificing accuracy for interpretability. The appli can be used are various, and each approach is typically develope and, as a consequence, it explicitly or implicitly delineates its ov tion. The aim of this article is to provide a classification of the m respect to the notion of explanation and the type of black box box type, and a desired explanation, this survey should help the Explanation Methods in Deep Learning: Users, Values, Concerns and Challenges^{*}

Gabriëlle Ras, Marcel van Gerven, Pim Haselager

Radboud University, Donders Institute for Brain, Cognition and Behaviour, Nijmegen, the Netherlands {g.ras, m.vangerven, w.haselager}@donders.ru.nl

Abstract

Issues regarding explainable AI involve four components: users, laws & regulations, explanations and algorithms. Together these components provide a context in which explanation methods can be evaluated regarding their adequacy. The goal of this chapter is to bridge the gap between expert users and lay users. Different kinds of users are identified and their concerns revealed, relevant statements from the General Data Protection Regulation are analyzed in the context of Deep Neural Networks (DNNs), a taxonomy for the classification of existing explanation methods is introduced, and finally, the various classes of explanation methods are analyzed to verify if user concerns are justified. Overall, it is clear that (visual) explanations can be given about various aspects of the influence of the input on the output. However, it is noted that explanation methods or interfaces for lay users are missing and we execute a which criteria ¹ mputational Intelligence, University of Granada, 18071 Granada, Spain vica, 28050 Madrid, Spain

(AI) has achieved a notable momentum that, if harnessed utions over many application sectors across the field. For this ire community stands in front of the barrier of explainability, brought by sub-symbolism (e.g. ensembles or Deep Neural type of AI (namely, expert systems and rule based models). in the so-called *eXplainable* AI (XAI) field, which is widely uctical deployment of AI models. The overview presented in the contributions already done in the field of XAI, including a or this purpose we summarize previous efforts made to define the a major focus on the audience for which the explainability propose and discuss about a taxonomy of recent contributions

Study probe: algorithm informed XAI Questions

Category of Methods	Explanation Method	Definition	Algorithm Examples	Question Type
Explain the model	Global feature importance	Describe the weights of features used by the model (includ- ing visualization that shows the weights of features)	[41, 60, 69, 90]	How
(Global)	Decision tree approximation	Approximate the model to an interpretable decision-tree	[11, 47, 52]	How, Why, Why not, What if
	Rule extraction	Approximate the model to a set of rules, e.g., if-then rules	[26, 93, 102]	How, Why, Why not, What if
Explain a prediction	Local feature importance and saliency method	Show how features of the instance contribute to the model's prediction (including causes in parts of an image or text)	[61, 74, 83, 85, 101]	Why
(Local)	Local rules or trees	Describe the rules or a decision-tree path that the instance fits to guarantee the prediction	[39, 75, 99]	Why, How to still be this
Inspect coun- terfactual	Feature influence or relevance method	Show how the prediction changes corresponding to changes of a feature (often in a visualization format)	[8, 33, 36, 51]	What if, How to be that, How to still be this
	Contrastive or counterfactual features	Describe the feature(s) that will change the prediction if perturbed, absent or present	[27, 91, 100]	Why, Why not, How to be that
Example based	Prototypical or representative examples	Provide example(s) similar to the instance and with the same record as the prediction	[13, 48, 50]	Why, How to still be this
	Counterfactual example	Provide example(s) with small differences from the instance but with a different record from the prediction	[37, 55, 66]	Why, Why not, How to be that

- User needs for XAI are represented as prototypical questions
- A question can be answered by one or multiple XAI methods
- An XAI method can be implemented by one or multiple XAI algorithms

An explanation is an answer to a question (Wellman, 2011; Miller 2018) The effectiveness of an explanation depends on the question asked (Bromberger, 1992)



Question: Why is this husky classified as wolf?



XAI method: local feature (pixels) contribution

XAI algorithms:

- LIME (Ribeiro et al. 2016)
- SHAP (Lundberg and Lee 2017)
- ...

Study probe: algorithm informed XAI Questions

Category of Methods	Explanation Method	Definition	Algorithm Examples	Question Type
Explain the model	Global feature importance	Describe the weights of features used by the model (includ- ing visualization that shows the weights of features)	[41, 60, 69, 90]	How
(Global)	Decision tree approximation	Approximate the model to an interpretable decision-tree	[11, 47, 52]	How, Why, Why not, What if
	Rule extraction	Approximate the model to a set of rules, e.g., if-then rules	[26, 93, 102]	How, Why, Why not, What if
Explain a prediction	Local feature importance and saliency method	Show how features of the instance contribute to the model's prediction (including causes in parts of an image or text)	[61, 74, 83, 85, 101]	Why
(Local)	Local rules or trees	Describe the rules or a decision-tree path that the instance fits to guarantee the prediction	[39, 75, 99]	Why, How to still be this
Inspect coun- terfactual	Feature influence or relevance method	Show how the prediction changes corresponding to changes of a feature (often in a visualization format)	[8, 33, 36, 51]	What if, How to be that, How to still be this
	Contrastive or counterfactual features	Describe the feature(s) that will change the prediction if perturbed, absent or present	[27, 91, 100]	Why, Why not, How to be that
Example based	Prototypical or representative examples	Provide example(s) similar to the instance and with the same record as the prediction	[13, 48, 50]	Why, How to still be this
	Counterfactual example	Provide example(s) with small differences from the instance but with a different record from the prediction	[37, 55, 66]	Why, Why not, How to be that

+

Input (data), output, performance

(Lim et al., 2009)

Methodology

- Interviewed 20 participants
- 16 Al products in IBM
- 1. Walk through the AI system
- 2. Common questions users might ask
- 3. Discuss each question card
- 4. General challenges to create XAI products



- What is the *source* of the data?
- How are the *labels/ground-truth* produced?



Methodology

- Interviewed 20 participants
- 16 Al products in IBM
- 1. Walk through the AI system
- 2. Common questions users might ask
- 3. Discuss each question card
- 4. General challenges to create XAI products



- What is the *source* of the data?
- How are the *labels/ground-truth* produced?



XAI question bank

Data	 What kind of data does the system learn from? What is the source of the data? How were the labels/ground-truth produced? * What is the source is the source? 	Why	 Why/how is this instance given this prediction? What feature(s) of this instance leads to the system's prediction? Why are [instance A and B] given the same prediction?
Data	 • * What is the sample size? • * What data is the system NOT using? • * What are the limitations/biases of the data? 		 Why/how is this instance NOT predicted? Why is this instance predicted P instead of Q? Why are [instance A and B] given different predictions?
	 What kind of output does the system give? What does the system output mean? 		 What would the system predict if this instance changes to? What would the system predict if this feature of the instance
Output	 How can I best utilize the output of the system ? * What is the scope of the system's capability? Can it do? 	What If	 changes to? What would the system predict for [a different instance]?
	 * How is the output used for other system component(s) ? How accurate/precise/reliable are the predictions? How often does the system make misteles? 	How to be that	 How should this instance change to get a different prediction? How should this feature change for this instance to get a different prediction?
Performance	 How often does the system make mistakes? In what situations is the system likely to be correct/incorrect? * What are the limitations of the system? * What kind of mistakes is the system likely to make? 	How to still be this	 What kind of instance gets a different prediction? What is the scope of change permitted to still get the same prediction?
	 * Is the system's performance good enough for 		• What is the [highest/lowest/] feature(s) one can have to still get the same prediction?
	 How does the system make predictions? What features does the system consider? * Is [feature X] used or not used for the predictions? 		 What is the necessary feature(s) present or absent to guarantee this prediction? What kind of instance gets this prediction?
How (global)	 What is the system's overall logic? How does it weigh different features? What rules does it use? 		 * How/what/why will the system change/adapt/improve/drift over time? (change)
	 How does [feature X] impact its predictions? * What are the top rules/features it uses? * What kind of algorithm is used? 	Others	 * How to improve the system? (change) * Why using or not using this feature/rule/data? (follow-up) * What does [ML terminology] mean? (terminological)
	• * How are the parameters set?		• * What are the results of other people using the system? (social)

XAI design challenge 1: Variability of XAI needs

Diverse end goals for explainability

- To gain further insights for the decision
- To appropriately evaluate AI's capability
- To adapt usage or control
- To improve AI performance
- Ethical responsibilities of AI products

To gain further insights for the decision



Why How to be that

66

Users need to know why the system is saying this will be late because the reason is going to determine what their next action is...If it's because of a weather event, so no matter what you do you're not going to improve this number, versus something small, if you just make a quick call, you can get that number down (I-5)

To appropriately evaluate Al's capability



Performance How

There is a calibration of trust, whether people will use it over time. But also saying hey, we know this fails in this way (I-6)

XAI design challenge 1: Variability of XAI needs

Diverse end goals for explainability

- To gain further insights for the decision
- To appropriately evaluate AI's capability
- To adapt usage or control
- To improve AI performance
- Ethical responsibilities of AI products

Also varying XAI needs: User group, usage point, algorithm and data type, decision context

XAI design challenge 2: Gaps between algorithmic output and human explanations

Human explanations are

- Selective
- Contrastive
- Interactive
- Tailored for recipients



Design attempt to mimic how people, especially domain experts, explain

XAI design challenge 3: "in the dark" design process

- Challenge navigating the technical capabilities
- finding the right pairing to put the ideas of what's right for the user together with what's doable given the tools or the algorithms
- **Communication barriers** between designers, data scientists and other stakeholders
- Cost of time and resource impeding buy-in
- It remains in this weird limbo where people know it's important. People see it happen. They don't know how to make it happen. And everybody's feeling their way in the dark with no lights.



XAI in Practice

Opportunities for technical XAI work

- Explain data limitations and generalizability
- Explain output of multiple models
- Explain system changes
- Multi-level global explanations
- Interactive counterfactual explanations
- Social explanations
- Personalized and adaptive explanations

Guidelines to address XAI user needs

Input: Provide comprehensive transparency of training data, especially the limitations

Output: Contextualize the system's output in downstream tasks and the users' overall workflow

Performance: Help users understand the limitations of the AI and make it actionable

Global model: Choose appropriate level of details to explain the model

Local decision: Provide resources for "why not"

Counterfactual: Consider opportunities as utility features for analytics or exploration

Liao et al. Questioning the AI: Informing Design Practices for Explainable AI User Experiences. CHI 2020

Supporting the process: Question-driven XAI design



Pain points to address:

- Identify application, user and interaction specific XAI needs
- "Sensitize" to XAI techniques by situating their capabilities and values
- "Boundary objects" to support designer-data scientists collaboration

Question-Driven XAI Design

Step 1

Identify user Au questions qu

Step 2 Analyze questions

Step 3 Map questions to modeling solutions

Step 4

Iteratively design and evaluate

Elicit user needs for XAI as questions

Also gather user intentions and expectations for asking the questions Cluster questions into categories and prioritize categories for the XAI UX to focus on

Summarize user intentions and expectations to identify key user requirements Map prioritized question categories to candidate XAI techniques as a set of functional elements that the design should cover

A mapping guide for supervised ML is provided for reference Create a design including the candidate elements identified in step 3

Iteratively valuate the design with the user requirements identified in step 2 and fill the gaps

Designers, users Designers, product team Designers, data scientists

Designers, data scientists, users

XAI Question Bank



How to select: identify user needs for XAI as questions

Liao et al. Questioning the AI: Informing Design Practices for Explainable AI User Experiences. CHI 2020 🞖

Question	Explanations	Example XAI techniques
Global how	 Describe what algorithm is used and what features are considered, if a user is only interested in a high-level view Describe the general model logic as feature impact*, rules* or decision-trees• (sometimes need to explain with a surrogate simple model) 	<u>ProfWeight*, Feature Importance*,</u> <u>PDP</u> *, <u>BRCG</u> + , <u>GLRM</u> + , <u>Rule List</u> + , <u>DT Surrogate</u> •
Why	 Describe what key features of the particular instance determine the model's prediction of it* Describe rules* that the instance fits to guarantee the prediction Show similar examples• with the same predicted outcome to justify the model's prediction 	<u>LIME</u> *, <u>SHAP</u> *, <u>LOCO</u> *, <u>Anchors</u> +, <u>ProtoDash</u> •
Why not	 Describe what changes are required for the instance to get the alternative prediction and/or what features of the instance guarantee the current prediction* Show prototypical examples* that had the alternative outcome 	<u>CEM</u> * , <u>Prototype counterfactual</u> * , <u>ProtoDash</u> * (on alternative class)
How to be that	 Highlight features that if changed (increased, decreased, absent, or present) could alter the prediction* Show examples with small differences but had a different outcome than the prediction* 	<u>CEM</u> *, <u>Counterfactuals</u> *, <u>DiCE</u> +
What if	 Show how the prediction changes corresponding to the inquired change 	<u>PDP, ALE, What-if Tool</u>
How to still be this	 Describe feature ranges* or rules* that could guarantee the same prediction Show examples that are different from the particular instance but still had the same outcome 	<u>CEM</u> *, <u>Anchors</u> +
Performance	 Provide performance metrics of the model Show confidence information for each prediction Describe potential strengths and limitations of the model 	Precision, Recall, Accuracy, F1, AUC Confidence <u>FactSheets, Model Cards</u>
Data	 Document comprehensive information about the training data, including the source, provenance, type, size, coverage of population, potential biases, etc. 	<u>FactSheets, DataSheets</u>
Output	 Describe the scope of output or system functions Suggest how the output should be used for downstream tasks or user workflow 	<u>FactSheets, Model Cards</u>

How to translate: support collaborative problem-solving between data scientists and designers with "*boundary objects*"

Liao et al. Question-Driven Design Process for Explainable Al User Experiences. (Under review)



Liao et al. Question-Driven Design Process for Explainable Al User Experiences. (Under review)

Concluding remarks

- Contextualize the tools
- Actionable frameworks: what users? What contexts?
- User-centered design processes
- Shared design vision drives model development
- Thinking "outside the box"

From XAI algorithms to XAI UX



With a toolbox: How to select? How to translate? How to expand?



From AI algorithms to AI UX



Many toolboxes: How to select? How to translate? How to expand?

Thank YOU!

...and thanks to

Rachel Bellamy, Amit Dhurandhar, Jonathan Dodge, Upol Ehsan, Bhavya Ghai, Werner Geyer, Daniel Gruen, Jaesik Han, Michael Hind, Stephanie Houde, David Millen, Aleksandra Mojsilović, Sarah Miller, Klaus Mueller, Michael Muller, Shweta Narkar, Milena Pribić, John Richards, Mark Riedl, Daby Sow, Chenhao Tan, Richard Tomsett, Kush Varshney, Dakuo Wang, Justin Weisz, Yunfeng Zhang

> Q. Vera Liao <u>vera.liao@ibm.com</u> <u>www.qveraliao.com</u> @QVeraLiao