Advances and Open Questions in Explainable AI (XAI):

A practical perspective from an HCI researcher

Q. Vera Liao IBM **Research** What are the most significant **advances in XAI**?

How might they influence academia vs. industry?



Advance 1: From academic research into a practitioners' toolbox



XAI in Academia



Advance 2: Towards interdisciplinary 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 they ask (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)

XAI in Practice

Advances impacting practices With a toolbox: How to select? How to translate?

- Taxonomies and guidelines
- Empirical design and evaluation



Towards real-world XAI: serving many domains and user groups

Taxonomies and guidelines



Categorize XAI techniques

Categorize Contexts and User Groups for XAI

Design and evaluation of XAI: HCI work



Pre-ML

Evaluate and "critique" XAI

Translate and evaluate XAI in specific domains

What are the most significant advances in XAI? How might they influence industry?

- Advance 1: From academic research to a **practitioners' toolbox**
- •Advance 2: Interdisciplinary perspectives inform the **selection** and **translation** of the toolbox





Not a solved problem yet... more on open questions!

What are the most significant open questions in XAI?

How might they influence academia vs. industry?

Open Question 1, with a toolbox: How to select? How to translate?

Top-down guidelines are not enough due to limitations in, e.g. granularity and scalability

Possible paths:

- Allow users to select: intelligent and interactive XAI
- Tackle the design process: user centered design of XAI



Question-Driven XAI Design

Step 1

Identify user questions

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 8

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) 	ProfWeight*, Feature Importance*, <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 	PDP, ALE, What-if Tool
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</u> , <u>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)

Open Question 2: How to expand the toolbox?





Toolbox of XAI

How to operationalize the missing element and make it an actionable tool?

What's missing? From AI explanations to explainability



Explainability goal: sense-making for better decisions



Algorithmic explanation



Socially situated explainability





"Sense-making is not just about opening the closed box of AI, but also about who is around the box, and the socio-technical factors that govern the use of the AI system and the decision. Thus the 'ability' in explainability does not lie exclusively in the guts of the AI system"

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

Towards socially situated explainability: an actionable framework

Custo	mer: Scout Inc.	Product: Access Management (SaaS)	Product ID (PID): 43523X	
Recommendation: Sell at \$100 per account per month					
Justification: the AI system considered the following components					
[0] Qu	ota goals	$[\mathbf{o}]$ Comparative pricing: what similar cus	stomers pay	[o] Cost: \$55 /account/month	
For this customer, 3 members of your team received pricing recommendations in past sales. However, 1 out 3 have sold at the recommended price. Click to see more details.					



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

What are the most significant advances and open questions in XAI? How might they influence industry?

- •Open question 1: How to actually **select** and **translate** the toolbox into good user experiences?
- Open question 2: How to **expand** the toolbox to support explainability/understanding?



Have answers? Submit to our virtual workshop on Operationalizing Human-Centered Perspectives in Explainable AI at CHI 2021

https://hcxai.jimdosite.com/ (Deadline February 14)

Thank YOU!

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