

Should We Trust (X)AI? Design Dimensions for Structured Experimental Evaluations

Fabian Sperrle
Univeristy of Konstanz
fabian.sperrle@uni.kn

Mennatallah El-Assady
University of Konstanz
mennatallah.el-assady@uni.kn

Grace Guo
Georgia Institute of Technology
gguo31@gatech.edu

Duen Horng Chau
Georgia Institute of Technology
polo@gatech.edu

Alex Endert
Georgia Institute of Technology
endert@gatech.edu

Daniel Keim
University of Konstanz
daniel.keim@uni.kn

Abstract

This paper systematically derives design dimensions for the structured evaluation of explainable artificial intelligence (XAI) approaches. These dimensions enable a descriptive characterization, facilitating comparisons between different study designs. They further structure the design space of XAI, converging towards a precise terminology required for a rigorous study of XAI. Our literature review differentiates between comparative studies and application papers, revealing methodological differences between the fields of machine learning, human-computer interaction, and visual analytics. Generally, each of these disciplines targets specific parts of the XAI process. Bridging the resulting gaps enables a holistic evaluation of XAI in real-world scenarios, as proposed by our conceptual model characterizing bias sources and trust-building. Furthermore, we identify and discuss the potential for future work based on observed research gaps that should lead to better coverage of the proposed model.

1 Introduction

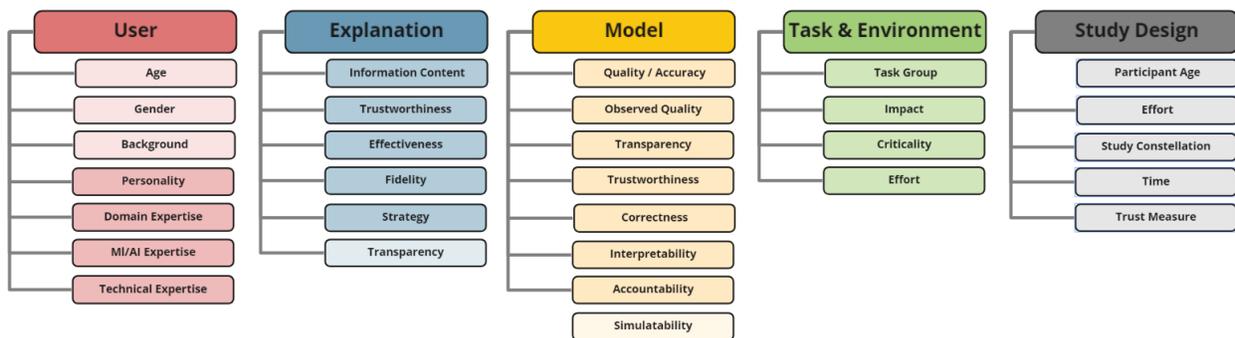


Figure 1: The design dimensions of (X)AI evaluation can be organized in the five groups *user*, *explanation*, *model*, *task & environment*, and *study design*. The individual dimensions have been collected from related work and are defined in section 4.

With the recent leaps in the performance and utility of Artificial Intelligence (AI) across a variety of applications, the demand for understanding their decision-making rationale is on the rise. *Explainable Artificial Intelligence* (XAI) is the study of making the decision-making processes of AI models explainable. Explanations not only can help foster trust among novice users but are also valuable tools when *discovering*, *improving*, *controlling*, or *justifying* [2] the machine learning models powering AI. Consequently, many different approaches to explaining AI have emerged in recent years [1]. XAI encompasses explanations throughout the whole *process* of machine learning from the raw data to presenting the discovered relations and patterns to the user. Within this process, XAI methods focus on explaining the data, the AI model, or presenting the output of the XAI method to the user [64].

A central task for XAI is calibrating trust in the context of complex machine learning models and processes that are not always intelligible. The fact that it is often difficult for humans to comprehend the inner workings of models raises many questions towards methods claiming to provide explanations: Are they valid? Do they calibrate user trust appropriately, or introduce bias? For which data types and tasks are they applicable, and in which environments? In focused application contexts, some approaches can evaluate the general propagation of these effects to derive suitable architectures. However, in the general context of evaluating (X)AI and finding answers to these questions requires comparing different ways to generate, design, and present explanations to different user types. While catering explanations to their intended audience concerning, for example, complexity and information density might seem straightforward, it is not always easy due to the complexity of various psychological processes. Additionally, Miller et al. state that “most of us as AI researchers are building explanatory agents for ourselves, rather than for the intended users” [48] and that XAI is more likely to be successful “if evaluation of these models is focused more on people than on technology.” [48]

Doshi-Velez and Kim [12] identify that interpretability is mostly evaluated in the context of a concrete application, or assumed to be “given” thanks to the use of a particular model class, stating that “To large extent, both evaluation approaches rely on some notion of “you’ll know it when you see it”” [12] According to Guidotti et al. [22], “In the state of the art a small set of existing interpretable models is recognized: decision tree, rules, linear models. These models are considered easily understandable and interpretable for humans” [22]. Poursabzi-Sangdeh et al. [55], however, evaluated such white-box linear models and found that transparency can be overwhelming, possibly due to information overload.

In this paper, we attempt to generate a comprehensive overview of the design dimensions for structured experimental evaluation of XAI methods. To that end, we contribute a literature review of the past five years of (X)AI evaluation research in the human-computer-interaction and visualization communities and report on their considerations for XAI studies and their design. In addition, we review a number of application papers that do not feature comparative study designs, but present user evaluations showcasing the effectiveness of XAI methods. From this review, we generalize and contribute five groups of design dimensions for future evaluation studies: *personal characteristics*, *explanation*, *model*, *task and environment*, and *study setup*. From this literature review we derive not only research gaps, but also opportunities and implications for future work. Further, we contribute a dependency model outlining all actors that must be considered when designing experiments evaluating trust in (X)AI. This model also describes the processes of bias propagation and trust building in (X)AI.

2 Background and Related Work

Already in 1994, Muir presented “a theoretical model of human trust in machines.” [52]. She collected multiple definitions of trust in humans, showcasing the different aspects of trust and suggesting that “trust is a multidimensional construct” [52] and that calibration and re-calibration of trust are necessary over time. More recently, Robbins surveyed psychological literature on trust, finding it to be fragmented with multidimensional definitions. He instead suggest defining trust in terms of “actor A’s beliefs, actor B’s trustworthiness, the matter(s) at hand, and unknown outcomes.” [59]

In addition to these theoretical considerations on trust, Hancock et al. [23] provide a meta-analysis on factors influencing the trust of humans in robots. This work was later extended to automation in general by Schaefer et al. [61]. Neither of the works places emphasis on how these factors can be evaluated. Their dimensions might, however, also be relevant for trust of humans in AI and can be grouped similarly to those presented in section 4.

Recent work from computer science provided a conceptual framework for designing XAI [70]. The authors reviewed work from psychology and philosophy and suggested how XAI should be designed to avoid cognitive biases. These guidelines are applied and evaluated in an application for clinical decision making. More generally, Doshi-Velez and Kim [12] provide a taxonomy for the evaluation of interpretable machine learning systems, identifying three types of studies: application-grounded, human-grounded, functionally-grounded. They discuss when each type of study might be most appropriate, but do not elaborate on the individual design dimensions of the different types. Hoffman et al. [24] provide four criteria for XAI evaluation: *goodness*, *satisfaction*, *comprehension*, and *performance* and focus on how these criteria can best be measured.

3 Literature Review

Many different buzzwords have been mentioned as goals for (X)AI in previous work: intelligibility, justifiability, or interpretability, to name just a few. However, it is not always immediately obvious how these goals can be achieved, and how success can be measured. In order to identify dimensions that have previously been evaluated and to distill guidelines for the evaluation of (X)AI, we conducted a literature review.

Publication	User				Explanation					Model					T & E				Misc.					
	Domain Expertise	ML/AI Expertise	Tech. Expertise	Personality	Available?	Information Content	Trustworthiness	Effectiveness	Fidelity	Strategy	Agnostic?	Quality / Accuracy	Observed Quality	Transparency	Trustworthiness	Correctness	Interpretability	Accountability	Task Group	Impact	Criticality	Effort	# Participants	Venue
Cai et al. (2019) [6]					✓	■					✓								X	0	1	1	1070	IUI
Cheng et al. (2019) [9]	■	■			■						X				■				A	3	3	■	202	CHI
Dodge et al. (2019) [10]					✓	■				■	✓				■				A	1	1	1	160	IUI
Dominguez et al. (2019) [11]					✓	■					✓	■			■				A	1	1	■	121	IUI
Eiband et al. (2019) [14]					✓	■	■	■	■	■	✓				■				A	1	1	1	30	CHI EA
Kouki et al. (2019) [31]	■			■	✓	■		■			✓	■		■	■				X	1	1	1	198	IUI
Millecamp et al. (2019) [46]	■		■		■						✓	■			■		■		A	1	1	1	71	IUI
Richter et al. (2019) [58]					✓			■			✓								A	■	■		65	IUI
Schaffer et al. (2019) [62]		■			■		■				✓								X	1	1	1	551	IUI
Springer and Whittaker (2019) [65]					■						✓			■	■				X	1	1	1	74	IUI
Yin et al. (2019) [72]					X						X	■	■		■				S	1	■	1	1994	CHI
Yin et al. (2019) [72]					X						X	■	■		■				S	1	■	1	757	CHI
Yin et al. (2019) [72]					X						X	■	■		■				S	1	■	1	1042	CHI
Zhou et al. (2019) [76]					■				■		X				■				D	1	1		22	CHI EA
Bigras et al. (2018) [3]					✓	■					✓				■				A	3	2	■	20	CHI EA
Kleinerman et al. (2018) [30]					✓	■					✓		■	■	■				X	1	1	1	59	RecSys
Rader et al. (2018) [56]	■				✓						X					■	■	■	A	1	1		681	CHI
Yu et al. (2017) [73]					X						X	■			■	■			A	2	2	1	21	IUI
Chang et al. (2016) [8]					✓	■	■				✓								X	1	1	1	220	RecSys
Kizilcec (2016) [29]					✓						X		■	■					X				103	CHI
Musto et al. (2016) [53]					✓	■	■				X								X	1	1	1	308	RecSys

Table 1: Synthesis of the most important dimensions mentioned in previous work on (explainable) artificial intelligence. Little squares indicate that a variable was artificially manipulated to a *fixed* ■ value, *measured* ■, constituted a *condition* ■, or a combination thereof ■.

3.1 Methodology

Scope We have collected the proceedings from high-quality computer science journals and conferences. We include conference papers from ACM Computer-Human Interaction (CHI), ACM Intelligent User Interfaces (IUI), ACM Recommender Systems (RecSys) and IEEE VIS (VIS). Additionally, we include journal articles from IEEE Transactions on Visualization and Computer Graphics (TVCG) and Extended Abstracts published at CHI. Furthermore, we retrieved all publications from the International Conference on Machine Learning (ICML) and the ACM Conference on Fairness, Accountability, and Transparency (FAT*). For all venues, we considered the years 2015 to 2019 (2018 and 2019 for FAT*) to focus on recent developments.

Paper Selection Once we had gathered the proceedings, we performed a keyword search for *trust*, *interpretable*, *interpretability*, *explanation*, *explainability*, *transparency* and *interactive machine learning* on the titles and abstracts of published works, retrieving an initial set of papers. We manually evaluated all potential papers of interest and excluded those that do not deal with some form of machine learning or artificial intelligence, or that do not perform user evaluation. For that reason, we excluded all papers from the FAT* and ICML from this review. Due to the large

amount of papers, papers that did not include interactivity and those that covered relatively fewer dimensions were also excluded from our review.

Coding We coded all papers in an iterative process and began with an initial set of eight randomly selected papers. After extracting all relevant dimensions and coding the initial paper set, we distilled and refined coding guidelines until an agreement between coders was reached. We then continued coding the remainder of papers with a single coder. Whenever we encountered new potential dimensions that had not been mentioned in papers previously coded, we conferred and decided whether to include them into adapted coding guidelines. During the coding process, significant differences between pure application papers and those with comparative study designs became apparent. We consequently decided to code application papers using separate guidelines (created using the same process) and present both types of papers separately in the following sections. This methodology allows for a more focused comparison of papers from a given paper type.

Concept Definitions During paper coding, we did not attempt to resolve potential conflicts, ambiguities, or overlaps between concept definitions but coded them as presented by the authors. As a consequence, the results of our literature review present a “union” of the definitions for concepts like trustworthiness or interpretability. Refining these concepts and converging to a common vocabulary presents an opportunity for future work that will be elaborated on further in section 6.

Presentation The results of our literature review are summarized in Table 1 (comparative study designs) and Table 2 (application papers). The tables highlight four different groups of design dimensions for structured experimental evaluation and sort them according to our trust building model introduced in section 5. **Personal** contains both standard personal characteristics, as well as dimensions on experience. **Explanation** and **Model** group dimensions of the respective elements in the XAI pipeline. **Task & Environment** focuses on the implications of using a given (X)AI system in a specific environment. Due to space constraints we do not include all dimensions that have been mentioned in literature. For example, *controllability* and *truthfulness* that were both mentioned only once were excluded. Instead, we focus on the most common dimensions and those that allow us to draw conclusions about the state of the field.

The tables also highlight the number of study participants, as well as the publication venue of the papers. For all dimensions, coloured boxes indicate whether they were *study conditions* , *measured*  in a study, or *fixed*  to an artificial value. Cases where dimensions varied as conditions where also measured,  is used. ✓ and ✗ indicate yes and no respectively, and are used to show whether explanations were available and whether the system was model agnostic from the point of view of the study participant. Possible values for task groups, as well as impact and criticality will be introduced in subsection 4.4. Darker color (0 — 5) indicates higher task complexity, impact or criticality, respectively. Similarly, we classify effort and user expertise that the system in an application paper was designed for as low , medium  or high . Some application papers claim that the presented systems are designed with a specific goal or property in mind, but do not evaluate their respective design decisions and are highlighted accordingly .

3.2 Comparative Studies

Table 1 contains 21 studies from 19 publications. While most publications are only present in the table once, Yin et al. (2019) [72] provide three large-scale studies on the same subject and have thus been included three times. In the remainder of this section we briefly summarize the main findings of our literature review.

3.2.1 Summary of Findings

17/21 studies include *explanations* in their study design. Out of these 17, the availability or absence of explanations is a study condition. Most work only evaluates perceived *trustworthiness of the machine learning model* (15/21) but not the *trustworthiness* of the explanation (3/17 papers that include explanations). While this evaluation of trust in the model is essential, we note a distinct lack of evaluation of the trust in the model explanation. Such explanation evaluations are important in the light of trust-building and bias propagation, as modeled in section 5.

The inclusion of *expertise* as a measured dimension appears to be a relatively recent development, with 7/8 studies having been published in 2019. Furthermore, the only reviewed study incorporating *personality* traits was published in the same year. Knowledge about such user details should influence the *information content*, the most utilized design dimension from the explanation group (9/17). This emphasizes the vast opportunities for presentation of explanations, including varying the level of detail or adding personalization.

Only two studies investigate manipulating the fidelity of an explainer. Worryingly, Eiband et al. [14] find little difference in trust towards real or placebic explanations. Similarly, there seems to be little distinction between the reported understanding of participants, and real understanding that is proven through, for example, little tasks and quizzes. This should inspire future research in that direction to avoid misleading users and miscalibrating their trust.

3.2.2 Discussion of Findings

Many of the experiments evaluate trust in recommendations or social feeds. Those experiments mostly feature low impact and low criticality, making them appropriate for non-expert users. Nonetheless, evaluations of trust in higher-impact settings are needed, especially considering the typically higher criticality of expert-user applications (see Table 2). Future studies are needed that draw from related work from psychology to adequately simulate scenarios that are more appropriate to real-world usage, instead of measuring individual variables in isolation. Better simulation of actual usage conditions and environments is likely to affect study results, especially when impact and criticality are high.

As mentioned above, almost none of the studies evaluate the impact of explanation fidelity. While fidelity is arguably important for interpretability and trust building, its necessity varies depending on the target audience. In expert systems, explanations highlight the models decision-making processes and can uncover training issues or biases. Here, it is essential that all explanations be high-fidelity and follow the inner workings of a given model closely. For the explanation of social feeds or movie recommendations for casual users, however, explanation by example might be more intuitive and effective. As education about the model is not the primary goal, designers have more freedom when creating explanations. Nonetheless, ethical questions remain as wrong explanations can easily mislead users.

In a similar direction, background, age and gender of studies are not always reported, especially when they were conducted through online crowd-sourcing platforms. In addition to user expertise, these dimensions are likely to have a significant impact, though. In particular it is not clear how well studies conducted exclusively with participants from the US generalize to other user groups with large cultural differences. Such cultural differences are also likely to influence optimal information content of explanations. For example, previous work has found differences in the preference for personalized explanations depending on their cultural background [19].

3.3 Application Papers

3.3.1 Summary of Findings

Application papers, almost by definition, describe the design of an XAI model and any accompanying evaluations. Only one paper (Brooks et al. [4]) makes the availability of the explanation a study condition. Most papers assess the explanations by applying them to a dataset as case studies or proof-of-concept demonstrations. Furthermore, most papers do not conduct any significant amount of testing on the explanation itself. For example, only one paper, (Ming et al. [49]), discussed the fidelity of the explanation developed.

The systems presented in the reviewed application papers tend to support more complex tasks (such as model refinement or comparison) than the ones evaluated in the reviewed comparative studies. More than half of the papers also designed explanations for users with high ML expertise (18/35). This ties in to the low impact of most of these papers, since the explanations will necessarily only be relevant to machine learning experts rather than a broader demographic of users. Interestingly, accountability seems to be a more recent trend in XAI. Only one paper (Cabrera et al. [5]) discussed accountability in terms of fairness and mitigating bias.

3.3.2 Discussion of Findings

Many of the reviewed application papers mention particular dimensions of (X)AI, such as trustworthiness, as design goals. However, these dimensions are rarely evaluated for in any user testing or case studies included in the papers. Without such evaluation, it would be harder to verify that the design criteria were indeed satisfied by the system created and that, for example, users indeed found the system to be trustworthy. As a consequence it would be difficult, going forward, to propose a set of guidelines for how (X)AI systems can be designed to meet certain criteria better. The same is true for user evaluations that are performed with a particularly low number of participants.

Finally, as mentioned above, many of the explanations presented in application papers are designed to be used by machine learning experts. In particular when the model being explained is used in the field of deep learning the machine learning experts using the explanations are often considered to be domain experts as well, regardless of the actual data domain. This suggests a potential area of research into designing explanations geared towards machine learning novices or individuals with different domain expertise.

Publication	User				Explanation							Model					T & E				Misc.				
	Domain Expertise	ML/AI Expertise	Tech. Expertise	Personality	Available?	Information Content	Trustworthiness	Effectiveness	Fidelity	Strategy	Iterative	Agnostic?	Quality / Accuracy	Observed Quality	Transparency	Trustworthiness	Correctness	Interpretability	Accountability	Task Group	Impact	Criticality	Effort	# Participants	Venue
Brooks et al. (2015) [4]		2				■						✓	■	□				□		R	2	1	■		VIS
Cabrera et al. (2019) [5]	3	3	3		✓							✓	□						□		D	1	2		VIS
Cavallo and Demiralp (2019) [7]	■	■	■		✓					✓	✓						■	■		U	2	1		12	VIS
El-Assady et al. (2018) [15]	■	■	■		✓					✓	✓	■	□		■			□		R	2	1	3		6 VIS
El-Assady et al. (2019) [18]	■	■	■		✓					✓	✓	■	□	□		■				C	2	1	1		6 VIS
El-Assady et al. (2019) [17]	■	■	■		✓					✓	✓	■	□		■					R	2	1	1		6 VIS
Hohman et al. (2019) [26]		3			✓							×						□		U	2	1			VIS
Kahng et al. (2018) [27]	3	3	3		✓							×						□		R	2	1			VIS
Kahng et al. (2019) [28]		1			✓					✓	×							□		S	3	2	1		VIS
Krause et al. (2017) [32]	3	3			✓							×	□	□						R	2	1	1		VIS
Kumpf et al. (2018) [33]	3				✓					✓	✓		□							S	3	2	2		VIS
Kwon et al. (2018) [34]	3	3			✓					✓	✓									A	2	1			VIS
Kwon et al. (2019) [35]	3		3		✓					✓	✓	■		□				□		S	3	4			VIS
Lin et al. (2018) [37]		3			✓							×	■							D	2	1	2		VIS
Liu et al. (2017) [40]		3			✓					✓	×	■								R	2	1			VIS
Liu et al. (2018) [42]		3			✓	□						×								D	2	1			VIS
Liu et al. (2018) [39]	3				✓					✓	×				■					A	2	4		14	VIS
Liu et al. (2018) [44]		3			✓					✓	×	■								D	2	1			VIS
Liu et al. (2018) [43]		3			✓		■	□				×						□		D	2	1			VIS
Liu et al. (2018) [41]		3			✓							×	■							D	2	1			VIS
Ma et al. (2019) [45]		3			✓							×	■							D	2	1	3		VIS
Ming et al. (2017) [49]		3			✓		■					×	■		□		□			D	2	1	1		VIS
Ming et al. (2019) [50]	3				✓			□				×	■		□			□		R	2	1		9	VIS
Muhlbacher et al. (2018) [51]	3				✓					✓	✓	■	□					□		R	2	1			VIS
Pezzotti et al. (2017) [54]		3			✓					✓	×	□						□		D	2	1			VIS
Ren et al. (2017) [57]	■	3			✓							×	■	□						C	2	1	1	24	VIS
Sacha et al. (2018) [60]		3			✓					✓										A	2	1			VIS
Spinner et al. (2019) [64]	■				✓		■			✓	✓							□		R	3	1	3	9	VIS
Stahnke et al. (2016) [66]	■				✓		□	■				×	□							U	2	1	3		VIS
Stoffel et al. (2015) [67]		3			✓					✓	✓									R	2	1			VIS
Strobelt et al. (2018) [68]		3			✓							×						□		D	3	1			VIS
Strobelt et al. (2019) [69]		3			✓					✓	✓									R	2	1			VIS
Wang et al. (2019) [71]		3			✓					✓	✓	■								R	2	1			VIS
Zhang et al. (2019) [74]		3			✓					✓	×	■		□				□		R	2	1			VIS
Zhao et al. (2019) [75]		2			✓		■					×		□				□		U	3	1	1		VIS

Table 2: The most important dimensions mentioned in previous application work on (explainable) artificial intelligence. Little orange squares indicate that a system was designed with a goal or property of (X)AI in mind, and that the property was evaluated ■ or not evaluated . *Measured* ■ variables and experimental *conditions* ■ are also shown.

4 Design Dimensions for Experimental Evaluations

In this section, we synthesize design dimensions for the structured experimental evaluation of explainable artificial intelligence from the literature review presented above. Where possible or necessary, we provide definitions. As many goals and properties of (X)AI have been defined in the literature but were not yet evaluated in the reviewed literature, we expand the dimensions with these definitions. All dimensions reuse the colors from Figure 1. Higher opacity indicates dimensions that appear in Table 1 or Table 2, while lower opacity is used for all remaining dimensions.

Previous work from Doshi-Velez and Kim [12] characterizes (X)AI along the dimensions of *global* and *local interpretability*, *time limitation* and the nature of *user experience*. [12]. We have also identified those dimensions from our literature review and report them below. Guidotti et al. [22] have identified *reliability*, *robustness*, *causality*, *scalability*, and *generality* as desired dimensions for machine learning models [22]. As these are high-level concepts, they have not yet been experimentally evaluated. They can, however, likely be approximated by the dimensions reported below.

4.1 User Attributes

As our literature review emphasized the evaluation of trust in (X)AI, user attributes play an important role. This design dimension is characterized by the question: **Who was the (X)AI method designed for?** Within the dimension, we distinguish between immutable personal characteristics and personal experience that is dependent on the circumstances. When designing experiments that modulate these dimensions, researchers can draw from extensive related work from psychology and the humanities on trust-building, explanation processes, and conversational explanations.

Personal Characteristics Personal characteristics are immutable.

- **Age** — The age group that the tool or system was designed for. We did not encounter work specifically targeting a certain age group. This provides opportunities for future research, for example in evaluating trust of teenagers in social media recommendations.
- **Gender** — The gender that the system was designed for. Typically, we expect systems to be designed for fifty percent female users. Some studies report gender-specific trust measurements. [30]
- **Background** — Cultural differences have a large influence on how we cooperate with peers, including machines, and how likely we are to follow or reject recommendations. Previous work from psychology found significant differences between some groups, but not others [19]. Opportunities for future work include verifying whether these findings transfer to trust in (X)AI.
- **Personality** — A dimension that is only mentioned in few studies and likely correlates with information captured by the *background* dimension. Personal characteristics of interest include, among others, the propensity to trust, differences between trust in humans and machines, prejudice built from previous experience, confidence or self-esteem.

Experience Objective assessment of experience is a challenging task, not only due to the Dunning-Kruger effect [13] causing non-experts to be notoriously bad at rating their own experience. Instead, study designs should rely on asking questions about the number of years in a given field and testing participants' knowledge with questions. Participants can then be classified as *novice*, *intermediate*, *proficient*, or *expert* to make studies more easily comparable.

- **Domain Expertise** — The visualization community often considers *expert* domain knowledge from, for example, medicine, linguistics, or biology. However, “casual users” also have relevant domain knowledge, for example, in music or movies. Studies should investigate whether there are significant differences in trust between these two user groups and whether results from one are directly informative for designs targeting the respective other group.
- **Technical Expertise** — Technical experience includes general familiarity with computers or automation, as well as awareness of potential issues that may arise. Users that are more familiar with technology are generally expected to be more proficient at using (X)AI systems.
- **ML Expertise** — More specific than technical expertise, machine learning expertise is concerned with the familiarity with and understanding of the specific machine learning algorithms used.

4.2 Explanations

Depending on the user group, explanations might be necessary or not. If they are presented to users, care has to be taken to calibrate trust and avoid biases. This dimension is thus characterized by these questions: **Are system decisions explained? If so, how?** This group of dimensions can draw from a significant body of related work from social sciences.

Not only do social sciences provide models of explanation, they also characterize expectations towards the explanation process [47].

- **Availability** — Many study designs include conditions without explanations as baselines. This is especially important when it is unclear whether explanations have an influence on a given variable or dimension.
- **Information Content** — Once the general usefulness of explanations in a given scenario has been demonstrated, study designers have vast opportunities in varying the information content of explanations. Subdimensions include, among many others, the information density, personalization of explanations, the use of emotional or factual statements.
- **Trustworthiness** — The trustworthiness assigned to explanations by study participants. The trustworthiness of an explanation can be explicitly affected by manipulating the correctness of the explanation or, more subtly, the tone in which it is presented.
- **Effectiveness** — Some studies measure the effectiveness of an explanation. This dimension is mostly used to capture the convincingness of an explanation to perform a given action, not in explaining some complex underlying theory.
- **Fidelity** — This dimension captures whether, and how well, an explanation actually explains the models' decision-making process (high fidelity), or just contains some information that is presented in the style of an explanation but does not correspond to the model in any way (low fidelity). High fidelity explanation methods are fundamental for effective (X)AI.
- **Strategy** — Three major reasoning strategies are known from the social sciences and used in (X)AI: inductive (example-based), deductive (theory-based), and abductive (inductive reasoning in the absence of all facts; iterative process once more knowledge becomes available) reasoning [16].
- **Transparency** — An explanation method is transparent when all its decision making processes can be observed and understood by users. While early work in (X)AI equated transparency with an explanation, later work found that transparency might be overwhelming [55].

4.3 Models

So far, dimensions have focused on users and model explanations. This group of dimensions characterizes models in detail and answers the question: **Which AI models are used in this process?**

- **Agnostic** — Some systems are model-agnostic from the point of view of the study participant in the sense that the users do not know what model is powering the system if there is one at all. Designing a model-agnostic system or providing model details to the participant has implications for user awareness, primes them by setting expectations, and thus influences trust and biases.
- **Observed Quality** — The quality of the model that users interact with, typically represented in terms of the accuracy that users could observe during the study on the actual data points used.
- **Quality** — The actual quality of the model that users interact with. This quality is typically represented by the *accuracy* measured on the held-out test data. Showing this number to study participants before or during the study sets their expectations, starting a new trust calibration process whenever observed quality and presented model quality differ.
- **Transparency** — A model is transparent when all its decision-making processes can be observed and understood by users. While early work in (X)AI equated transparency with an explanation, later work found that transparency might be overwhelming [55].
- **Correctness** — This dimension describes user-perceived correctness (as opposed to model quality) in terms of how well the system output aligns with users' expectations [56].
- **Interpretability** — We define a system to be interpretable when users can understand why it behaves in a given way under given circumstances. In that sense, interpretability can be considered an inductive process, where users first create a mental model of the system and then verify whether the system is consistent with that mental model, making it interpretable. Lipton [38] has previously surveyed interpretability and suggests “that interpretability is not a monolithic concept, but in fact reflects several distinct ideas.”
- **Accountability** — Accountability “measure[s] the extent to which participants think the system is fair and they can control the outputs the system produces.” [56]
- **Trustworthiness** — This high-level dimension is based on multiple other dimensions. A model can be considered trustworthy when it is correct (according to user beliefs) and interpretable.
- **Simulatability** — A model is simulatable when users can successfully predict the model output for a given input.

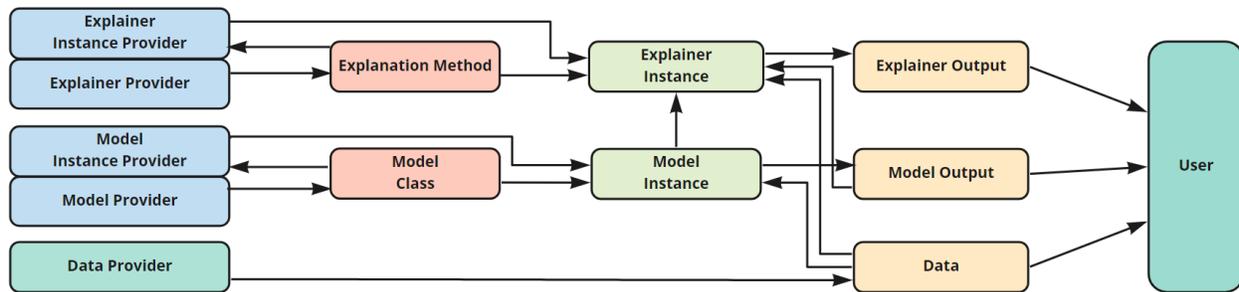


Figure 2: Dependency Model for the XAI process. Bias propagates along the arrows, while trust is built based on the user’s interaction with the data, model, and/or explainer outputs, respectively, following the dependency arrows in reverse.

4.4 Tasks and Environment

This group captures dimensions answering the question: **How and where are models and explanations used?** The dimensions characterize typical user tasks and the costs of (mis-)using an AI system. The cost is a combination of *impact* and *criticality*. We introduce some examples of different impact-criticality combinations after defining the dimensions.

- **Task Group** — We identified seven general tasks with increasing difficulty: *application* ^A, *understanding* ^U, *simulation* ^S, *diagnosis* ^D, *refinement* ^R, *justification* ^J, *comparison* ^C. On the one hand, these groups simplify comparing system designs. On the other hand, it is apparent that users simply applying a machine learning model have different requirements towards explanations than those who have to diagnose problems, or justify decisions. Furthermore, previous work suggests that users are more likely to accept recommendations when working on complex tasks. [21]
- **Impact** — The amount of people impacted by decisions supported through (X)AI systems likely has an influence on user behaviour: the more people are affected the less risk is acceptable. We define five categories for characterizing impact: *none* ⁰, *one* ¹, *some* ², *many* ³, and *all* ⁴.
- **Criticality** — Criticality reports how severe the influence of an (X)AI system can be on those impacted. Possible values are *none* ⁰, *marginal* ¹, *significant* ², *troublesome* ³, *livelihood* ⁴, *extreme* ⁵.
- **Effort** — Independent of the impact and criticality, human actions motivated by (X)AI systems require a certain amount of effort when executed. We classify this effort as *low* ¹, *medium* ² or *high* ³. Related work from psychology shows significant differences in willingness to act depending on the effort required [36].

A few years ago, a bug in the voice assistant Alexa caused smart speakers in many people’s homes to play laughs randomly. We classify this incident as high-impact, low-criticality as it affected many people but caused no harm. In contrast, the failure of the autopilot of a car or plane would be high-criticality, and low- or high-impact, respectively.

4.5 Study Design

In addition to model- and user-specific dimensions, evaluation studies are highly dependent on the study design and setup. This section thus groups dimensions answering the question: **How was the study designed?**

- **Participant Age** — The age of the study participants. This age can differ from the age group that a system was designed for.
- **Exerted Effort** — The effort that was needed to complete the study. In order to produce reliable results that can be generalized, the gap between the effort needed in actual use and the effort exerted under study conditions should be minimized.
- **Study Constellation** — Under study constellation, we summarize all variables like the number of participants completing the study in parallel, whether participants were intentionally disturbed or distracted to create the desired effect, or how much help was available, for example, in pair analytics sessions. This dimension offers great potential for creating realistic study settings that replicate real-world usage conditions.
- **Time** — Many (X)AI systems are used under time pressure in day-to-day operations. Consequently, evaluation studies need to be run under realistic time limits to create a comparable environment.

- **Trust Measure** — There are different ways to measure trust. Poursabzi-Sangdeh et al. [55] use simulation and weight of advice for prediction tasks, but different systems might require different trust measures.

5 Structuring the Design Space

To generate a holistic view on the evaluation of XAI we strive to bring the design dimensions into context, structuring the design space of XAI studies and defining their scope and influences. In this section, we present a dependency model for XAI processes. This model describes the different stages and stakeholders of XAI. Each of the dimensions detailed in section 4 has a different impact on the model’s components. In particular, our contextual model can be utilized to describe biases that might arise within XAI processes, as well as bias propagation through the dependencies. In addition, we postulate that trust building occurs through the user’s interaction with the data, AI model, and explanations. Hence, trust building follows the dependency arrows of our model in reverse order. This section emphasizes that the aim of XAI studies should be to achieve a broad coverage, while not increasing the complexity of the modeling or losing too much detail. In the following, we describe the dependency model in more detail, and discuss the processes of bias propagation and trust building in (X)AI.

5.1 Dependency Model

The proposed, conceptual dependency model in Figure 2 covers the stakeholders in XAI systems and the building blocks they provide. This model highlights *dependencies* and is not designed to model possible *interactions* or *iterative feedback loops*. In the following, we describe our dependency model as depicted in Figure 2. We simply refer to AI/ML models as “*model*” in the remainder of this section.

The ***model provider*** is a person or entity creating a novel ***model class***. Such a model class is subject to the design dimensions presented in subsection 4.3. A concrete ***model instance*** is needed for the model to be practically used. Such an instance is created by the ***model instance provider***. In addition to a model implementation (omitted in the model to avoid unnecessary complexity unrelated to bias and trust propagation) the instance provider typically requires some ***training data***. This data comes from the ***data provider***. Note that all stakeholders in the model might be the same entity, or all be distinct. Once the model is trained, it can produce some ***model output***. Together with the training data and the model instance, this output forms the potential inputs for a model ***explainer instance***. Which inputs are actually used depends on the type of the explainer [64]. Analogous to models, explainers are particularly influenced by dimensions from subsection 4.2, and have an associated ***explainer provider***, ***explainer method*** and an ***explainer instance provider***. Depending on the system design, ***explainer output***, model output and training data, data might be available to the user, who is characterized by dimensions from subsection 4.1.

5.2 Bias Propagation

Despite systems and explanations being designed as deliberately as possible, they are still subject to external factors like where or by whom a machine learning model has been trained or deployed. The dependency model contains many stakeholders with potentially diverging goals and interests: a data provider might discriminate against foreigners for political reasons, a model instance provider against minorities, and a particular explainer might only be useful to experts. Whether it is willingly or unwillingly, such biases might hamper the trustworthiness of the complete (X)AI pipeline.

As these potential biases *propagate* through the XAI process, we use the dependency model to describe their influences. For example, it is impossible to obtain a fair and unbiased, high-quality model if the training data had a racial or gender bias. Increasing the transparency of a system, for example through explanations, can help to reveal such biases; however, it does not reduce them. On the contrary, explanations might miscalibrate user trust in a system and lessen bias awareness. This is even true for high-fidelity explanations that correctly represent the model’s decision-making process: if the model itself is insufficient, any local explanation can itself be correct, while still misleading users and not revealing model shortcomings. We, therefore, argue for acknowledging sources of potential biases and their effects on other stages and stakeholders in the XAI process using the dependency model. Mitigating such biases in-place and limiting their propagation through the model can reduce their harmful effects, as mentioned above.

In addition to those general biases that are present for all users, the (X)AI process is subject to user-specific biases. Those biases are based on users’ previous experiences and their knowledge. In our model, biases can be added to the process at any block in Figure 2. There, they will *increase* the existing biases and propagate along the depicted dependency arrows. The effects of bias propagation in (X)AI are an interesting area of future research and studies.

5.3 Trust Building

Humans usually do not build trust in abstract concepts, but concrete outputs that they see and can potentially interact with. Consequently, users first recognize that a model seems to work well and that the explanations seem to make sense. Once they have built trust in the explanations and model outputs, they start building trust in the model, before eventually trusting the model instance provider. As users propagate back through the XAI process, their personal biases apply, influencing trust building positively or negatively.

For example, a particular user might have experienced unreliability using deep learning classifiers resulting in a personal conviction that such models do not work. This person might also have had a good experience using models provided by a particular tech company, resulting in a positive personal bias towards the model provider. If in our case, this user provides their own data that they are familiar with, they can judge the trustworthiness of a model output given their expectations about the data. They can probably also judge the fidelity of an explainer based on its output. Using our dependency model, we can follow the trust-building process, taking into account positive and negative amplifications that are reinforced through biases. If the explanations provided by the XAI method are matching the user's expectation, they might foster trust in the explainer output and start building trust in the explainer instance. This effect might be reinforced by other factors, such as a positive experience using the same explainer instance on a different data point. Increasing the trust in the explainer instance might lead to the user starting to trust the explainer instance provider and/or the explanation method itself.

To answer the question of whether or not we should trust (X)AI, we have to take the cross-relational effects of the dependency model into account. In particular, when designing evaluations for (X)AI systems, we have to consider the coverage of the model. In our literature review, we broadly observed a disconnect between the different research communities; with HCI focusing more on trust-building in the presentation of X(AI), and the AI/ML community, generally concerned with the correctness of the models as a means to increase trust. We argue that a broad coverage of the XAI process is necessary, as observed in some application papers (within their focused scope). Moreover, we postulate that studies concerned with one part of the dependency model should abstain from partly including descriptions of other parts of the model without considering possible dependencies and cross-effects. Coming back to our example from above this means that a study participant should not be informed about the model providers (the tech company) if the study design is not set up to appropriately capture potential biases, their dependencies, and the resulting effects on trust.

6 Discussion and Implications

The model presented in the previous section has direct implications for study design: any components of the model that are mentioned in the study prototype, questionnaire, or some meta-information must be taken into account as potential sources of bias, distorting results. At the same time, they highlight the vast opportunities for future work, conducting studies that include or exclude those particular areas of the model. In the following, we will present further opportunities, as well as limitations of our work.

6.1 Opportunities

(1) Strive Towards Better Coverage of the XAI Process The dependency model highlights stakeholders in the different stages of the (X)AI process and their dependencies. Currently, the different communities tend to focus on different stages of the process when conducting evaluations: the machine learning community does not typically involve users in evaluations, and the HCI community tends to focus on the presentation of model outputs and explanations. Better coverage of the model in the form of studies spanning multiple stages is required to bring the field forward. While application studies (mostly from VIS) attempt to bridge this gap by providing rich model interactions, they are problem-specific and often only gather qualitative feedback. Consequently, there is great potential for collaboration between these communities to provide end-to-end testing of the (X)AI process, explaining the inner workings of machine learning models.

(2) Bridge the Gap to other Communities In addition to better connecting the different communities from computer science, evaluation of (X)AI can profit from collaborations with the social sciences, and previous sections have already occasionally alluded to particular related work from psychology. Significant bodies of work have investigated trust in inter-human relations, and made generalizations towards human-robot [23] and human-automation collaboration [61]. The respective experiments should be repeated to verify that their findings still apply for human-AI collaboration. Additionally, collaboration with psychologist is needed to create study scenarios that more closely resemble real-world usage conditions, rather than mostly relying on online crowd-sourced studies. This is especially important for evaluating use-case specific applications tackling high-impact and high-criticality issues.

(3) Apply a Clearly Defined Terminology We observed a tendency, among some papers that we studied, of stating some systems were designed with specific goals like interpretability in mind, but not evaluating whether the said properties were achieved. More work like the structuring review by Lipton [38] is needed to refine and merge the concepts that have already been proposed. This allows related fields to converge on common terminologies that are well-aligned with each community’s identified goals. Then, instead of defining more high-level goals for (X)AI, researchers would consider how the proposed dimensions can be measured effectively, and what approximations and proxies might be necessary.

(4) Acknowledge Biases and Propagation of Trust Researchers evaluating (X)AI should acknowledge the inherent biases in human trust-building and draw from related work in the social sciences. This should, in particular, influence the presentation of explanations. Results from psychology show that “sets of source factors (expertise, liking, trust, and similarity) and message factors (politeness, response efficacy, feasibility, absence of limitations, and confirmation)” [20] each influence how humans deal with advice. Especially *liking* and *trust* are likely to vary from individual to individual based on existing prejudices and biases.

(5) Trust in Explanations vs. Trust in Models Typical XAI systems aim to increase trust by providing explanations. Success is then often measured by evaluating the trust users have in the system. However, if the explanations are not transparent to the user and cannot be verified, they cannot ease doubts about the correctness of the underlying model. Instead, they simply shift the problem to a different stage of the XAI process by explaining a black box with a black box. Consequently, the trust in explanations and the correct calibration of trust in them should take a more central role in XAI research.

(6) Consider Explainer Fidelity Model complexity, especially of deep neural networks, is quickly increasing thanks to the availability of relatively cheap computing resources. The more difficult it becomes for humans to interpret these complex models, the more difficult it becomes to generate intuitive explanations. At the same time, some work from psychology on placebo explanations has already been successfully replicated in the context of XAI [36, 14], and it has been suggested that humans are eager to believe explanations they are provided with [25]. Consequently, a cornerstone of XAI research should be ensuring that explainers have high fidelity. Otherwise, the field risks producing explanations that “sound good” but are misleading users and exploiting miscalibrated trust.

(7) Incorporate Motivation for Use There can be various reasons for users to be interacting with (X)AI. These reasons can stem from intrinsic or extrinsic motivation, be spontaneous or persistent over time. All these factors influence the perceived impact and criticality of a given task, and consequently, the effort that users are willing to exert. Not only is it important to motivate system users, for example, using gameful design elements [63], but also to create study constellations that reflect those motivations.

6.2 Limitations

To set the scope of our work, in this section, we acknowledge the limitations of our literature review and discuss potential alternatives. First, our review is based on a keyword search on the title and abstract of impactful venues. Future work could extend this review into a survey by expanding on the methodology and including more related work through the inclusion of forward and backward references. As previously elaborated by Lipton [38], not all concepts in XAI are clearly defined. We do not attempt a disambiguation of terms used. Instead, our coding is directly based on the concepts mentioned in the respective reviewed work. Further, we only collect and review relevant papers and do not provide a complex meta-evaluation taking reported significant differences and effect sizes into account. Instead, we focus on the reported dimensions. The dependency framework presented in section 5 is preliminary and focuses on those entities that are important for trust building and bias propagation. Future work should provide extensions covering model implementations (a common source of errors from our experience) and afforded interaction possibilities.

6.3 Future Work

Our literature review revealed the independent goals of various sub-domains in the computer science community. Visual analytics, mostly working with expert systems, is concerned with explanations teaching users the inner workings of models and educating them in machine learning, while recommender systems are often tailored to be convincing. This work presents general design dimensions that are applicable to various domains. Tailoring those dimensions to the specific goals of the domains facilitates converging towards a common, shared vocabulary as introduced above.

As Table 1 revealed, the majority of comparative studies assess the trustworthiness of models. Most of the reviewed works asked study participants to rate the trustworthiness on Likert-scale questionnaires. Instead of adding this step of indirection, Yin et al. [72] utilize simulatability and weight of advice as proxies for trust. Avoiding post-usage questionnaires and relying on implicit trust measures lowers the risk of users’ preconceptions and biases influencing their trust-responses. Future work should investigate which proxies for trust are best suited for which tasks.

7 Conclusion

We have presented a literature review of the past five years of explainable artificial intelligence in the visualization and human-computer-interaction communities. From our review, we have distilled design dimensions for the user-centered evaluation of (X)AI methods and systems. Comparing those design dimensions and goals typically mentioned for (X)AI, we identified research gaps and opportunities for future work. So far, many design dimensions have barely been utilized in evaluations. This is especially true for abstract concepts like interpretability or accountability. We have also presented a dependency model highlighting the different stages of the (X)AI process and showing how bias and trust propagate in (X)AI systems. Together with the design dimensions, this model guides future evaluations of (X)AI systems.

References

- [1] A. Abdul, J. Vermeulen, D. Wang, B. Y. Lim, and M. Kankanhalli. Trends and Trajectories for Explainable, Accountable and Intelligible Systems: An HCI Research Agenda. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, CHI '18, pages 582:1–582:18, New York, NY, USA, 2018. ACM. ISBN 978-1-4503-5620-6. doi: 10.1145/3173574.3174156. URL <http://doi.acm.org/10.1145/3173574.3174156>. event-place: Montreal QC, Canada.
- [2] A. Adadi and M. Berrada. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access*, 6:52138–52160, 2018. doi: 10.1109/ACCESS.2018.2870052.
- [3] E. Bigras, M.-A. Jutras, S. Sénécal, P.-M. Léger, M. Fredette, C. Black, N. Robitaille, K. Grande, and C. Hudon. Working with a Recommendation Agent: How Recommendation Presentation Influences Users' Perceptions and Behaviors. In *Extended Abstracts of the 2018 CHI Conference on Human Factors in Computing Systems*, CHI EA '18, pages LBW098:1–LBW098:6, New York, NY, USA, 2018. ACM. ISBN 978-1-4503-5621-3. doi: 10.1145/3170427.3188639. URL <http://doi.acm.org/10.1145/3170427.3188639>. event-place: Montreal QC, Canada.
- [4] M. Brooks, S. Amershi, B. Lee, S. M. Drucker, A. Kapoor, and P. Simard. FeatureInsight: Visual support for error-driven feature ideation in text classification. In *2015 IEEE Conference on Visual Analytics Science and Technology (VAST)*, pages 105–112, Oct. 2015. doi: 10.1109/VAST.2015.7347637.
- [5] A. A. Cabrera, W. Epperson, F. Hohman, M. Kahng, J. Morgenstern, and D. H. Chau. FairVis: Visual Analytics for Discovering Intersectional Bias in Machine Learning. *arXiv:1904.05419*, Apr. 2019. URL <http://arxiv.org/abs/1904.05419>. arXiv: 1904.05419.
- [6] C. J. Cai, J. Jongejan, and J. Holbrook. The Effects of Example-based Explanations in a Machine Learning Interface. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*, IUI '19, pages 258–262, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6272-6. doi: 10.1145/3301275.3302289. URL <http://doi.acm.org/10.1145/3301275.3302289>. event-place: Marina del Ray, California.
- [7] M. Cavallo and C. Demiralp. Clustrophile 2: Guided Visual Clustering Analysis. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):267–276, Jan. 2019. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2018.2864477. URL <https://ieeexplore.ieee.org/document/8440035/>.
- [8] S. Chang, F. M. Harper, and L. G. Terveen. Crowd-Based Personalized Natural Language Explanations for Recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*, RecSys '16, pages 175–182, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4035-9. doi: 10.1145/2959100.2959153. URL <http://doi.acm.org/10.1145/2959100.2959153>. event-place: Boston, Massachusetts, USA.
- [9] H.-F. Cheng, R. Wang, Z. Zhang, F. O'Connell, T. Gray, F. M. Harper, and H. Zhu. Explaining Decision-Making Algorithms Through UI: Strategies to Help Non-Expert Stakeholders. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, CHI '19, pages 559:1–559:12, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-5970-2. doi: 10.1145/3290605.3300789. URL <http://doi.acm.org/10.1145/3290605.3300789>. event-place: Glasgow, Scotland Uk.
- [10] J. Dodge, Q. V. Liao, Y. Zhang, R. K. E. Bellamy, and C. Dugan. Explaining Models: An Empirical Study of How Explanations Impact Fairness Judgment. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*, IUI '19, pages 275–285, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6272-6. doi: 10.1145/3301275.3302310. URL <http://doi.acm.org/10.1145/3301275.3302310>. event-place: Marina del Ray, California.
- [11] V. Dominguez, P. Messina, I. Donoso-Guzmán, and D. Parra. The Effect of Explanations and Algorithmic Accuracy on Visual Recommender Systems of Artistic Images. In *Proceedings of the 24th International Conference on*

- Intelligent User Interfaces*, IUI '19, pages 408–416, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6272-6. doi: 10.1145/3301275.3302274. URL <http://doi.acm.org/10.1145/3301275.3302274>. event-place: Marina del Ray, California.
- [12] F. Doshi-Velez and B. Kim. Towards A Rigorous Science of Interpretable Machine Learning. *arXiv:1702.08608 [cs, stat]*, Feb. 2017. URL <http://arxiv.org/abs/1702.08608>. arXiv: 1702.08608.
- [13] D. Dunning. Chapter five - The Dunning–Krugger Effect: On Being Ignorant of One’s Own Ignorance. In J. M. Olson and M. P. Zanna, editors, *Advances in Experimental Social Psychology*, volume 44, pages 247–296. Academic Press, Jan. 2011. doi: 10.1016/B978-0-12-385522-0.00005-6. URL <http://www.sciencedirect.com/science/article/pii/B9780123855220000056>.
- [14] M. Eiband, D. Buschek, A. Kremer, and H. Hussmann. The Impact of Placebic Explanations on Trust in Intelligent Systems. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems - CHI EA '19*, pages 1–6, New York, New York, USA, 2019. ACM Press. ISBN 978-1-4503-5971-9. doi: 10.1145/3290607.3312787. URL <http://dl.acm.org/citation.cfm?doid=3290607.3312787>.
- [15] M. El-Assady, R. Sevastjanova, F. Sperrle, D. Keim, and C. Collins. Progressive Learning of Topic Modeling Parameters: A Visual Analytics Framework. *IEEE Transactions on Visualization and Computer Graphics*, 24(1): 382–391, Jan. 2018. ISSN 1077-2626. doi: 10.1109/TVCG.2017.2745080. URL <http://ieeexplore.ieee.org/document/8019825/>.
- [16] M. El-Assady, W. Jentner, R. Kehlbeck, U. Schlegel, R. Sevastjanova, F. Sperrle, T. Spinner, and D. Keim. Towards xai: Structuring the processes of explanations. pages 1–9, 2019.
- [17] M. El-Assady, R. Kehlbeck, C. Collins, D. Keim, and O. Deussen. Semantic Concept Spaces: Guided Topic Model Refinement using Word-Embedding Projections. *IEEE Transactions on Visualization and Computer Graphics*, pages 1–1, 2019. doi: 10.1109/TVCG.2019.2934654.
- [18] M. El-Assady, F. Sperrle, O. Deussen, D. Keim, and C. Collins. Visual Analytics for Topic Model Optimization based on User-Steerable Speculative Execution. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):374–384, Jan. 2019. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2018.2864769. URL <https://ieeexplore.ieee.org/document/8467535/>.
- [19] A. M. Evans and J. I. Krueger. The Psychology (and Economics) of Trust. *Social and Personality Psychology Compass*, 3(6):1003–1017, 2009. ISSN 1751-9004. doi: 10.1111/j.1751-9004.2009.00232.x. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1751-9004.2009.00232.x>.
- [20] B. Feng and E. L. MacGeorge. The Influences of Message and Source Factors on Advice Outcomes. *Communication Research*, 37(4):553–575, Aug. 2010. ISSN 0093-6502. doi: 10.1177/0093650210368258. URL <https://doi.org/10.1177/0093650210368258>.
- [21] F. Gino and D. A. Moore. Effects of task difficulty on use of advice. *Journal of Behavioral Decision Making*, 20(1):21–35, 2007. ISSN 1099-0771. doi: 10.1002/bdm.539. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/bdm.539>.
- [22] R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, F. Giannotti, and D. Pedreschi. A Survey of Methods for Explaining Black Box Models. *ACM Comput. Surv.*, 51(5):93:1–93:42, Aug. 2018. ISSN 0360-0300. doi: 10.1145/3236009. URL <http://doi.acm.org/10.1145/3236009>.
- [23] P. A. Hancock, D. R. Billings, K. E. Schaefer, J. Y. C. Chen, E. J. de Visser, and R. Parasuraman. A Meta-Analysis of Factors Affecting Trust in Human-Robot Interaction. *Human Factors*, 53(5):517–527, Oct. 2011. ISSN 0018-7208. doi: 10.1177/0018720811417254. URL <https://doi.org/10.1177/0018720811417254>.
- [24] R. R. Hoffman, S. T. Mueller, G. Klein, and J. Litman. Metrics for Explainable AI: Challenges and Prospects. *arXiv:1812.04608 [cs]*, Dec. 2018. URL <http://arxiv.org/abs/1812.04608>. arXiv: 1812.04608.
- [25] F. Hohman, A. Head, R. Caruana, R. DeLine, and S. M. Drucker. Gamut: A Design Probe to Understand How Data Scientists Understand Machine Learning Models. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, CHI '19, pages 579:1–579:13, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-5970-2. doi: 10.1145/3290605.3300809. URL <http://doi.acm.org/10.1145/3290605.3300809>. event-place: Glasgow, Scotland Uk.
- [26] F. Hohman, H. Park, C. Robinson, and D. H. Chau. Summit: Scaling Deep Learning Interpretability by Visualizing Activation and Attribution Summarizations. *arXiv:1904.02323 [cs]*, Apr. 2019. URL <http://arxiv.org/abs/1904.02323>. arXiv: 1904.02323.
- [27] M. Kahng, P. Y. Andrews, A. Kalro, and D. H. Chau. ActiVis: Visual Exploration of Industry-Scale Deep Neural Network Models. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):88–97, Jan. 2018. ISSN

- 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2017.2744718. URL <https://ieeexplore.ieee.org/document/8022871/>.
- [28] M. Kahng, N. Thorat, D. H. P. Chau, F. B. Viegas, and M. Wattenberg. GAN Lab: Understanding Complex Deep Generative Models using Interactive Visual Experimentation. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):310–320, Jan. 2019. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2018.2864500. URL <https://ieeexplore.ieee.org/document/8440049/>.
- [29] R. F. Kizilcec. How Much Information?: Effects of Transparency on Trust in an Algorithmic Interface. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, CHI '16, pages 2390–2395, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-3362-7. doi: 10.1145/2858036.2858402. URL <http://doi.acm.org/10.1145/2858036.2858402>. event-place: San Jose, California, USA.
- [30] A. Kleinerman, A. Rosenfeld, and S. Kraus. Providing Explanations for Recommendations in Reciprocal Environments. In *Proceedings of the 12th ACM Conference on Recommender Systems*, RecSys '18, pages 22–30, New York, NY, USA, 2018. ACM. ISBN 978-1-4503-5901-6. doi: 10.1145/3240323.3240362. URL <http://doi.acm.org/10.1145/3240323.3240362>. event-place: Vancouver, British Columbia, Canada.
- [31] P. Kouki, J. Schaffer, J. Pujara, J. O'Donovan, and L. Getoor. Personalized Explanations for Hybrid Recommender Systems. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*, IUI '19, pages 379–390, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6272-6. doi: 10.1145/3301275.3302306. URL <http://doi.acm.org/10.1145/3301275.3302306>. event-place: Marina del Ray, California.
- [32] J. Krause, A. Dasgupta, J. Swartz, Y. Aphinyanaphongs, and E. Bertini. A Workflow for Visual Diagnostics of Binary Classifiers using Instance-Level Explanations. In *2017 IEEE Conference on Visual Analytics Science and Technology (VAST)*, pages 162–172, Phoenix, AZ, Oct. 2017. IEEE. ISBN 978-1-5386-3163-8. doi: 10.1109/VAST.2017.8585720. URL <https://ieeexplore.ieee.org/document/8585720/>.
- [33] A. Kumpf, B. Tost, M. Baumgart, M. Riemer, R. Westermann, and M. Rautenhaus. Visualizing Confidence in Cluster-Based Ensemble Weather Forecast Analyses. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):109–119, Jan. 2018. ISSN 1077-2626. doi: 10.1109/TVCG.2017.2745178. URL <http://ieeexplore.ieee.org/document/8019883/>.
- [34] B. C. Kwon, B. Eysenbach, J. Verma, K. Ng, C. De Filippi, W. F. Stewart, and A. Perer. Clustervision: Visual Supervision of Unsupervised Clustering. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):142–151, Jan. 2018. ISSN 1077-2626. doi: 10.1109/TVCG.2017.2745085. URL <http://ieeexplore.ieee.org/document/8019866/>.
- [35] B. C. Kwon, M.-J. Choi, J. T. Kim, E. Choi, Y. B. Kim, S. Kwon, J. Sun, and J. Choo. RetainVis: Visual Analytics with Interpretable and Interactive Recurrent Neural Networks on Electronic Medical Records. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):299–309, Jan. 2019. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2018.2865027. URL <https://ieeexplore.ieee.org/document/8440842/>.
- [36] E. J. Langer, A. Blank, and B. Chanowitz. The mindlessness of ostensibly thoughtful action: The role of "placebic" information in interpersonal interaction. *Journal of Personality and Social Psychology*, 36(6):635–642, 1978. ISSN 1939-1315(Electronic),0022-3514(Print). doi: 10.1037/0022-3514.36.6.635.
- [37] H. Lin, S. Gao, D. Gotz, F. Du, J. He, and N. Cao. RCLens: Interactive Rare Category Exploration and Identification. *IEEE Transactions on Visualization and Computer Graphics*, 24(7):2223–2237, July 2018. doi: 10.1109/TVCG.2017.2711030.
- [38] Z. C. Lipton. The Mythos of Model Interpretability. *Queue*, 16(3):30:31–30:57, June 2018. ISSN 1542-7730. doi: 10.1145/3236386.3241340. URL <http://doi.acm.org/10.1145/3236386.3241340>.
- [39] J. Liu, T. Dwyer, K. Marriott, J. Millar, and A. Haworth. Understanding the Relationship Between Interactive Optimisation and Visual Analytics in the Context of Prostate Brachytherapy. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):319–329, Jan. 2018. ISSN 1077-2626. doi: 10.1109/TVCG.2017.2744418. URL <http://ieeexplore.ieee.org/document/8017652/>.
- [40] M. Liu, J. Shi, Z. Li, C. Li, J. Zhu, and S. Liu. Towards Better Analysis of Deep Convolutional Neural Networks. *IEEE Transactions on Visualization and Computer Graphics*, 23(1):91–100, Jan. 2017. ISSN 1077-2626. doi: 10.1109/TVCG.2016.2598831. URL <http://ieeexplore.ieee.org/document/7536654/>.
- [41] M. Liu, S. Liu, H. Su, K. Cao, and J. Zhu. Analyzing the Noise Robustness of Deep Neural Networks. *arXiv:1810.03913 [cs, stat]*, Oct. 2018. URL <http://arxiv.org/abs/1810.03913>. arXiv: 1810.03913.
- [42] M. Liu, J. Shi, K. Cao, J. Zhu, and S. Liu. Analyzing the Training Processes of Deep Generative Models. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):77–87, Jan. 2018. ISSN 1077-2626. doi: 10.1109/TVCG.2017.2744938. URL <http://ieeexplore.ieee.org/document/8019879/>.

- [43] S. Liu, P. Bremer, J. J. Thiagarajan, V. Srikumar, B. Wang, Y. Livnat, and V. Pascucci. Visual Exploration of Semantic Relationships in Neural Word Embeddings. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):553–562, Jan. 2018. doi: 10.1109/TVCG.2017.2745141.
- [44] S. Liu, J. Xiao, J. Liu, X. Wang, J. Wu, and J. Zhu. Visual Diagnosis of Tree Boosting Methods. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):163–173, Jan. 2018. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2017.2744378. URL <https://ieeexplore.ieee.org/document/8017582/>.
- [45] Y. Ma, T. Xie, J. Li, and R. Maciejewski. Explaining Vulnerabilities to Adversarial Machine Learning through Visual Analytics. *IEEE Transactions on Visualization and Computer Graphics*, pages 1–1, 2019. doi: 10.1109/TVCG.2019.2934631.
- [46] M. Millecamp, N. N. Htun, C. Conati, and K. Verbert. To Explain or Not to Explain: The Effects of Personal Characteristics when Explaining Music Recommendations. In *Proceedings of the 24th International Conference on Intelligent User Interfaces, IUI '19*, pages 397–407, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6272-6. doi: 10.1145/3301275.3302313. URL <http://doi.acm.org/10.1145/3301275.3302313>. event-place: Marina del Ray, California.
- [47] T. Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence*, 267: 1–38, Feb. 2019. ISSN 0004-3702. doi: 10.1016/j.artint.2018.07.007. URL <http://www.sciencedirect.com/science/article/pii/S0004370218305988>.
- [48] T. Miller, P. Howe, and L. Sonenberg. Explainable AI: Beware of Inmates Running the Asylum Or: How I Learnt to Stop Worrying and Love the Social and Behavioural Sciences. In *arXiv:1712.00547 [cs]*, Dec. 2017. URL <http://arxiv.org/abs/1712.00547>. arXiv: 1712.00547.
- [49] Y. Ming, S. Cao, R. Zhang, Z. Li, Y. Chen, Y. Song, and H. Qu. Understanding Hidden Memories of Recurrent Neural Networks. In *2017 IEEE Conference on Visual Analytics Science and Technology (VAST)*, pages 13–24, Phoenix, AZ, Oct. 2017. IEEE. ISBN 978-1-5386-3163-8. doi: 10.1109/VAST.2017.8585721. URL <https://ieeexplore.ieee.org/document/8585721/>.
- [50] Y. Ming, H. Qu, and E. Bertini. RuleMatrix: Visualizing and Understanding Classifiers with Rules. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):342–352, Jan. 2019. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2018.2864812. URL <https://ieeexplore.ieee.org/document/8440085/>.
- [51] T. Muhlbacher, L. Linhardt, T. Moller, and H. Piringer. TreePOD: Sensitivity-Aware Selection of Pareto-Optimal Decision Trees. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):174–183, Jan. 2018. ISSN 1077-2626. doi: 10.1109/TVCG.2017.2745158. URL <http://ieeexplore.ieee.org/document/8019878/>.
- [52] B. M. Muir. Trust in automation: Part I. Theoretical issues in the study of trust and human intervention in automated systems. *Ergonomics*, 37(11):1905–1922, Nov. 1994. ISSN 0014-0139. doi: 10.1080/00140139408964957. URL <https://doi.org/10.1080/00140139408964957>.
- [53] C. Musto, F. Narducci, P. Lops, M. De Gemmis, and G. Semeraro. ExpLOD: A Framework for Explaining Recommendations Based on the Linked Open Data Cloud. In *Proceedings of the 10th ACM Conference on Recommender Systems, RecSys '16*, pages 151–154, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4035-9. doi: 10.1145/2959100.2959173. URL <http://doi.acm.org/10.1145/2959100.2959173>. event-place: Boston, Massachusetts, USA.
- [54] N. Pezzotti, B. P. F. Lelieveldt, L. v. d. Maaten, T. Höllt, E. Eisemann, and A. Vilanova. Approximated and User Steerable tSNE for Progressive Visual Analytics. *IEEE Transactions on Visualization and Computer Graphics*, 23(7):1739–1752, July 2017. doi: 10.1109/TVCG.2016.2570755.
- [55] F. Poursabzi-Sangdeh, D. G. Goldstein, J. M. Hofman, J. W. Vaughan, and H. Wallach. Manipulating and Measuring Model Interpretability. Feb. 2018. URL <http://arxiv.org/abs/1802.07810>. arXiv: 1802.07810.
- [56] E. Rader, K. Cotter, and J. Cho. Explanations As Mechanisms for Supporting Algorithmic Transparency. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems, CHI '18*, pages 103:1–103:13, New York, NY, USA, 2018. ACM. ISBN 978-1-4503-5620-6. doi: 10.1145/3173574.3173677. URL <http://doi.acm.org/10.1145/3173574.3173677>. event-place: Montreal QC, Canada.
- [57] D. Ren, S. Amershi, B. Lee, J. Suh, and J. D. Williams. Squares: Supporting Interactive Performance Analysis for Multiclass Classifiers. *IEEE Transactions on Visualization and Computer Graphics*, 23(1):61–70, Jan. 2017. ISSN 1077-2626. doi: 10.1109/TVCG.2016.2598828. URL <http://ieeexplore.ieee.org/document/7539404/>.
- [58] R. M. Richter, M. J. Valladares, and S. C. Sutherland. Effects of the Source of Advice and Decision Task on Decisions to Request Expert Advice. In *Proceedings of the 24th International Conference on Intelligent User Interfaces, IUI '19*, pages 469–475, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6272-6. doi:

- 10.1145/3301275.3302279. URL <http://doi.acm.org/10.1145/3301275.3302279>. event-place: Marina del Ray, California.
- [59] B. G. Robbins. What is Trust? A Multidisciplinary Review, Critique, and Synthesis. *Sociology Compass*, 10(10):972–986, 2016. ISSN 1751-9020. doi: 10.1111/soc4.12391. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/soc4.12391>.
- [60] D. Sacha, M. Kraus, J. Bernard, M. Behrisch, T. Schreck, Y. Asano, and D. A. Keim. SOMFlow: Guided Exploratory Cluster Analysis with Self-Organizing Maps and Analytic Provenance. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):120–130, Jan. 2018. ISSN 1077-2626. doi: 10.1109/TVCG.2017.2744805. URL <http://ieeexplore.ieee.org/document/8019867/>.
- [61] K. E. Schaefer, J. Y. C. Chen, J. L. Szalma, and P. A. Hancock. A Meta-Analysis of Factors Influencing the Development of Trust in Automation: Implications for Understanding Autonomy in Future Systems. *Human Factors*, 58(3):377–400, May 2016. ISSN 0018-7208. doi: 10.1177/0018720816634228. URL <https://doi.org/10.1177/0018720816634228>.
- [62] J. Schaffer, J. O’Donovan, J. Michaelis, A. Raglin, and T. Höllerer. I Can Do Better Than Your AI: Expertise and Explanations. In *Proceedings of the 24th International Conference on Intelligent User Interfaces, IUI ’19*, pages 240–251, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6272-6. doi: 10.1145/3301275.3302308. URL <http://doi.acm.org/10.1145/3301275.3302308>. event-place: Marina del Ray, California.
- [63] R. Sevastjanova, H. Schäfer, J. Bernard, D. Keim, and M. El-Assady. Shall we play? extending the visual analytics design space through gameful design concepts. pages 1–9, 2019.
- [64] T. Spinner, U. Schlegel, H. Schäfer, and M. El-Assady. explAIner: A Visual Analytics Framework for Interactive and Explainable Machine Learning. *IEEE Transactions on Visualization and Computer Graphics*, pages 1–1, 2019. doi: 10.1109/TVCG.2019.2934629.
- [65] A. Springer and S. Whittaker. Progressive Disclosure: Empirically Motivated Approaches to Designing Effective Transparency. In *Proceedings of the 24th International Conference on Intelligent User Interfaces, IUI ’19*, pages 107–120, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-6272-6. doi: 10.1145/3301275.3302322. URL <http://doi.acm.org/10.1145/3301275.3302322>. event-place: Marina del Ray, California.
- [66] J. Stahnke, M. Dörk, B. Müller, and A. Thom. Probing Projections: Interaction Techniques for Interpreting Arrangements and Errors of Dimensionality Reductions. *IEEE Transactions on Visualization and Computer Graphics*, 22(1):629–638, Jan. 2016. doi: 10.1109/TVCG.2015.2467717.
- [67] F. Stoffel, L. Flekova, D. Oelke, I. Gurevych, and D. A. Keim. Feature-based visual exploration of text classification. In *Symposium on Visualization in Data Science at IEEE VIS*, 2015.
- [68] H. Strobel, S. Gehrmann, H. Pfister, and A. M. Rush. LSTMVis: A Tool for Visual Analysis of Hidden State Dynamics in Recurrent Neural Networks. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):667–676, Jan. 2018. ISSN 1077-2626. doi: 10.1109/TVCG.2017.2744158. URL <http://ieeexplore.ieee.org/document/8017583/>.
- [69] H. Strobel, S. Gehrmann, M. Behrisch, A. Perer, H. Pfister, and A. M. Rush. Seq2seq-Vis: A Visual Debugging Tool for Sequence-to-Sequence Models. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):353–363, Jan. 2019. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2018.2865044. URL <https://ieeexplore.ieee.org/document/8494828/>.
- [70] D. Wang, Q. Yang, A. Abdul, and B. Y. Lim. Designing Theory-Driven User-Centric Explainable AI. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems, CHI ’19*, pages 601:1–601:15, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-5970-2. doi: 10.1145/3290605.3300831. URL <http://doi.acm.org/10.1145/3290605.3300831>. event-place: Glasgow, Scotland Uk.
- [71] J. Wang, L. Gou, H.-W. Shen, and H. Yang. DQNViz: A Visual Analytics Approach to Understand Deep Q-Networks. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):288–298, Jan. 2019. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2018.2864504. URL <https://ieeexplore.ieee.org/document/8454905/>.
- [72] M. Yin, J. Wortman Vaughan, and H. Wallach. Understanding the Effect of Accuracy on Trust in Machine Learning Models. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems, CHI ’19*, pages 279:1–279:12, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-5970-2. doi: 10.1145/3290605.3300509. URL <http://doi.acm.org/10.1145/3290605.3300509>. event-place: Glasgow, Scotland Uk.
- [73] K. Yu, S. Berkovsky, R. Taib, D. Conway, J. Zhou, and F. Chen. User Trust Dynamics: An Investigation Driven by Differences in System Performance. In *Proceedings of the 22Nd International Conference on Intelligent User Interfaces, IUI ’17*, pages 307–317, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-4348-0. doi:

10.1145/3025171.3025219. URL <http://doi.acm.org/10.1145/3025171.3025219>. event-place: Limassol, Cyprus.

- [74] J. Zhang, Y. Wang, P. Molino, L. Li, and D. S. Ebert. Manifold: A Model-Agnostic Framework for Interpretation and Diagnosis of Machine Learning Models. *IEEE Transactions on Visualization and Computer Graphics*, 25(1): 364–373, Jan. 2019. doi: 10.1109/TVCG.2018.2864499.
- [75] X. Zhao, Y. Wu, D. L. Lee, and W. Cui. iForest: Interpreting Random Forests via Visual Analytics. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):407–416, Jan. 2019. doi: 10.1109/TVCG.2018.2864475.
- [76] J. Zhou, Z. Li, H. Hu, K. Yu, F. Chen, Z. Li, and Y. Wang. Effects of Influence on User Trust in Predictive Decision Making. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems*, CHI EA '19, pages LBW2812:1–LBW2812:6, New York, NY, USA, 2019. ACM. ISBN 978-1-4503-5971-9. doi: 10.1145/3290607.3312962. URL <http://doi.acm.org/10.1145/3290607.3312962>. event-place: Glasgow, Scotland Uk.