

# Predictability and Comprehensibility in Post-Hoc XAI Methods: A User-Centered Analysis

ANAHID JALALI, AIT - Austrian Institute of Technology, Austria

BERNHARD HASLHOFER, AIT - Austrian Institute of Technology, Austria

SIMONE KRIGLSTEIN, AIT - Austrian Institute of Technology, Austria and Masaryk University, Czech Republic

ANDREAS RAUBER, Vienna University of Technology, Austria

Post-hoc explainability methods aim to clarify predictions of black-box machine learning models. However, it is still largely unclear how well users comprehend the provided explanations and whether these increase the users' ability to predict the model behavior. We approach this question by conducting a user study to evaluate comprehensibility and predictability in two widely used tools: LIME and SHAP. Moreover, we investigate the effect of counterfactual explanations and misclassifications on users' ability to understand and predict the model behavior. We find that the comprehensibility of SHAP is significantly reduced when explanations are provided for samples near a model's decision boundary. Furthermore, we find that counterfactual explanations and misclassifications can significantly increase the users' understanding of how a machine learning model is making decisions. Based on our findings, we also derive design recommendations for future post-hoc explainability methods with increased comprehensibility and predictability.

CCS Concepts: • **Human-centered computing** → **User studies**.

Additional Key Words and Phrases: Interpretable Machine Learning, XAI, XAI Evaluation, User Study

## ACM Reference Format:

Anahid Jalali, Bernhard Haslhofer, Simone Kriglstein, and Andreas Rauber. 2022. Predictability and Comprehensibility in Post-Hoc XAI Methods: A User-Centered Analysis. 1, 1 (September 2022), 17 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

## 1 INTRODUCTION

The opacity of machine learning models is a well-known problem in application areas such as health care systems, financial services, or industrial applications [2], where transparency and accountability are fundamental requirements. When models are not transparent, users and machine learning experts (ML) have difficulty explaining how models arrive at their predictions. Therefore, ongoing research in the field of Interpretable Machine Learning (IML), also known as eXplainable Artificial Intelligence (XAI), focuses on implementing methods that either examine the inner structure of models (model-specific) or explain predictions of a trained model based on a training dataset (post-hoc) [21].

Well-known post-hoc explanation techniques are: Local Interpretable Model Agnostic Explanations (LIME) [26], DeepLift [31], COVAR [9], and Shapely Additive exPlanations (SHAP) [15]. They provide different types of explanations

---

Authors' addresses: Anahid Jalali, [anahid.jalali@ait.ac.at](mailto:anahid.jalali@ait.ac.at), AIT - Austrian Institute of Technology, Vienna, Austria; Bernhard Haslhofer, AIT - Austrian Institute of Technology, Vienna, Austria, [bernhard.haslhofer@ait.ac.at](mailto:bernhard.haslhofer@ait.ac.at); Simone Kriglstein, AIT - Austrian Institute of Technology, Vienna, Austria and Masaryk University, Brno, Czech Republic, [simone.kriglstein@ait.ac.at](mailto:simone.kriglstein@ait.ac.at); Andreas Rauber, Vienna University of Technology, Vienna, Austria, [andreas.rauber@ifs.tuwien.ac.at](mailto:andreas.rauber@ifs.tuwien.ac.at).

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2018 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

1

to the user and aim at fulfilling the goals of interpretability, which are: following a model’s prediction for a given dataset, being easy to comprehend, and being efficient [28]. However, as [2] stated, there is ambiguity in the definition of interpretability and how it should be measured and evaluated. The authors further argued that application-grounded evaluation is the most appropriate as “it assesses interpretability in the end goal with the end-users”.

Several user studies already evaluated model interpretations and explanations: [26, 27] measured comprehensibility and trust by asking users to explain the best model, the most suitable features, as well as model behaviors and irregularities. Other studies also focused on measuring comprehensiveness [23, 24], usefulness [19, 26], and trustworthiness of explanations [13, 23, 24, 29]. However, they all evaluated a single XAI method and aimed to improve its accuracy for a particular task or manipulated explanations to mislead users and measure their bias when presenting a fidelity or model accuracy score for a given task. Furthermore, they studied mainly annotation tasks on image or text datasets.

Our work is motivated by [29], who measured comprehensibility and trust based on the users’ interaction time in a text classification task. However, it remains unclear how users make their judgments: do they blindly follow provided model explanations or make their judgments based on the real-world meaning of data points, which were words in this case? Furthermore, current studies do not infer recommendations that could inform future, improved model explainability methods.

The need for XAI user studies has been pointed out by [12], who argue that XAI techniques must align with the mental model of ML-practitioners. Also, [18] stated that “the ultimate goal is for people (experts and/or users) to understand the models, and it is, therefore, essential to involve human feedback and reasoning as a requisite component for design and evaluation of interpretable-ML systems.”

We acknowledge the work of Jacovi et al. [11] on formalizing human-AI trust, in which they stated “*defining trust as the user’s attempt to predict the impact of the model behavior under risk and uncertainty [10] is a goal but not necessarily a symptom of trust*”. Moreover, Mohseni et al. [20], defines one of the desired properties of explainer systems as “predictability”, which is the ability of these systems to “support building a mental model of the system that enables user to predict system behavior.” Therefore, in this work, we aim at measuring the *comprehensibility* and the *predictability*, to compare two well-known and widely used XAI approaches, SHAP and LIME. For this purpose, we first refine our notion of *comprehensibility* and *predictability* in classification tasks as follows:

**Comprehensibility:** denotes the user’s ability to transfer information on feature contributions obtained from model explanations across samples of the same class.

Previous studies [5, 14, 26, 29] assume trust as “trust in model correctness” and evaluate the user’s ability to guess a sample’s label correctly, given model explanations. We additionally consider both notions of predictability [20] and simulation [8], which is the user’s ability to guess a model’s prediction on a new sample correctly, and broaden the definition of predictability as follows.

**Predictability:** denotes the explainer’s ability to support the users with predicting model predictions, which can be correct or incorrect, on a new sample given model explanations for a sample that the model predicted correctly with high confidence.

In the following, we aim to evaluate SHAP and LIME explanations qualitatively and quantitatively as part of a user study. We formulate our research questions as follows:

**RQ1.** To what extent do users comprehend the explanations provided by different XAI methods, and are they able to predict the decision made by the model?

When interacting with these tools, we noticed that it is essential to understand how features impact a model decision and that this is easier to comprehend when a model is more confident about a decision. We tested these hypotheses and answered the above question in a comparative study, which we describe in more detail in Section *RQ1: Comprehensibility*. We found evidence that supports our hypotheses for SHAP.

Furthermore, we observed that users who understand model predictions for the two different classes found it easier to classify unexplained and unlabeled samples. Therefore, we hypothesize that users predictability increases when they can classify new samples using the explanations from samples of different classes. In Section *RQ1: Predictability*, we elaborate on this in more detail and show that users are able to predict the decisions made by the model, using both SHAP and LIME equally.

Given these findings, we further investigate the effect of counterfactual and misclassified samples on the users ability to predict the model's decision. More precisely, we consider the following research question:

**RQ2.** To what extent can visualizations of counterfactual and misclassified samples improve the user's predictability?

We answered this question by testing the following hypothesis: adding explanations of misclassified and counterfactual samples can improve the predictability and support anticipating the model's behavior. We describe our experiment in more details in Section *RQ2: Improving Predictability with Visualizations* and report evidence that supports that hypothesis for both SHAP and LIME. We found that users have higher predictability with LIME than with SHAP explanations.

Moreover, throughout our experiments, we asked users to provide their subjective feedback on given explanations. Their responses allowed us to answer our third and final research question:

**RQ3.** To what extent can visualizations of local XAI explanations guide users in finding global explanations?

Throughout our experiments, we collected user feedback for both LIME and SHAP, and asked the participants to explain possible shortcomings and sources of confusion. As detailed in Section *RQ3: Qualitative Analysis*, we observed that negative feedback provides valuable insight into how users interpret explanations. For example, the users complained about inconsistent explanations of values and different scaling of visualizations. Our experiments showed that this reduced the ability of the user to understand the behavior of a model.

The findings of our comparative user study can substantially contribute to the design of new or refinement of existing XAI methods, and therefore, we propose a set of design recommendations in Section *Discussion*, which we will confer in more detail.

## 2 RELATED WORK

Previous studies, such as [2], [37], [35] or [32], discuss the various properties of explainability and define evaluation criteria for qualitative (e.g. measuring explanations' comprehensiveness, trustworthiness) and quantitative (e.g. measuring explanations' accuracy, fidelity and consistency) approaches.

The need for an explainability baseline for evaluating the quality of explanations was raised by [19], who measured the trustworthiness of explanations and quantitatively compared LIME with Grad-Cam explanations on three different baselines: human-attention mask, segmentation mask, and human judgment ratings. Their findings indicate human biases in ratings and significant differences in evaluation scores, which they assume to be caused by "clear non-uniform distribution of weights in human attention masks".

Our work builds on previous user studies, which evaluated model interpretations and explanations. We summarize them into Table 1 and roughly divide them into three categories: the first category ([1, 24, 29]) evaluates explainability quantitatively and measures the success of the user with and without explanations. The second category ([8, 13, 23, 24, 26, 27]) measures comprehensiveness and trustworthiness of explanations quantitatively by aligning meaningful and meaningless or manipulated explanations with human logic. They point out that manipulated explanations can increase user trust in biased models and lead to mistrust when explanations do not match the decision-making process of users. The third category ([12, 30, 33, 36]) measures the comprehensibility and trust of users by considering system interactions. These approaches focus on understanding how users perceive information, diagnose and refine the AI systems. Their qualitative results are often used as recommendations for designing the XAI approaches to improve users’ trust.

Thus, previous studies focused mainly on image or text data, misleading the users with manipulated explanations in binary annotation tasks. User considerations of the decision-making process and design recommendations on improving explanations were mainly out of scope. For more detailed information on each of the mentioned works, please refer to the summary in table 1.

Therefore, we do not limit ourselves to a quantitative evaluation but also evaluate XAI approaches qualitatively and consider user feedback to understand the reasons behind possible confusion in decision making. Moreover, we include explanations of counterfactual and misclassified samples to test users’ predictability using the explanations.

Table 1. A summary of existing studies for evaluating XAI approaches.

Paper	Evaluation			Explanations			
	Metric	Approach	Data	Manipulated		Show	
				Model	Examples	Preds	Labels
[1]	Compreh.	Quan.	Image	×	×	✓	×
[8]	Simulation (Trust)	Quan.	Text	×	✓	✓	×
[12]	Trust	Quan. Qual.	Tabular	×	×	✓	×
[13]	Trust	Quan.	Tabular	✓	✓	×	×
[19]	Trust	Quan.	Image Text	×	×	✓	×
[23]	Trust	Quan.	Image	×	✓	×	×
[24]	Trust Compreh.	Quan.	Text	×	✓	×	×
[26]	Trust Compreh.	Quan.	Image Text	✓	×	✓	×
[27]	Compreh.	Quan.	Image	×	×	✓	×
[29]	Trust Compreh.	Quan.	Text	✓	×	✓	×
[30]	Trust	Quan.	Text	×	×	✓	×
[33]	Compreh.	Qual.	Image	×	×	✓	×
[36]	Trust	Qual.	Tabular (Medical)	×	×	✓	×
Our Work	Predictability Compreh.	Quan. Qual.	Tabular	×	×	✓	✓

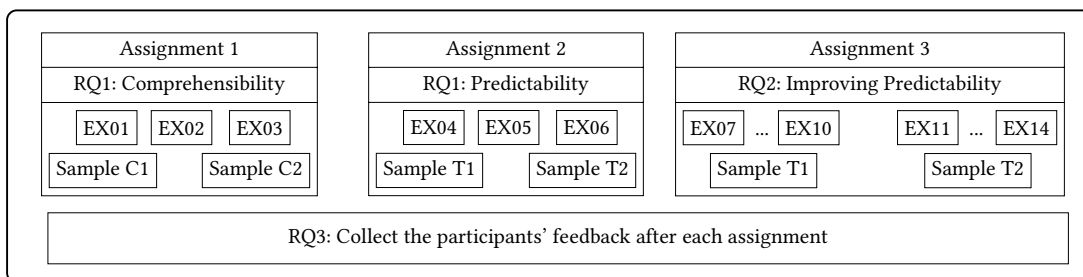


Fig. 1. An overview of the assignments of our user study. We asked participants to work on assignments and specific tasks, answering a specific research question in the main survey. Assignment one and two contained three samples and their explanations, depicted as *EXxx*, provided by either LIME or SHAP. The users studied the samples and explanations and had to answer questions for a test sample, depicted as *Sample Cx* or *Tx*

### 3 METHODOLOGY

We conducted a user study to evaluate comprehensibility and predictability in explanations provided by two widely-used XAI methods: SHAP and LIME. Our experimental setup follows a between-subject design, with the XAI method as the primary varying condition. That means we exposed each participant to only one condition (LIME or SHAP), which shortened the duration of experimental sessions for each participant (c.f., [16]).

Figure 1 depicts our overall experimental setup, which started with participant recruitment and a pre-test survey phase. As part of the central survey, we defined, for each research question, several assignments and tasks to be solved by the participants. The first assignment measured the user’s comprehensibility for a given XAI method, and the second one the explainer’s predictability. The third assignment investigated the effect of adding explanations of misclassified and counterfactual samples on explainer’s predictability. In each task, we presented visualizations of SHAP- or LIME-explanations to the participant and asked them to answer four questions, in which they had to interpret the visualizations for the given test sample.

*Dataset and Implementation.* We chose the Boston Housing dataset [6] because of its simplicity and transparency. This dataset estimates the median price of apartments in Boston. We transferred this regression task into a classification task by categorizing the estimated prices into three classes: 1) low-price, 2) medium-price, and 3) high-price while preserving the feature correlations with the target variable. We used five features that give our model the highest accuracy: average number of rooms, pupil and teacher ratio, air pollution level, crime rate, and the zone where an apartment is located.

In an initial trial experiment, we found that participants tend to project their interpretation of feature labels (e.g., crime rate) onto an explanation instead of interpreting the information provided by either LIME or SHAP. Therefore, we anonymized the feature names to *F1*, *F2*, *F3*, *F4*, and *F5*, respectively.

We further min-max normalized each feature and trained a machine learning model using a 3-layered fully connected dense neural network (64 units, Relu function, and a softmax at the output layer). We optimized the model’s trained weights with Stochastic Gradient Descent (0.001 learning rate). The model had 93% accuracy, and median prediction probability of 70.21%, 46.81%, and 72.87% for low-price, medium-price, and high-price classes, respectively. However, we did not provide the participants with this information.

To set up our experiment, we used Python 3.7, and for the explainability visualizations from LIME and SHAP, we used the *LimeTabularExplainer* (*lime* library version 0.2.0.1), and the *KernelExplainer* (*shap* library version 0.34.0).

*Participants.* We recruited participants with ML and Data Science experience having technical backgrounds in Computer Science, Mathematics, and Physics. The participants were scientists and practitioners from eight different institutions in five countries, collected via their LinkedIn profile. Overall, 47 participants took part in our experiment, and we randomly assigned one XAI approach, either SHAP or LIME, to each participant. We compensated participants with €20 for their approximately 1-hour effort of taking part in the experiment.

We had to remove one user, who answered “I do not know” and stated afterwards that he was not focused and could not participate in this study. This data cleaning step left us with a total of 46 participants: 30 male and 16 female with an average age of 31, 24 users evaluated LIME, and 22 users evaluated SHAP. When asked about their experience with explainability approaches and interpreting machine learning; they often stated that they interpret models by looking at feature importance plots of decision trees or on the coefficients of linear regression models. On the other hand, only twelve users had experience with XAI approaches such as LIME, Layer-wise relevance-propagation (LRP), heatmaps, or Google’s Language Interpretability Tool (LIT) [34]. Therefore, we can assume that most of our participants (34 of 46) were non-experts in XAI.

*Survey Procedure.* We implemented our survey using the [25] platform. After the users gave their consent to the overall experimental design, they started the pre-test survey, in which we asked them about their demographic information, background, and data science experience. Then, we measured their experience by presenting them with 12 data science and machine learning know-how questions such as bias and variance trade-off or distribution functions.

Afterwards, we acquainted the participants with the overall survey procedure in an initial training phase, in which we explained the dataset, the tasks, the visualizations provided by each explainability approach, and the structure of the assignments. Then, in the second part of the survey, the participants started working on the three assignments, each comprising two tasks with four questions. Thus, each participant had to answer 24 questions related to interpretability in total.

Finally, we present the participants with the Nasa Task Load Index (TLX) [7] questionnaire to obtain insight into the mental, physical and temporal demand of the survey as well as the participants’ success, effort, and frustration.

*Data Coding.* For the quantitative part of our analysis, we assigned scores to each multiple-choice answer and computed the sum of all answers to compare assignment results. We gave each correct answer a score of 2, each wrong answer a -1, and “I do not know” answers a 0. This scoring scheme allows us to distinguish participants who tried to answer the questions seriously but possibly wrong from those who just checked “I don’t know.”

For the qualitative evaluation, we interviewed the participants and transcribed their feedback throughout each interview session. We followed the Mayring qualitative analysis decoding rules described in [17] to categorize the transcribed participant feedback.

#### 4 RQ1: COMPREHENSIBILITY

This section provides answers to our first research question, RQ1, which seeks to understand the relationship between the comprehensibility of a set of three explanations for the user and the prediction confidence of a machine learning model.

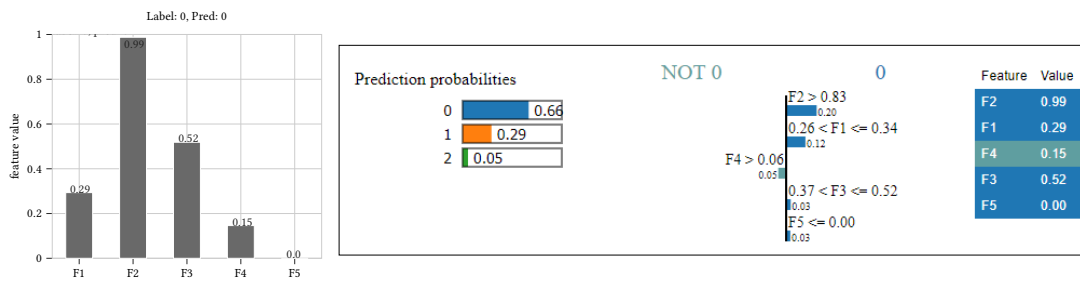


Fig. 2. Example explanation (EX01) provided to the user. On the left hand side, it shows a sample (an apartment) that the model classified as being low-price (label 0). The x-axis of the bar plot are the sample’s attributes and the y-axis are the values. On the right hand side, a LIME explanation, describes the decision of the model.

For this purpose, we randomly assigned each user to a XAI method (LIME or SHAP) and presented them with three explained samples (EX01, EX02, and EX03), each representing an apartment that the machine learning model classified as belonging to the class *low-price* (coded as 0). Figure 2 illustrates one of these three samples and shows how we presented and explained it to the users. We have chosen these samples based on their information and ensured that several features contributed to its low-price classification.

Moreover, we chose two additional test samples from the Boston Housing dataset, C1 and C2, with features similar to the explained samples, which the model correctly classified as low-price. However, the prediction confidence for C1 was higher than the confidence for C2, indicating that C1 is further away from the model’s decision boundary. We use the model’s Probability Distribution Delta (PDD) to quantify the model’s confidence.

The users studied the explanations and tried to use the information they learned from EX01, EX02, and EX03 to answer the following multiple-choice questions for C1 and then C2:

- (1) Choose two features that highly influence the prediction of class “low-price (0)”.
- (2) How does the value of F2 (together with F1 and F5) influence the model’s decision on class “low-price (0)”?
- (3) How does the value of F1 affect the probability of class “low-price (0)”?
- (4) How does the value of F3 in this sample (w.r.t. explanations) affect the probability of class “low-price (0)”?

With the above questions, we wanted to measure the user’s understanding of whether they can interpret how feature values can increase or decrease the probability of being classified as “low price”. Therefore, we also formulated control questions to avoid random answers and ensure participants focused on the tasks. For the third question, for instance, the answer should match the answer to the second question. If this is not the case, we know that an answer is random and that we should remove it from our analysis. However, this was not the case with our users, and they did not randomly answer the questions.

We compute scores for all the responses, giving an overall comprehensibility score for each user and each assignment. We also compared the individual question scores of C1 and C2 to measure quantitatively whether the presented visualization helps the user comprehend each feature’s contribution to the model’s decision.

Furthermore, we code the participant’s interview responses into three categories: (i) C1 was more difficult than C2, (ii) C2 was more difficult than C1, and (iii) C1 and C2 were equally challenging.

*Results.* We collected responses from 46 participants, each answering the questions above for either LIME (24) or SHAP (22), and computed the overall comprehensibility score for each user. Figure 3 shows the minimum, maximum,

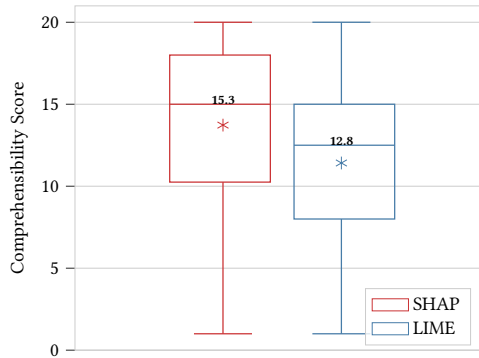


Fig. 3. LIME and SHAP’s comprehensibility score. We included the mean values with stars, and included the median values above the median lines of each box-plot. No significant difference is visible between these two XAI methods.

sample median, as well as the first and third quartiles of the comprehensibility scores for both methods. Post-hoc comparison of mean values using a two-sample t-test revealed that there is no significant difference ( $t=-1.56$ ,  $p=0.12$ ) between the mean comprehensibility score of SHAP (13.73) and LIME (11.42). This result shows that LIME and SHAP are equally comprehensible by the participants.

Next, we tested whether a model’s confidence in the prediction, which we can measure by considering the PDD between possible classes, affects the users’ comprehensibility. Recall that C1 is further away from the model’s decision boundary than C2. Since the responses to these tasks represent variables from repeated measures groups, we follow [3] and compare means by first calculating an adjustment factor for each user. We computed that factor by subtracting the participant’s means (pMean) from the mean of both C1 and C2. We then add these adjustment factors to our participants’ actual comprehensibility scores, comparing C1 and C2. Again, we did not see a significant difference ( $t=-0.80$ ,  $p=0.42$ ) between the mean comprehensibility scores of C1 (5.42) and C2 (6.0). However, as shown in Figure 4, for SHAP, we measured a significant decrease ( $t= 5.54$ ,  $p=0.00$ ) between the scores of C1 (8.18) and C2 (5.55). These results show that SHAP’s explainability visualizations were less comprehensible to the participants when the model’s prediction was closer to the decision boundary. Moreover, we see a significant difference ( $t=4.51$ ,  $p=0.00$ ) between the mean comprehensibility score of SHAP (8.18) and LIME (5.42) for C1. That difference indicates that SHAP’s users’ comprehensibility was higher than LIME’s for C1. The low score can be explained by the “I do not know” answers on the F3 feature contribution.

In our first intermediate analysis, we noticed the difference between SHAP and LIME’s comprehensibility scores, and therefore, to obtain further qualitative insight into the effect of the decision boundary distance on the participants’ comprehensibility, we started asking them the following question right after the first assignment was over:

Q1. Which task did you find more difficult (between C1 and C2)? And why?

We collected the feedback from 20 participants and categorized their answers into three groups: i) “C1 more difficult than C2”, ii) “C2 more difficult than C1”, and iii) “C1 and C2 similarly difficult”.

Fifteen participants stated that they found answering the questions for C1, which has a higher distance from the decision boundary than C2, more complicated than C2. One participant, for instance, stated “*I did not know how to answer the questions and work with visualizations for C1. But for C2, it was clear for me how to use the visualizations and find my answers*”. This result indicates a learning effect involved and that participants became familiar with the



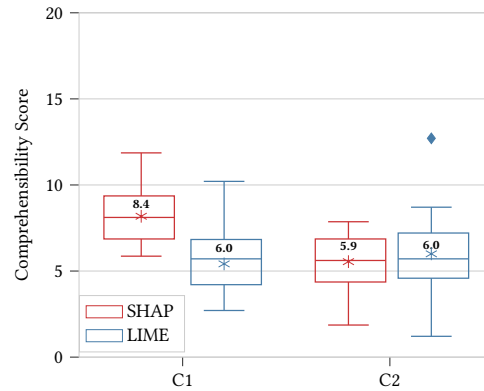


Fig. 4. C1 and C2 scores for both LIME and SHAP. For LIME, we observe no significant difference between C1 and C2 scores. On the other hand, we observe a significant decrease from C1 to C2 for SHAP’s comprehensibility tasks.

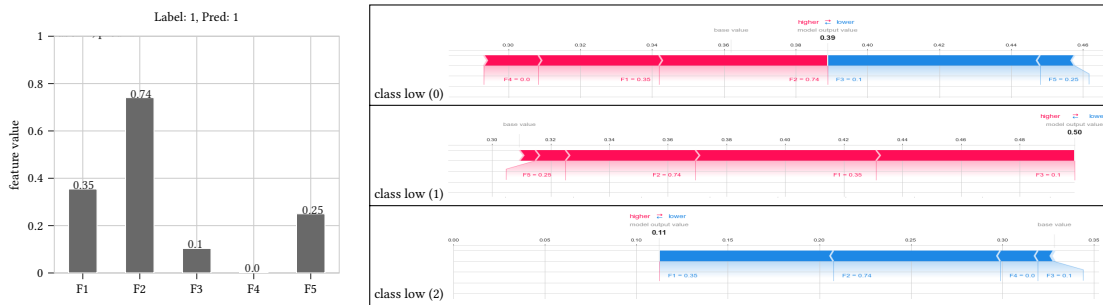


Fig. 5. Example SHAP explanation (EX05) shown to the user as part of Assignment 2. The sample is classified as medium-price (label 1) and depicted as bar plot with its actual label and the model’s prediction next its plot.

visualizations over time while answering these samples’ questions. We compared their feedback with their scores and noticed that none of these 15 participants answered the questions for C2 correctly, which indicates a mismatch between perceived and actual comprehensibility.

### 5 RQ1: PREDICTABILITY

After we reported on the *comprehensibility* aspect of the first research question, RQ1, in the previous section, we now focus on the *predictability* aspect.

For that purpose, we presented the second assignment to the participants (see Figure 1) to analyze whether they can predict the model’s behavior and detect the misclassification using the information they receive from the explanations. Figure 5 shows one of the three SHAP explanations we presented to the user. In contrast to the previous assignment, each was classified differently as low-, medium-, or high-price. In two tasks, we also introduced misclassified samples from medium- to low-price (T1) and from high- to medium-price (T2).

We asked the participants to answer the following multiple-choice questions for the test samples in each task:

- (1) What would the model predict when considering only the value of F1.
- (2) What would the model predict when we consider the value of F2 and F1?

Table 2. Scoring rules for the last question in Assignment 2. The rules are based on what the user detected.

score	detected correct label	detected prediction correctly
4	yes	yes
1	no	yes
1	yes	no
-2	no	no
1	User suspects the misclassification when the model misclassifies a sample	
0	User answers with "I do not know"	

- (3) What is the effect on class probabilities if we consider F1, F2, and F3 (or F5) values?  
(4) What does the model predict for this sample?

The fourth question is about the model’s behavior and follows a different scoring mechanism than the other questions. Table 2 lists the possible answers and the scoring rules we applied. Following our comprehensibility scoring, we give each correct answer (yes) a score of 2, and each wrong answer (no) a score of -1. For example, a participant that guesses the label correctly but predicts the model prediction wrong, receives a score of 1 (2 + -1). However, participants who suspect the misclassification receive a score of 1 for correctly guessing the model’s failure.

*Results.* Analogous to the comprehensibility score reported in the previous section, we computed a predictability score for both XAI methods. As shown in Figure 6 (RQ1), and confirmed by a two-sample t-test (t-value=-0.553, p=0.582) we found no significant difference between the mean predictability scores of SHAP (5.64) and LIME (6.38), which were both drawn from Gaussian distributions that were not significantly different from each other.

We further apply the Mann-Whitney test to analyze the difference between LIME and SHAP answer category distribution. We categorize the answers into three groups; first, the category of correct guesses on model prediction, second are the incorrect prediction guesses, and third is the neutral category "I do not know". For SHAP with 22 participants and their answers for two tasks, this categorization results in 15 counts for the first category, 20 counts for the second category and 9 counts for neutral category. Categorization of LIME with 24 participants results in 20 counts for the first category, 23 counts for the second category and 5 counts for neutral category. We find no significant difference between LIME and SHAP answer category distribution for this assignment (t-value=922.5, p=0.128).

We also analyzed the participant’s answers to see whether they could predict the model’s behavior using this set of explained samples. We compared the first question of samples T1 and T2, which only considers the values of feature F1. We noticed that participants using LIME answered correctly more often than those using SHAP and that SHAP users struggled to find a threshold to decide whether the feature value increased the classification probability. LIME explanations, on the other hand, provide such a threshold range.

For the second question, we noticed that participants using SHAP scored better than those using LIME. Based on LIME participants’ feedback, the effect of F1 and F2 values on medium-price class was unclear to them.

As part of the third question, we asked the participants about the impact of feature value F3. We wanted to know whether it pushes a decision towards a class and away from another class. The average scores of both participants groups, SHAP and LIME, were low for this question, with a mean score of 0.27 and 0.66, respectively. Fourteen participants

using SHAP answered wrong, and 4 participants answered with “I do not know”. On the other hand, 8 Participants using LIME answered wrong, and 6 participants answered with “I do not know”.

Finally, the participants performed equally well when answering the fourth question for SHAP and LIME. Breaking down the scores and looking deeper, we noticed that SHAP participants often predicted the label of a sample correctly (15 from 22 participants) but failed to predict the model’s behavior. On the other hand, LIME participants detected model misclassifications more often and predicted labels correctly (15 out of 24 participants).

## 6 RQ2: IMPROVING PREDICTABILITY WITH VISUALIZATIONS

We now address our second research question, RQ2, and examine how visualizations of misclassified and counterfactual samples can improve the users’ predictability.

For that purpose, Assignment 3 uses the same set of questions and test samples as Assignment 2 but presents explanations about the misclassified and correctly classified samples to the user. This experiment allows us to compare predictability scores across assignments and measure how they are affected by additional visualizations. We illustrate the tasks and the assignment in Figure 1. All chosen samples have low PDD values and are close to the model’s decision boundary.

*Results.* After checking that sample scores for both methods follow a Gaussian distribution, we compared the mean scores of both methods using a two-sample t-test and found that LIME’s predictability score was significantly higher (t-value=-2.263, p=0.029) than that of SHAP. This result is also visible in Figure 6, which also shows that LIME’s visualization of misclassified samples has a more substantial effect on improving the explainer’s predictability.

We further compared the total sum of the achieved predictability scores of Assignment 3 with those of the previously conducted Assignment 2 and see, as shown in Figure 6, significant improvements for LIME (t-value=-4.387, p=0.0), but not for SHAP (t-value=-1.97, p=0.055). Since the responses in the Assignments 2 and 3 also represent variables from repeated measures groups, analogous to the comprehensibility score in RQ1, we follow Field et al., [3] method of comparing means for repeated measures by first calculating an adjustment factor for each user and adding this factor to their predictability scores. As shown in Figure 7, we see a significant improvement of the participants for both LIME (t-value=-4.301, p=0.0) and SHAP (t-value=-2.245, p=0.030). This result indicates that presenting explanations around the model’s decision boundary to the user helps the user to understand the model’s decision-making.

Analogous to our analysis in predictability RQ1, we apply the Mann-Whitney test to analyze the difference between LIME and SHAP category of answers distribution and find no significant difference (t-value=961.0, p=0.203). The categorization of SHAP answers results in 23 counts for the first category, 17 counts for the second category and 4 counts for neutral category. Categorization of LIME results in 31 counts for the first category, 15 counts for the second category and 2 counts for neutral category. We further compare the category distribution between LIME answers in RQ1 and RQ2 and notice the distributions are significantly different (t-value=702.0, p=0.0). We find the same when comparing SHAP’s answer category distribution between RQ1 and RQ2 (t-value=685.5, 0.004).

Moreover, we observed that participants correctly identified sample labels and predicted the correct class for the misclassified test sample. From only 4 participants using LIME and 7 participants using SHAP, we see an improvement to 15 participants using LIME and 12 participants using SHAP, who answered the third assignment correctly.

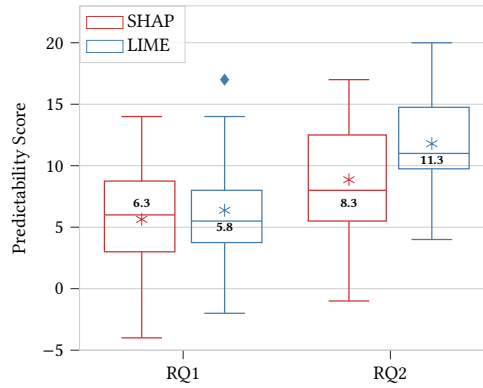


Fig. 6. LIME's and SHAP's predictability score before and after the users study the explained counter-factual and misclassified samples. LIME shows significant improvement.

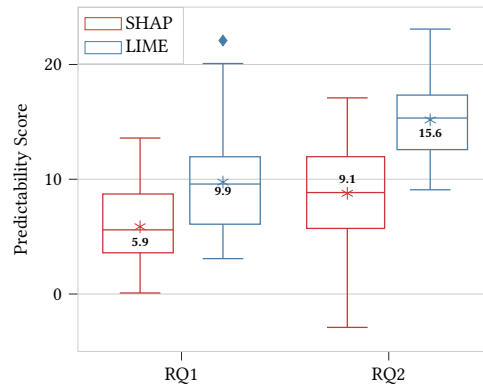


Fig. 7. LIME's and SHAP's calculated adjusted mean of predictability score before and after each user studies the explained counter-factual and misclassified samples. Both methods show significant improvements.

## 7 RQ3: QUALITATIVE ANALYSIS

In this section, we answer our third research question, RQ3, studying the participants' feedback to provide a guideline for improving the design of XAI methods with local explanations.

We seek to answer this question by qualitatively analyzing the participants' input. Therefore, at the end of the third assignment, we presented the participants with the following question:

Q2. How confident were you when answering the questions? When did you feel less confident?

We followed Mayring's qualitative coding rules[17] and coded the participants' responses into the following categories: i) high, ii) average, and iii) low self-confidence. This categorization resulted in balanced classes: 17 participants had high self-confidence using LIME (8) and SHAP (9). 13 participants had average self-confidence using LIME (6) and SHAP (7), and 10 participants had low self-confidence using LIME (10) and SHAP (6).

We further compared the participants' scores with their self-confidence category; we found that SHAP's participants' confidence category matches with their scores, indicating that participants who had high self-confidence using SHAP also correctly answered the assignments.

However, we observed no correlation between the participants' self-confidence and the LIME scores. The median comprehensibility score (median score = 14.0) for 7 participants with average self-confidence was higher than both high-confidence (median score = 10.0) and low-confidence (median score = 11.0). LIME's predictability scores for both assignments two and three stayed the same and did not decrease with the participants' confidence category.

We compared the relationship between the participants' expertise (years of experience in ML and data science) with their scores for each category. We noticed a positive correlation for participants using SHAP, indicating more experienced participants could better interpret the explanations than participants with lower expertise.

However, participants who used LIME and had low expertise achieved higher scores than the more experienced users. We considered the time participants needed to answer the questions and noticed that those who took more time achieved a higher score. The participants with high expertise and higher confidence often took less time to answer LIME's questions and had lower scores than the other user's who needed longer to answer the questions.

We noticed that participants using LIME agreed that misclassification information increased their confidence in their answers. However, participants using SHAP stated that the visualizations did not improve their confidence, and 6 participants stated that the misclassification information provided by SHAP did not help them at all and confused them more. Their feedback also correlates with their scores.

We continued our qualitative analysis and presented the participants with the following questions:

Q3. How much did the visualizations (of the second assignment) help to get an insight into the model and how it comes to its decisions?

We categorized the helpfulness and negative feedback of the participants into categories i) helpful to answer the test sample, ii) only helps to understand the explained samples, and iii) not helpful at all. We further clustered the negative feedback from the second and third categories to construct the improvement guidelines based on participants' needs.

Our coding resulted in having 36 of 46 participants (18 LIME and 17 SHAP) stating that the information was not enough to scale to a new sample. One participant who used LIME found the visualizations helpful to scale the information for a new sample. Ten participants (6 SHAP and 4 LIME) stated that they could not use the information and answered with very low confidence. We moved further to our fourth question to understand whether the misclassification and counter-factual samples help the participants gain more insight into the model's decision-making process.

Q4. How much did the explanations of counter-factual and misclassified samples help to answer the questions of the third assignment?

We cluster the participants' feedback based on their similarities to understand what confused them and caused them to fail in solving the assignments. We clustered their feedback into three categories; i) inconsistency of the explained values, ii) missing information from the visualizations, and iii) not at all helpful.

We first present the negative feedback from participants who used LIME; Two participants stated that assignments two and three were confusing. One had a low expertise rank, and the other had average expertise and no XAI experience. However, both scores significantly increased after the visualizations of misclassified samples and counter-factual explanations. One participant stated that *"the way all the information was presented at once made it very difficult to understand the differences of the samples"* (*user\_id=27*), and *"I do not understand what the visualizations are trying to point*

out, they are all very similar to each other and identifying the correlation of the feature values on the target classes was not intuitive at all. I always used the bar plot and compared my sample with the labeled samples." (user\_id=39).

Four LIME users were in the second\_category (missing information from the visualizations) and stated that they needed more explanations of more misclassified samples. They also mentioned that *"the visualizations helped me understand why and when a misclassification might happen. Still, it was not enough to assume the model's behavior confidently."* (user\_id=1). These participants also had significantly higher scores in the third assignment, indicating that the visualizations helped identify misclassifications regardless of their subjective feedback.

Finally, only one LIME user stated that the information was completely confusing. the user could not make sense of the inequality range given by LIME's visualization and why some features had the same inequality range, even though they were classified into two different classes. When we explained that the inequalities depend on the feature interactions and why these changes stay the same for one feature, users admitted that they understood it. Still, it was not intuitive for the user at the assignment time.

We conclude that the reason behind most users' confusion was the unexplained, inconsistent range of inequalities for samples from the same class. Moreover, LIME's local explanations only present that a feature value reduces the probability of the predicted class but does not reveal which class's probability increases. SHAP plots present this information.

We move on to negative feedback given by participants who used SHAP; No participant stated that SHAP's visualizations were not at all helpful. Only two stated that more explanations of misclassified samples would have helped them scale the explanations for a new sample (category (ii)) and had a non-significant lower score  $t=1.0$ ,  $p=0.5$ ) for assignment three. Moreover, eight participants, who also achieved a higher score for the third assignment, mentioned that the visualizations confused them very much: *"There were too many bars and numbers and finding the contribution of features to each class and the effect of feature values on other classes was very demanding and tiring"* (user\_id=26). All these participants also stated that the inconsistency of the plot's scaling fooled them, and they had to invest more time to find the real contribution of feature values towards each class.

We noticed that for SHAP, it still plots a long bar for each feature when all features have a low contribution. If the user does not look at the probability scales carefully, they might assume that all these features are equally contributing highly towards the respective class.

## 8 DISCUSSION & DESIGN RECOMMENDATION

We now summarize the key findings of our user study and propose a set of design recommendations that can substantially contribute to the design of new or refinement of existing XAI methods:

- *Explanations should be consistent.* This may seem obvious, but when evaluating LIME, we noticed that many participants were confused by inconsistent inequality ranges for the same feature in different samples of the same class (e.g., model predicted two instances as low-price because one instance had  $0.26 < F1\text{-value} < 0.36$ , while another instance had  $F1\text{-value} < 0.26$ ). Such inconsistencies could be explained by showing the correlation between feature values.
- *Explanations should have fixed scales.* Participants working with SHAP tended to interpret feature explanations incorrectly for features with smaller scales. The visualizations "fooled" them, and they assumed that a feature with a higher bar has a greater influence on a sample, even though its SHAP value was much smaller than in another sample.

- *Explanations should also provide counterfactual examples.* Our quantitative and qualitative results show that participants predict model behavior better using the explanations if counterfactual samples are presented. These additional explanations could also help users understand the differences between the classes more clearly.
- *Explanations should contain misclassified samples close to the decision boundary.* From our experiments, we learned that participants predict model’s decisions with explanations better when they see misclassified examples and understand why the model made a wrong decision. This can be achieved by presenting samples close to a model’s decision boundary, which of course, are often subject to misclassification.
- *Explanations should contain correctly predicted samples close to the decision boundary.* Our results also indicate that participants achieved higher scores when being presented with explanations on correct predictions of samples. These samples can be identified by selecting correctly classified samples with low prediction confidence.

Our work is, of course, currently limited to comparing the local explanations of two XAI approaches, LIME and SHAP. Also, the Boston housing dataset we used in our experiments has been simplified to tabular data with a few dimensions. Moreover, users neither had the option to choose the samples themselves nor did we allow them to interact with the model and different outputs of the XAI approaches. However, we argue that these restrictions were necessary to reduce the number of confounding variables in our user study, which could influence our variables by events that are not causally related. We also believe that our experimental setup provides the necessary degree of generalizability to be transferred to the evaluation of other explanation tools.

Potential future work could expand our approach to other types of data (e.g., acoustic or sensor data) and other emerging XAI techniques (e.g., LORE [4]), we have not yet considered in our explanations. Moreover, one could compare local and global explanation designs and study how they affect users’ comprehensibility and predictability. Another potentially interesting research direction is to investigate how active learning techniques could be used to support users in comprehending and predicting model decisions [22].

Overall, we believe that user studies should become an integral part of the improvement of existing and the development of new XAI methods. Since explanations must ultimately be understood by users, their perceptions and interpretations of explanations should also be systematically analyzed and understood. This can improve the XAI methods and also the user interaction with these methods.

## 9 CONCLUSION

In this paper, we conducted a user study to investigate how well users comprehend the explanations and predict model behavior provided by two widely used tools, LIME, and SHAP. We measured the comprehensibility and predictability participants gained after interpreting a given set of local explanations to increase the comprehensiveness of the captured information for a new un-explained sample. We formed our first research question to measure comprehensibility and predictability. We showed that the comprehensibility of SHAP explanations significantly decreases for samples close to the decision boundary. Second, we studied the information participants require to gain a more global interpretation of the model behavior and to increase their predictability using explanations. We observed that explaining misclassified and counterfactual samples to the participants can significantly improve their predictability (especially with LIME explanations). They recognized the model behavior for unexplained samples close to the model decision boundary. Furthermore, our qualitative analysis of participants’ feedback indicated that they require information such as justifying the explained values (LIME inequalities or SHAP values) to correctly interpret the model behavior and move towards a more global interpretation of the model decision boundary. Finally, we learned that the users’ confidence in interpreting

the explanations strongly relies on the diversity and quantity of the explained samples. The more different instances were studied by participants, the more accurately they could interpret the outputs of the explainability approaches and predict the decisions of the model.

## 10 ACKNOWLEDGMENTS

We thank all the volunteers, and all the reviewers, who wrote and provided helpful comments on previous versions of this document. We specially thank our colleagues, Clemens Heistracher and Denis Katic for their constructive feedback on the structure of this work. We further like to thank Dr. Jasmin Lampert for her constructive feedback and insight to this work. We also thank the Austrian Research Promotion Agency (FFG) for funding this work, which is a part of the industrial project DeepRUL, project ID 871357.

## REFERENCES

- [1] Ahmed Alqaraawi, Martin Schuessler, Philipp Weiß, Enrico Costanza, and Nadia Berthouze. 2020. Evaluating Saliency Map Explanations for Convolutional Neural Networks: A User Study. In *Proceedings of the 25th International Conference on Intelligent User Interfaces (Cagliari, Italy) (IUI '20)*. Association for Computing Machinery, New York, NY, USA, 275–285. <https://doi.org/10.1145/3377325.3377519>
- [2] Diogo V Carvalho, Eduardo M Pereira, and Jaime S Cardoso. 2019. Machine learning interpretability: A survey on methods and metrics. *Electronics* 8, 8 (2019), 832.
- [3] Andy P Field, Jeremy Miles, and Zoë Field. 2012. *Discovering statistics using R*. SAGE publications, London, England. 361–365 pages.
- [4] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Dino Pedreschi, Franco Turini, and Fosca Giannotti. 2018. Local rule-based explanations of black box decision systems. *arXiv preprint arXiv:1805.10820* (2018).
- [5] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. 2018. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)* 51, 5 (2018), 1–42.
- [6] David Harrison Jr and Daniel L Rubinfeld. 1978. Hedonic housing prices and the demand for clean air. *Journal of environmental economics and management* 5, 1 (1978), 81–102.
- [7] SG Hart et al. 1988. Development of NASA-TLX: Results of empirical and theoretical research.” inP. A. Hancock and N. Meshkati (eds.), *Human Mental Workload*.
- [8] Peter Hase and Mohit Bansal. 2020. Evaluating Explainable AI: Which Algorithmic Explanations Help Users Predict Model Behavior? *Association for Computational Linguistics (ACL) (2020)*.
- [9] Stefan Haufe, Frank Meinecke, Kai Gørgen, Sven Dähne, John-Dylan Haynes, Benjamin Blankertz, and Felix Bießmann. 2014. On the interpretation of weight vectors of linear models in multivariate neuroimaging. *Neuroimage* 87 (2014), 96–110.
- [10] Robert R Hoffman. 2017. A taxonomy of emergent trusting in the human–machine relationship. *Cognitive systems engineering: The future for a changing world (2017)*, 137–164.
- [11] Alon Jacovi, Ana Marasović, Tim Miller, and Yoav Goldberg. 2021. Formalizing trust in artificial intelligence: Prerequisites, causes and goals of human trust in AI. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 624–635.
- [12] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jennifer Wortman Vaughan. 2020. Interpreting Interpretability: Understanding Data Scientists’ Use of Interpretability Tools for Machine Learning. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems (Honolulu, HI, USA) (CHI '20)*. Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3313831.3376219>
- [13] Himabindu Lakkaraju and Osbert Bastani. 2020. “How Do I Fool You?”: Manipulating User Trust via Misleading Black Box Explanations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (New York, NY, USA) (AIES '20)*. Association for Computing Machinery, New York, NY, USA, 79–85. <https://doi.org/10.1145/3375627.3375833>
- [14] Zachary C Lipton. 2018. The Mythos of Model Interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue* 16, 3 (2018), 31–57.
- [15] Scott M. Lundberg and Su-In Lee. 2017. A Unified Approach to Interpreting Model Predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (Long Beach, California, USA) (NIPS '17)*. Curran Associates Inc., Red Hook, NY, USA, 4768–4777.
- [16] David W. Martin. 2007. Doing Psychology Experiments. (2007), 148–170.
- [17] Philipp Mayring. 2004. Qualitative content analysis. *A companion to qualitative research* 1, 2 (2004), 159–176.
- [18] Sina Mohseni. 2019. Toward Design and Evaluation Framework for Interpretable Machine Learning Systems. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. 553–554.
- [19] Sina Mohseni and Eric D Ragan. 2020. A human-grounded evaluation benchmark for local explanations of machine learning. *arXiv preprint arXiv:1801.05075* (2020).



- [20] Sina Mohseni, Niloofar Zarei, and Eric D Ragan. 2021. A multidisciplinary survey and framework for design and evaluation of explainable AI systems. *ACM Transactions on Interactive Intelligent Systems (TiiS)* 11, 3-4 (2021), 1–45.
- [21] Christoph Molnar, Giuseppe Casalicchio, and Bernd Bischl. 2020. Interpretable Machine Learning—A Brief History, State-of-the-Art and Challenges. *Koprinska I. et al. (eds) ECML PKDD Workshops. ECML PKDD 2020. Communications in Computer and Information Science, vol 1323. Springer, Cham.* (2020).
- [22] Ishani Mondal and Debasis Ganguly. 2020. ALEX: Active Learning based Enhancement of a Classification Model’s EXplainability. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management.* 3309–3312.
- [23] Mahsan Nourani, Samia Kabir, Sina Mohseni, and Eric D Ragan. 2019. The effects of meaningful and meaningless explanations on trust and perceived system accuracy in intelligent systems. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, Vol. 7. 97–105.
- [24] Andrea Papenmeier, Gwenn Englebienne, and Christin Seifert. 2019. How model accuracy and explanation fidelity influence user trust. *AI IJCAI Workshop on Explainable Artificial Intelligence* (2019).
- [25] Qualtrics. [n.d.]. Copyright Year: 2021, Location: Provo, Utah, USA. <https://www.qualtrics.com>
- [26] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "Why should I trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.* 1135–1144.
- [27] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018. Anchors: High-Precision Model-Agnostic Explanations.. In *Proceedings of the 32nd AAAI International Conference on Artificial Intelligence*, Vol. 18. 1527–1535.
- [28] Stefan Rüping. 2006. Learning interpretable models. (2006). <http://dx.doi.org/10.17877/DE290R-8863>
- [29] Philipp Schmidt and Felix Biessmann. 2019. Quantifying interpretability and trust in machine learning systems. *AAAI-19 Workshop on Network Interpretability for Deep Learning* (2019).
- [30] Donghee Shin. 2021. The effects of explainability and causability on perception, trust, and acceptance: Implications for explainable AI. *International Journal of Human-Computer Studies* 146 (2021), 102551.
- [31] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning Important Features through Propagating Activation Differences. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70* (Sydney, NSW, Australia) (*ICML '17*). JMLR.org, Sydney, NSW, Australia, 3145–3153.
- [32] Kacper Sokol and Peter Flach. 2020. Explainability Fact Sheets: A Framework for Systematic Assessment of Explainable Approaches. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) (*FAT\* '20*). Association for Computing Machinery, New York, NY, USA, 56–67. <https://doi.org/10.1145/3351095.3372870>
- [33] Thilo Spinner, Udo Schlegel, Hanna Schäfer, and Mennatallah El-Assady. 2019. explAIner: A visual analytics framework for interactive and explainable machine learning. *IEEE transactions on visualization and computer graphics* 26, 1 (2019), 1064–1074.
- [34] Ian Tenney, James Wexler, Jasmijn Bastings, Tolga Bolukbasi, Andy Coenen, Sebastian Gehrmann, Ellen Jiang, Mahima Pushkarna, Carey Radebaugh, Emily Reif, and Ann Yuan. 2020. The Language Interpretability Tool: Extensible, Interactive Visualizations and Analysis for NLP Models. (2020), 107–118. <https://www.aclweb.org/anthology/2020.emnlp-demos.15>
- [35] Sana Tonekaboni, Shalmali Joshi, Melissa D. McCradden, and Anna Goldenberg. 2019. What Clinicians Want: Contextualizing Explainable Machine Learning for Clinical End Use. In *Proceedings of the 4th Machine Learning for Healthcare Conference (Proceedings of Machine Learning Research, Vol. 106)*, Finale Doshi-Velez, Jim Fackler, Ken Jung, David Kale, Rajesh Ranganath, Byron Wallace, and Jenna Wiens (Eds.). PMLR, Ann Arbor, Michigan, 359–380. <http://proceedings.mlr.press/v106/tonekaboni19a.html>
- [36] Danding Wang, Qian Yang, Ashraf Abdul, and Brian Y. Lim. 2019. *Designing Theory-Driven User-Centric Explainable AI*. Association for Computing Machinery, New York, NY, USA, 1–15. <https://doi.org/10.1145/3290605.3300831>
- [37] Jianlong Zhou, Amir H Gandomi, Fang Chen, and Andreas Holzinger. 2021. Evaluating the quality of machine learning explanations: A survey on methods and metrics. *Electronics* 10, 5 (2021), 593.