
Learning to Reject Low-Quality Explanations via User Feedback

Luca Stradiotti

DTAI lab & Leuven.AI,
KU Leuven, Belgium
luca.stradiotti@kuleuven.be

Dario Pesenti

CIMeC,
University of Trento, Italy
dario.pesenti@unitn.it

Stefano Teso

CIMeC and DISI,
University of Trento, Italy
stefano.teso@unitn.it

Jesse Davis

DTAI lab & Leuven.AI,
KU Leuven, Belgium
jesse.davis@kuleuven.be

Abstract

Machine Learning predictors are increasingly being employed in high-stakes applications such as credit scoring. Explanations help users unpack the reasons behind their predictions, but are not always “high quality”. That is, end-users may have difficulty interpreting or believing them, which can complicate trust assessment and downstream decision-making. We argue that *classifiers should have the option to refuse handling inputs whose predictions cannot be explained properly* and introduce a framework for **learning to reject low-quality explanations** (LtX) in which predictors are equipped with a *rejector* that evaluates the quality of explanations. In this problem setting, the key challenges are how to properly define and assess explanation quality and how to design a suitable rejector. Focusing on popular attribution techniques, we introduce ULER (User-centric Low-quality Explanation Rejector), which learns a simple rejector from human ratings and per-feature relevance judgments to mirror *human* judgments of explanation quality. Our experiments show that ULER outperforms both state-of-the-art and explanation-aware learning to reject strategies at LtX on eight classification and regression benchmarks and on a new human-annotated dataset, which we publicly release to support future research.

1 Introduction

Machine Learning (ML) predictors are increasingly deployed in *high-stakes* decision-making applications, such as medical diagnosis and credit scoring [1–3]. In these domains, incorrect predictions can lead to severe consequences [4]. To promote trust, *Learning to Reject* (LtR) strategies have emerged that allow models to defer predictions to human experts if the model has an elevated risk of making a misprediction [5–7]. Traditional LtR approaches typically abstain when the model is uncertain about its prediction or the input differs from the observed training data [8, 9].

Currently, LtR focuses on predictive performance which neglects a critical aspect of decision-making: *explanation quality* [10, 11], cf. Fig. 1 (left). In many applications, it is equally important that models provide clear and convincing explanations for their predictions [12]. Without addressing explanation quality, a model might make predictions that cannot be satisfactorily explained [13]. We argue that low-quality explanations can affect trust assessment and downstream decisions [12, 14–16]. As a consequence, we believe models should *offload predictions that they cannot properly explain* to human stakeholders, thus ensuring that predictions are based on human-validated reasoning and

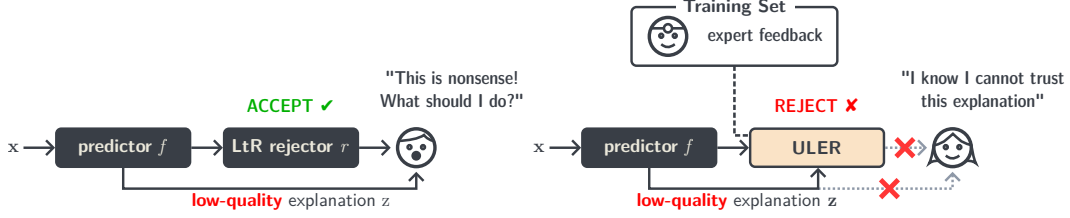


Figure 1: **Illustration of ULER.** Learning to Reject (LtR) is unconcerned with the quality of machine explanations (left). ULER *instead addresses Learning to Reject Low-Quality Explanations (LtX)*, which requires to reject predictions that cannot be explained properly to stakeholders, improving trust assessment and down-stream decision quality (right).

preserving the overall trustworthiness of the system. This perspective aligns with the Four Principles of Explainable Artificial Intelligence [13], an official document from the U.S. government, which emphasizes the importance that an AI system recognizes and declares its knowledge limits. According to the authors, “safeguarding answers so that a judgment is not provided when it may be inappropriate to do so” can prevent “misleading, dangerous, or unjust outputs”. E.g., consider a general practitioner that uses an AI system to assist in diagnosing malignant melanoma. When examining a suspicious lesion, the AI correctly advises against further action, citing the size of the lesion as a key factor, which is irrelevant in the doctor’s opinion. Distrusting the AI’s explanation, the doctor decides to proceed with additional examinations, resulting in unnecessary costs and delays

To formalize this notion, we introduce the *Learning to Reject Low-Quality Explanations (LtX)* problem where a model should abstain from making a prediction when it can only provide an unsatisfactory explanation from the user’s perspective, cf. Fig. 1 (right). This is a challenging problem that current techniques cannot adequately address. On the one hand, LtR focuses only on prediction quality but just because a model can offer a correct prediction does not imply it can offer an acceptable explanation for it. On the other hand, existing metrics for evaluating explanations do so on the basis of properties of the model. Consequently, these may not agree with a human’s assessment of the quality of the explanation.

To address the LtX problem, we propose ULER (User-centric Low-quality Explanation Rejector) to train a novel type of rejector to assess the quality of an explanation from a user’s perspective. It does so by leveraging expert annotations comprising quality judgments and per-feature relevance judgments. ULER consists of two main steps. First, to avoid having to collect a large number of explanation judgments, we apply a novel quality-aware augmentation strategy that exploits the human annotations to augment the training set. Second, we fit the rejector to evaluate the explanations’ quality using the augmented quality judgment labels. Empirically, we demonstrate that ULER outperforms many popular LtR strategies as well as approaches to estimate the quality of the explanation on both the machine and human side. We also show the effectiveness of ULER on a novel larger-scale dataset of machine explanations accompanied by human annotations obtained running a user study, which we make available to support the analysis and development of LtX solutions.

Contributions: Summarizing, we: (i) Introduce the problem of *learning to reject low-quality explanations (LtX)*, filling a significant gap in current LtR strategies, which ignore explanation quality altogether. (ii) Design ULER, a rejector that uses modest amounts of human annotations – including explanation ratings and per-feature relevance judgments – to learn an effective rejection policy. (iii) Empirically evaluate ULER on both popular data sets and on a novel human-annotated task, showcasing its benefits over standard LtR and state-of-the-art explanation quality metrics. (iv) Provide the first larger-scale (1050 examples, 5 annotations each) data set of human-annotated machine explanations as well as a template for running the associated collection campaign.

2 Preliminaries

We begin by introducing the basic framework used throughout the paper. We consider a *predictor* f that maps inputs $\mathbf{x} \in \mathcal{X}$ to a target value $f(\mathbf{x}) \in \mathcal{Y}$. Here, \mathcal{X} is a d -dimensional feature space and \mathcal{Y} a discrete ($\mathcal{Y} = \{1, \dots, C\}$) or continuous ($\mathcal{Y} = \mathbb{R}$) target space. When the target is discrete, we

view the predictor as a probabilistic *classifier* that assigns a predictive distribution $P(Y|X = \mathbf{x})$ to each input \mathbf{x} ; predictions are obtained via MAP inference, that is $f(\mathbf{x}) = \operatorname{argmax}_{c \in \mathcal{Y}} P(Y = c|\mathbf{x})$ [17]. When the target is continuous, we view it as a *regressor* $f(\mathbf{x}) = \mathbb{E}[Y|X = \mathbf{x}]$.

In the following, we assume the predictor is paired with an *explainer* e which produces a local explanation $\mathbf{z} = e(f(\mathbf{x}))$ of individual prediction $f(\mathbf{x})$. Specifically, we focus on *feature importance* explanations, perhaps the most well-known and widespread class of explanations [18–25]. These associate a *relevance score* $z_i \in \mathbb{R}$ to each input feature x_i that quantifies its relative contribution for the prediction. For example, in loan approval, \mathbf{z} might indicate that an application \mathbf{x} was rejected (*i.e.*, $f(\mathbf{x}) = 0$) because a specific feature x_{income} , which is too low, “votes” against approval by assigning it a negative value (*i.e.*, $z_{\text{income}} < 0$). We refer to the pair $(f(\mathbf{x}), \mathbf{z})$ as the model *output*, since each prediction $f(\mathbf{x})$ is returned to the user along with its corresponding explanation \mathbf{z} .

Learning to reject. To promote trust, a **Learning to Reject** (LtR) model combines a predictor f with a *rejector* r . The role of the rejector is to offload difficult predictions to a human expert [26, 7, 5]. Formally, it does so by extending the target space \mathcal{Y} to include an additional symbol \mathbb{R} indicating the model abstains from making a prediction [6, 27]. Two classes of rejection strategies have been studied in the literature. **Ambiguity rejection** occurs when the predictor f is too uncertain about a particular input \mathbf{x} , *e.g.*, due to class overlap or poor choice of the predictor’s hypothesis space [27, 28]. **Novelty rejection** checks if \mathbf{x} falls in a region where there is little or no training data [29, 30]. Although such rejection strategies improve the model’s reliability [31, 32], they focus solely on predictor’s performance [5] and ignore cases where the explanations themselves are unsatisfactory to the user.

Metrics of Explanation Quality. Several metrics have been proposed to evaluate the quality of an explanation [33]. Most of them depend solely on the relationship between the explanation and the predictor and, as such, can be computed accurately using information gathered during inference and/or training. For example, **faithfulness** [34–36] measures whether an explanation accurately reflects the model’s reasoning process, and it is typically computed by assessing whether the features with high relevance are sufficient and necessary for the prediction. Another key metric is **stability** [37, 38], which measures the degree to which different (possibly conflicting) explanations can be provided for a given prediction. Despite their utility, recent works [39, 40] have shown that *these metrics do not align with human judgment*, highlighting the need for alternatives. An exception is PASTA, a novel perceptual quality metric that mimics human preferences across multiple dimensions [39] and that we compare against in our experiments (Section 4). A deeper discussion of all these is provided in Appendix B. Although several metrics of explanation quality exist, none have been integrated into rejection strategies to guide the rejector’s decisions. To address this gap, next we introduce a novel framework that incorporates user-perceived explanation quality into the rejection process.

3 Learning to Reject Low-Quality Explanations

We introduce the *Learning to Reject Low-Quality Explanations* (LtX) problem where a rejector acts as a filter based on the user-perceived explanation quality [41, 42, 16]. Consequently, the rejector in this setting operates on \mathbf{z} as opposed to $f(\mathbf{x})$ or \mathbf{x} as in a standard LtR setting. Formally, a model with reject option in the LtX setting is defined as follows.

Definition 1 (A LtX Model) *An LtX model m consists of three components: a predictor f , an explainer e and a rejector r . Given (test) instance \mathbf{x} , m computes $f(\mathbf{x})$ and corresponding explanation $e(f(\mathbf{x}))$. Then, m applies the rejector r to $e(f(\mathbf{x}))$ to assign a score representing the quality of the explanation \mathbf{z} with lower scores being associated with worse explanations. If the score is below a threshold τ , the model abstains from providing the prediction and the corresponding explanation to the user. Formally, m is defined as:*

$$m_{(f,e,r)}(\mathbf{x}) = \begin{cases} \mathbb{R} & \text{if } r(\mathbf{z}) < \tau \\ (f(\mathbf{x}), \mathbf{z}) & \text{otherwise} \end{cases} \quad (1)$$

Our key contribution is to learn a rejector that abstains when e provides a low quality explanation from the user’s perspective. Obtaining such a rejector is challenging for three reasons. First, LtR strategies determine when the model should abstain based on where the predictor is likely to make a mistake. However, the predictor may still output a correct prediction even when the corresponding explanation is unreliable, and as such they cannot be used as-is. Second, existing metrics to evaluate

explanations focus only on the model’s internal functioning and are not able to measure the quality of the explanation from the user’s perspective, as we will show empirically in [Section 4](#). Third, training a standard LtR model only requires standard supervised dataset consisting of instances and their target values. In contrast, LtX requires human-judgment labels about the explanations of each prediction which are usually not available and may be time-consuming to obtain.

3.1 Rejecting low-quality explanations with ULER

We propose a novel approach for the LtX problem called ULER (User-centric Low-quality Explanation Rejector) that addresses the aforementioned challenges by (i) collecting a small set of user annotated explanations, (ii) employing a feedback-driven data augmentation strategy, and (iii) training a rejector that estimates the user-perceived quality of an explanation. We detail these steps next.

The rejector’s training data. ULER assumes access to two sources of expert feedback. The first one is a set of explanations and corresponding *human quality judgments*. This set can be formalized as $\mathcal{D} = \{(\mathbf{z}_1, y_{\mathbf{z}_1}), \dots, (\mathbf{z}_n, y_{\mathbf{z}_n})\}$, where \mathbf{z} are the explanations, and $y_{\mathbf{z}} \in \{0, 1\}$ their corresponding human quality judgments (0 = low-quality, 1 = high-quality). This feedback is essential for training an LtX rejector that is aligned with expert judgments of explanation quality. Yet, such annotations can be expensive to acquire and therefore typically available in modest amounts [43, 39].

For this reason, ULER exploits *per-feature human labels*, a second and more detailed source of information, to augmentation the available quality judgments. The per-feature labels indicate, for each explanation \mathbf{z} in \mathcal{D} , what relevance scores the user deems incorrect, if any. Formally, we indicate as $\mathcal{W}_{\mathbf{z}}$ (resp. $\mathcal{C}_{\mathbf{z}}$) the indices of the features whose relevance the user deems *wrong* (resp. *correct*). In [Section 4.2](#), we show how to design an annotation campaign to obtain both kinds of feedback.

Augmenting the data. The *augmentation step* works by perturbing each $(\mathbf{z}, y_{\mathbf{z}}) \in \mathcal{D}$ using a stochastic transformation that leverages the per-feature labels while keeping $y_{\mathbf{z}}$ fixed. If explanation \mathbf{z} is deemed to be high-quality, we stipulate that slightly perturbing the small number of incorrect feature relevance scores will not affect the explanation label. The situation for low-quality explanations is reversed: we can perturb the *correct* entries of \mathbf{z} and obtain a different low-quality explanation. Formally, for each explanation \mathbf{z} we create K new explanations \mathbf{z}_{aug} sharing the same human-judgment label $y_{\mathbf{z}}$ as:

$$\mathbf{z}_{aug} \sim \mathcal{N}(\mathbf{z}, \epsilon_0 \mathbf{s} \times \Sigma) \quad (2)$$

Here, ϵ_0 is a hyperparameter controlling the overall magnitude of the perturbations, Σ is a diagonal matrix whose elements are the per-feature standard deviations across all explanations in \mathcal{D} and is responsible for rescaling perturbations compatibly with the data distribution, and \mathbf{s} is a binary vector used to selectively perturb the features depending on $y_{\mathbf{z}}$, $\mathcal{W}_{\mathbf{z}}$, and $\mathcal{C}_{\mathbf{z}}$. In practice, if the explanation is low-quality, the entries of \mathbf{s} corresponding to the indices in $\mathcal{C}_{\mathbf{z}}$ are set to 1 and those in $\mathcal{W}_{\mathbf{z}}$ to 0. The opposite happens if the explanation is high-quality. Our experiments support the small annotation cost of the augmentation step, as empirically shown in [Section 4.1](#).

Learning the rejector. The two key components of the rejector are a binary classifier r and a threshold τ . ULER trains a binary classifier on the augmented data \mathcal{D}_{aug} . ULER is agnostic to the specific choice of classifier: any model class that associates a score with its prediction is possible. Empirically, we find that simple models (e.g., kernel SVMs [44]) work well. τ determines how often a prediction and explanation are offered by m . Lower values of τ mean that m will operate more autonomously (i.e., return more prediction-explanation pairs) albeit with the risk that some explanations are low quality. Higher values mean the model is more cautious and only offers predictions-explanation pairs when its more certain about the quality of the explanation but at the cost of offloading more decisions to the user. Hence, this value should be carefully tuned, e.g., on validation to navigate this tradeoff. Two natural strategies are to set τ such that (i) it achieves a specific rejection rate on the validation data (e.g., one aligned with a user’s capacity to make decisions) or (ii) its rejection rate is equal to the proportion of low-quality explanations in the training set.

3.2 Benefits and Limitations

We remark that ULER is designed to identify and offload predictions associated with unsatisfactory explanations, as doing so is crucial for ensuring an accurate decision making. As such, ULER only yields a marginal improvement in predictive performance for the accepted inputs by rejecting

explanations of incorrect predictions, specifically when these errors occur due to a mismatch between the expert’s and the model’s reasoning. One option is however to combine ULER with state-of-the-art LtR strategies designed for this different goal. One limitation of ULER is that, just like PASTA [39], it relies on high-quality human annotations. We argue that this is necessary in high-stakes applications, but also that good annotations are likely to be available anyhow as in these settings expert users *have* to oversee machine decisions at all times [41, 45, 46], and can therefore consistently supply high-quality responses. Our experiments in Section 4 indicate that ULER is quite sample efficient, as it outperforms the SOTA while using less than 1000 annotations, and that augmentation boosts the performance of the rejector. Finally, our study focuses on tabular data rather than images or text. Working with a larger number of features may increase the sample complexity of the rejector. A possible solution is to adapt ULER to work in a rich pre-trained embedding space, as done by PASTA.

4 Empirical Evaluation

We answer empirically the following research questions: **(Q1)** Does ULER reject more low-quality explanations than the competitors? **(Q2)** What inputs does ULER need to learn to evaluate the explanation quality? **(Q3) (User study)** Is ULER capable of mimicking human judgments?

Competitors. We compare ULER against *eight* representative rejection strategies from two groups: (i) standard LtR strategies, and (ii) explanation-aware strategies. All strategies yield a score for each input; the $\rho\%$ inputs with the lowest score are rejected, where $\rho\%$ is the *rejection rate*.

We consider *three standard LtR strategies*. RandRej is a baseline that assigns a random score to each input. NovRej_X rejects inputs based on their novelty [29]: it first computes their distance to the k -th nearest training instances and converts these into scores using a monotonically decreasing function, e.g., $1/(1+x)$, such that farthest inputs get lower scores. PredAmb uses prediction’s confidence as score [5]. For binary classification tasks, confidence is computed as the margin of the class probabilities $|P(Y=1|\mathbf{x}) - P(Y=0|\mathbf{x})|$ [28]. For regression tasks, the conditional variance for each input is computed and then the score is obtained applying a monotonically decreasing function, e.g., $1/(1+x)$, such that higher-variance predictions obtain lower scores [47].

We consider *five novel but natural explanation-aware strategies*. Three leverage machine-side explanation metrics as scores, one for each category in [33]. Specifically, StabRej looks at the stability of the explanation [35], measuring the similarity among the different explanations that can be generated for the same prediction. FaithRej assesses the faithfulness [36] of an explanation by measuring how well the explanation identifies features that are truly causally relevant for the prediction. CompRej measures the complexity [48] of an explanation *i.e.*, the cognitive load it enforces on a user; since low-complexity explanations are preferred, the score is obtained applying a monotonically decreasing transformation, e.g., $1/(1+x)$, to the metric value. PASTARej uses an adaptation of the state-of-the-art human-side PASTA-metric to score each explanation [39]. Since our focus is on tabular data, we drop the embedding network and fit only the scoring network using the explanations as input to learn the human-judgment. Full details on all metrics are provided in Appendix B. Finally, NovRej_Z mirrors NovRej_X but works in the explanation space, testing whether the perceived low-quality explanations correspond to outlier explanations.

Evaluation metrics. Ideally, at test time, a user wants to observe only predictions that are accompanied by high-quality explanations. A good rejector should therefore minimize the number of low-quality explanations it shows to the user (*accepted set*), and maximize the ones for which it abstains (*rejected set*). Thus, we report the percentage of low-quality explanations in the accepted and rejected sets when varying the rejection rate. Additionally, we measure the rejector’s ability to effectively rank low-quality explanations below high-quality ones, making them more likely to be rejected, by reporting the AUROC, which is also standard in novelty rejection [29, 9, 49].

Setup. We employ the following procedure: for each dataset, we (i) split \mathcal{D} into \mathcal{D}_{train} , \mathcal{D}_{val} and \mathcal{D}_{test} (70%/10%/20%), (ii) fit the rejectors on \mathcal{D}_{train} and optimize their hyperparameters on \mathcal{D}_{val} , (iii) vary the rejection rate $\rho\%$ from 1% to 25%, and (iv) compute the metrics outlined in the previous paragraph on \mathcal{D}_{test} . To improve robustness, we repeat steps (i)–(iv) 10 times and report the average results. All experiments were implemented in Python and executed on an Intel i7-12700 machine with 64 GB RAM. The experiments required approximately two days to complete.

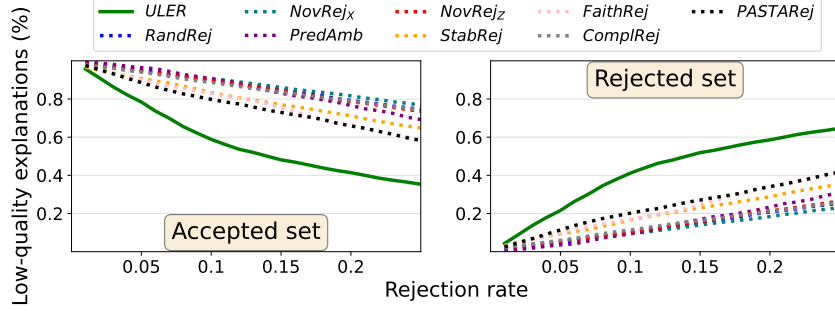


Figure 2: ULER rejects on average more low-quality explanations than all competitors. Average percentage of low quality explanations in the accepted and rejected set for all the considered strategies over the 8 datasets for 25 rejection rates $\rho\%$. For all the considered rejection rates, ULER consistently rejects more low-quality explanations than all competitors.

Model selection. All explanations are computed using *KernelSHAP* [21] with 100 samples and the predictor’s training set as background. We choose *KernelSHAP* as it is one of the most well-known and widely used explainers. To further support our findings, we also include results using *LIME* [20] in Appendix C.3. For ULER, we train an SVM to assess explanation quality. As mentioned, we optimize ULER’s and the competitors’ hyperparameters via grid search on \mathcal{D}_{val} , see Appendix C.2 for details.

4.1 Q1 and Q2: Benchmark Datasets

Datasets. We evaluate all competitors on *eight* widely used benchmarks datasets [50] using simulated human judgments. Since our approach works for any type of prediction function, we select four classification tasks and four regression tasks covering several application domains, including healthcare (*parkinson*), economics (*creditcard*, *adult*), industry (*appliances*), law (*compas*), etc. (*wine*, *news*, *churn*). Full details about the datasets are provided in Appendix C.1.

Simulating human judgments. We simulate human quality judgments $Y_{\mathbf{z}}$ and identify features with incorrect relevance scores using a ML oracle \mathcal{O} . Specifically, we train a predictor \mathcal{O} and use its explanations $\mathbf{z}_{\mathcal{O}}$ as a surrogate for those that an expert would provide. Then we train the proper predictor (that is, f) and classify its explanations \mathbf{z} as low- or high-quality depending on how much they correlate with the oracle’s explanation. In practice, for each classification (resp. regression) task, we train a Random Forest classifier (resp. regressor) to serve as the oracle \mathcal{O} and a linear SVC (SVR) as the proper predictor. All predictors use the default scikit-learn implementations [51]. We select predictors with different inductive biases to mirror real-world scenarios where human’s predictions may differ from model outputs. Both predictors are evaluated on a disjoint test set consisting of 2000 instances: the oracle achieves an average balanced accuracy (resp. MSE) of 0.76 (resp. 0.008), while the model of 0.69 (resp. 0.020); see Table 3 (Appendix) for full performance and training details. To verify the robustness of our findings, we replicate these experiments using different predictors in Appendix C.4.

Then, explanations for both the oracle and the predictor are generated on \mathcal{D} . An explanation \mathbf{z} is labeled as low-quality ($y_{\mathbf{z}} = 0$) if the correlation with the corresponding oracle’s explanation $\mathbf{z}_{\mathcal{O}}$ falls below a threshold $\tau_{\mathbf{z}}$, and as high-quality ($y_{\mathbf{z}} = 1$) otherwise. We fix $\tau_{\mathbf{z}} = 0.25$ as this ensures datasets with varying amount of low-quality explanations (1%-48%). Additionally, for each explanation \mathbf{z} , we construct the set of “wrong” relevance scores $\mathcal{W}_{\mathbf{z}}$ by selecting the scores in \mathbf{z} that deviate most from the corresponding scores in the oracle explanation $\mathbf{z}_{\mathcal{O}}$. Intuitively, if \mathbf{z} is low-quality, $\mathcal{W}_{\mathbf{z}}$ should include those entries that account for most of the difference between $\mathbf{z}_{\mathcal{O}}$ (which is high-quality by construction) and \mathbf{z} . To this end, we first compute the difference in relevance $|z_i - z_{\mathcal{O},i}|$ for each i , and then include in $\mathcal{W}_{\mathbf{z}}$ the indices i ’s with the highest difference and that cumulatively account for $u\%$ of the L_1 distance between $\mathbf{z}_{\mathcal{O}}$ and \mathbf{z} . Conversely, if \mathbf{z} is high-quality, a human may still identify some scores as incorrect. Here, we similarly construct $\mathcal{W}_{\mathbf{z}}$ but select features that cumulatively account for only $1 - u\%$ of the total L_1 -difference. We set $u\%$ to 0.75 in the experiments. Since we had sufficient data, we could afford to use non-overlapping sets to train the rejector and the predictor, although doing so is not strictly necessary.

Table 1: ULER **outperforms the competitors at separating low-quality from high-quality explanations**. Average AUROC for all the rejection strategies over the 8 datasets and its standard deviation. ULER consistently obtains the best results in all datasets.

	Classification				Regression			
	compas	creditcard	adult	churn	news	wine	parkinson	appliances
ULER	0.75 ± 0.04	0.87 ± 0.02	0.85 ± 0.04	0.92 ± 0.01	0.91 ± 0.01	0.93 ± 0.03	0.87 ± 0.01	0.82 ± 0.01
RandRej	0.52 ± 0.05	0.50 ± 0.02	0.53 ± 0.06	0.49 ± 0.02	0.50 ± 0.02	0.51 ± 0.07	0.49 ± 0.01	0.50 ± 0.02
NovRej _X	0.46 ± 0.04	0.58 ± 0.02	0.30 ± 0.05	0.36 ± 0.02	0.54 ± 0.01	0.51 ± 0.04	0.58 ± 0.02	0.56 ± 0.02
PredAmb	0.56 ± 0.03	0.46 ± 0.02	0.71 ± 0.03	0.85 ± 0.01	0.48 ± 0.04	0.50 ± 0.02	0.49 ± 0.02	0.50 ± 0.02
StabRej	0.69 ± 0.04	0.45 ± 0.02	0.53 ± 0.05	0.63 ± 0.02	0.62 ± 0.01	0.76 ± 0.04	0.53 ± 0.03	0.56 ± 0.02
FaithRej	0.63 ± 0.04	0.42 ± 0.02	0.71 ± 0.03	0.86 ± 0.01	0.64 ± 0.01	0.74 ± 0.05	0.49 ± 0.02	0.46 ± 0.02
ComplRej	0.69 ± 0.04	0.53 ± 0.05	0.45 ± 0.02	0.63 ± 0.02	0.76 ± 0.04	0.62 ± 0.01	0.56 ± 0.02	0.53 ± 0.03
PASTARej	0.52 ± 0.04	0.82 ± 0.03	0.66 ± 0.13	0.87 ± 0.02	0.61 ± 0.03	0.55 ± 0.10	0.61 ± 0.03	0.61 ± 0.03
NovRej _Z	0.46 ± 0.04	0.58 ± 0.02	0.57 ± 0.04	0.52 ± 0.02	0.50 ± 0.01	0.53 ± 0.05	0.57 ± 0.02	0.55 ± 0.01

(Q1) Comparison with competitors. Fig. 2 shows the percentage of low-quality explanations for the accepted and the rejected set as a function of the rejection rate $\rho\%$ averaged over the eight considered datasets. On average, ULER reduces the number of low-quality explanations in the accepted set by approximately 20% vs PASTARej, 21% vs FaithRej, 24% vs StabRej and ComplRej, 27% vs PredAmb, 31% vs RandRej, and 32% vs NovRej_X and NovRej_Z. Moreover, ULER rejects the highest number of low-quality explanations in around 94% of the experiments against all competitors. Finally, all the rejectors based on explanation metrics work better than the standard LtR strategies. This confirms that focusing on the prediction ambiguity or input novelty is not aligned with the objective of the LtX setting.

Table 1 reports the average AUROC per dataset. ULER performs better at separating low-quality from high-quality explanations for all the considered datasets and obtains an average improvement of 20% and 24% from the two runner-ups, respectively PASTARej and FaithRej. Importantly, this trend remain unchanged if we change the explainer (see Fig. 4 and Table 4 in the Appendix) or the oracle and the predictor (see Fig. 5): ULER *always compares favorably to the competitors in all cases*.

(Q2) ULER’s input space. To investigate which inputs the rejector needs to assess explanation quality, we consider three variants of our method, which train the rejector r using a different input space: ULER_{Z,X} uses both the explanation and its corresponding instance, ULER_{Z,Y} uses the explanation along with the prediction, and ULER_{Z,X,Y} uses the explanation, the instance, and the prediction.¹ For each variant, we first augment the explanations as described in Section 3.1. We then construct the training set by concatenating each (augmented) explanation with the input and/or prediction depending on variant. The rejector is subsequently trained on this enriched dataset to assess explanation quality.

Table 2 reports the average AUROC per dataset for ULER and each of the above variants. ULER consistently outperforms the other variants across all considered datasets. Interestingly, including the instances as part of the rejector’s input tends to decrease the performance: on average, ULER outperforms ULER_{Z,X} and ULER_{Z,X,Y} by approximately 16%. This is likely due to the limited number of human-judgment labels which makes it difficult for the rejector to learn the relationship between the explanations and the instances. Moreover, even concatenating only the prediction as in ULER_{Z,Y} results in a small performance hit (typically, around 3%), suggesting that explanations alone are often sufficient. The large drop in *news* is due to a sub-optimal performance of the predictor. For completeness, we also ablate the augmentation step. The average improvement is modest ($\approx 2\%$ across tasks) but consistent: ULER outperforms the ablated variant in 79% of the experiments. We argue augmentation is worth it in high-stakes applications, also because there the cost of obtaining per-feature feedback is small, cf. Section 3.2. When this is not the case, one could simply skip the augmentation: this variant of ULER still outperforms all the competitors, see Appendix C.5.

¹For classification tasks, we use the positive class probability rather than predicted label.

Table 2: ULER **outperforms its variants that additionally provide inputs and/or predictions as input to the rejector**. Average AUROC for ULER and three variants using different inputs to learn the quality of an explanation over the 8 datasets. ULER consistently achieves the highest AUROC across all datasets, showing that explanations alone suffice for the rejector to assess their quality effectively.

	Classification				Regression			
	compas	creditcard	adult	churn	news	wine	parkinson	appliances
ULER	0.75 ± 0.04	0.87 ± 0.02	0.85 ± 0.04	0.92 ± 0.01	0.91 ± 0.01	0.93 ± 0.03	0.87 ± 0.01	0.82 ± 0.01
ULER _{Z,X}	0.69 ± 0.04	0.79 ± 0.02	0.77 ± 0.05	0.87 ± 0.01	0.63 ± 0.02	0.63 ± 0.05	0.60 ± 0.02	0.69 ± 0.01
ULER _{Z,Y}	0.71 ± 0.05	0.86 ± 0.02	0.84 ± 0.04	0.90 ± 0.01	0.73 ± 0.02	0.90 ± 0.03	0.79 ± 0.01	0.81 ± 0.01
ULER _{Z,X,Y}	0.69 ± 0.04	0.79 ± 0.02	0.79 ± 0.04	0.87 ± 0.01	0.63 ± 0.02	0.64 ± 0.06	0.60 ± 0.02	0.70 ± 0.01

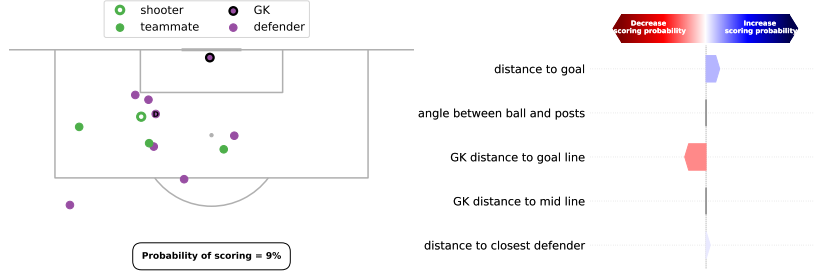


Figure 3: **Image from the user study** illustrating the snapshot (left), the predicted probability of scoring (bottom) and the associated KernelSHAP explanation (right). This suggests that the feature “distance to goal” slightly increases the probability, while “GK distance to goal line” decreases it.

4.2 Q3: ULER Predicts Human Judgments Better than the SOTA

Finally, we apply ULER to high-quality human ratings of machine explanations collected through a large-scale annotation campaign, using the crowd-sourcing platform Prolific (<https://www.prolific.com>) for recruitment.²

Our task was to explain the prediction of an expected goals (xG) model, which values the quality of a scoring opportunity in soccer as the probability that a shot results in a goal [52, 53]. Our choice stems from three considerations. First, Prolific enabled us to recruit subjects that possess the necessary domain expertise to perform the task, cf. Appendix D.3 for our vetting criteria. Second, all instances can be easily visualized, as shown in Fig. 3. Third, this is a real-world task with xG values being shown on TV and used in player recruitment.³ We collected annotations for 1050 explanations from five annotators each, for a total of 5250 annotations.

Obtaining the explanations. As a first step, we trained the predictor whose explanations we aim to annotate. Following standard practice in soccer analytics [52, 53], we learned an XGBoost ensemble classifier [54] to estimate the probability of a shot resulting in a goal. The training data consists of 21337 annotated shot events from the 2015-16 season in the top divisions of England, Spain, Germany and France [55]. For each shot, the location and the result (goal vs. no goal) are recorded. Additionally, a snapshot is available, capturing the locations of the players visible in the broadcast video at the moment the shot is taken, cf. Fig. 3 (left). From this data, we extract features that describe the positions of the shooter, goalkeeper, and nearest defender. Importantly, we include only features that are directly visualizable by the annotators in the snapshot. Explanations are generated on a separate set of 1050 shots from the 2015–16 season of the Italian top division on which the predictor achieves an AUROC of 0.81. All preprocessing and training details are provided in Appendix D.2.

Obtaining the annotations. Our goal is to obtain human-judgment labels on the explanation quality and per-feature feedback on the relevance scores. Given that subjective tasks are highly sensitive to interface design [56] and question framing [57], we designed our annotation protocol with the help

²The campaign has received approval from the Research Ethics Committee of the University of Trento (Protocol No. 2025-006ESA).

³The model used in our experiments is not as complex as deployed models.

of a psychologist and conducted several pilot studies to mitigate cognitive biases [58]. Participants ($N = 175$) were recruited via Prolific while annotations were collected through Google Forms. Each participant annotated 30 trials. In each trial, participants were shown a snapshot like in Fig. 3 depicting a shot and the corresponding prediction and explanation. The left side shows the position of all involved players and the ball, along with the model’s prediction. The right side shows the relevance scores of each feature as arrows indicating whether the feature increases or decreases the predicted probability of scoring. The features were chosen specifically to be easily interpretable and visually grounded, enabling intuitive assessment by the annotators. The annotators were requested to specify how much they agreed with the model’s prediction and, separately, with its explanation using two 5-point Likert-scale questions (1 = completely disagree, 5 = completely agree). Next, they were asked to optionally select individual features they believed were misused in the explanation, *i.e.*, had an incorrect relevance score, via a multiple-choice question. We validated our experimental design by tracking the consistency of individual annotations in two pilot studies: on average, annotators tended to assign consistent scores to the same explanation across repeated trials. Full details about the annotation protocol are provided in Appendix D.3.

Annotations preprocessing. To ensure annotations are high-quality, we filtered out participants who failed an attention check, those who rated all explanations with the same score, and those who did not flag any score as incorrect, leaving us with a total of 149 participants. For the same reason we also removed explanations with low inter-annotator agreement. We aggregated the explanation scores using the average and considered explanations with an average score lower than 3 as low-quality, and the others as high-quality [59, 60]. For feature-level feedback, we marked a relevance score as incorrect if the majority of annotators agreed that the corresponding feature was misused.

Results. Given the labels for explanation quality and feature-level feedback, we can now compare ULER and PASTARej. We consider only PASTARej as a competitor, as it is the only baseline that leverages human-judgment labels and emerged as the runner-up in previous experiments. ULER achieves an AUROC of 0.63 ± 0.05 , outperforming PASTARej, which scores 0.51 ± 0.09 . A paired t-test confirms that the difference is statistically significant ($p < 0.01$). These results indicate that learning human-perceived explanation quality is inherently challenging, especially in this subjective task. The overall low performance can be attributed to this increased variability. When annotations are provided by a single expert, consistency improves and performance approaches that observed in Table 2. Additionally, ULER rejects more low-quality explanations than PASTARej in approximately 84% of the experiments across varying rejection rates $\rho\%$ from 1% to 25%, thus confirming in superior ability in detecting low-quality explanations.

5 Related Work

Learning to Reject. The problem of deferring hard decisions has been studied in the context of *learning to reject*, *learning to defer* [61, 62], *learning under algorithmic triage* [63, 64], *learning under human assistance* [65, 66], and *learning to complement* [67, 68]; see [5] for a recent survey. These approaches all enable the machine to offload certain decisions to a human expert, but differ in what criterion they use. While some strategies entirely rely on the machine’s self-assessed uncertainty [8, 63, 69], others implement the rejection policy as a machine learning classifier and optimize it for joint team performance [61] or learn the classifier and the policy jointly [65–67]. None of them, however, considers the role of explanations in decision making, which we argue is central. Note that ULER is not meant as a replacement for existing strategies, as it has a different goal. On the contrary, it could and should be combined with them to ensure *both* incorrect predictions and unsatisfactory explanations are deferred. We will evaluate this generalization in future work.

Explainable AI (XAI) aims at designing mechanisms for properly justifying algorithmic decisions to end-users in non-technical terms [70]. We focus on (post-hoc) feature attribution techniques, which highlight what features influenced a prediction the most. Many high profile techniques belong to this group, *e.g.*, LIME [20], SHAP [71–74, 21], (integrated) input gradients [75, 18], and formal feature attributions [76]. ULER can assess the perceived quality of attributions irrespectively of how these are computed, and as such it works with any of them, see also our experiments with LIME in Appendix C.3. The only work that combines XAI and LtR is [77], which focuses on explaining the reasons behind rejection using counterfactuals, and as such is orthogonal to our work.

Evaluating explanations. There is a large body of work on evaluating explanation quality. Most metrics are “machine-side”, in that they only consider properties of the model and of how the explanation is computed (*e.g.*, faithfulness, stability, complexity) [36–38, 34, 78–80]. Our experiments show that these metrics cannot anticipate whether *users* will agree with or believe in a given explanation. In contrast, we learn our rejector to mimic human judgments of explanation quality. Closest to our work is PASTA [39], which however is not designed for rejection and underperforms in our experiments.

6 Conclusion

We have introduced the problem of *learning to reject low-quality explanations* (LtX) and proposed ULER, a simple yet effective technique for learning a high-quality rejector from a limited amount of expert feedback. Our empirical analysis showcases how, in contrast to other LtR approaches, ULER successfully identifies low-quality explanations in both synthetic and human-annotated tasks. In future work, we will extend our setup to learn the rejector and classifier jointly, so as to optimize their overall performance [65–67], and look into leveraging ULER’s rejector as a loss term for aligning machine explanations of confounded ML models for the purpose of debiasing [81].

Acknowledgments and Disclosure of Funding

This research was supported by the Flemish Government through the “Onderzoeksprogramma Artificiële Intelligentie (AI) Vlaanderen” programme [LS, JD], and by the KU Leuven Research Fund (iBOF/21/075) [JD]. This project has received funding from the European Union under Grant Agreement No. 101120763 – TANGO [ST, DP]. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Health and Digital Executive Agency (HaDEA). Neither the European Union nor the granting authority can be held responsible for them.

References

- [1] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen A. W. M. van der Laak, Bram van Ginneken, and Clara I. Sánchez. A survey on deep learning in medical image analysis. *Medical Image Anal.*, 42:60–88, 2017.
- [2] Filippo Pesapane, Marina Codari, and Francesco Sardanelli. Artificial intelligence in medical imaging: threat or opportunity? radiologists again at the forefront of innovation in medicine. *European Radiology Experimental*, 2, 10 2018. doi: 10.1186/s41747-018-0061-6.
- [3] Periklis Gogas and Theophilos Papadimitriou. Machine learning in economics and finance. *SSRN Electronic Journal*, 01 2023. doi: 10.2139/ssrn.3885538.
- [4] Constantine Kotropoulos and Gonzalo R. Arce. Linear classifier with reject option for the detection of vocal fold paralysis and vocal fold edema. *EURASIP J. Adv. Signal Process.*, 2009, 2009.
- [5] Kilian Hendrickx, Lorenzo Perini, Dries Van der Plas, Wannes Meert, and Jesse Davis. Machine learning with a reject option: A survey. *Machine Learning*, pages 1–38, 2024.
- [6] Claudio De Stefano, Carlo Sansone, and Mario Vento. To reject or not to reject: that is the question-an answer in case of neural classifiers. *IEEE Trans. Syst. Man Cybern. Part C*, 30(1): 84–94, 2000.
- [7] C. K. Chow. On optimum recognition error and reject tradeoff. *IEEE Trans. Inf. Theory*, 16(1): 41–46, 1970.
- [8] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Learning with rejection. In *Algorithmic Learning Theory*, 2016.
- [9] Weitang Liu, Xiaoyun Wang, John D. Owens, and Yixuan Li. Energy-based out-of-distribution detection. In *NeurIPS*, 2020.

- [10] Jenia Kim, Henry Maathuis, and Danielle Sent. Human-centered evaluation of explainable AI applications: a systematic review. *Frontiers Artif. Intell.*, 7, 2024.
- [11] European Parliament and Council of the European Union. Regulation (EU) 2016/679 of the European Parliament and of the Council. URL <https://data.europa.eu/eli/reg/2016/679/oj>.
- [12] Misgina Tsighe Hagos, Kathleen M. Curran, and Brian Mac Namee. Identifying spurious correlations and correcting them with an explanation-based learning. *CoRR*, abs/2211.08285, 2022.
- [13] P Jonathon Phillips, P Jonathon Phillips, Carina A Hahn, Peter C Fontana, Amy N Yates, Kristen Greene, David A Broniatowski, and Mark A Przybocki. Four principles of explainable artificial intelligence. 2021.
- [14] Leilani H. Gilpin, David Bau, Ben Z. Yuan, Ayesha Bajwa, Michael A. Specter, and Lalana Kagal. Explaining explanations: An overview of interpretability of machine learning. In *DSAA*, pages 80–89. IEEE, 2018.
- [15] Johannes Schneider, Christian Meske, and Michalis Vlachos. Deceptive xai: Typology, creation and detection. *SN Computer Science*, 5(1):81, 2023.
- [16] Himabindu Lakkaraju and Osbert Bastani. "how do I fool you?": Manipulating user trust via misleading black box explanations. In *AIES*, pages 79–85. ACM, 2020.
- [17] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [18] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR, 2017.
- [19] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)*, 51(5):1–42, 2018.
- [20] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "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*, pages 1135–1144, 2016.
- [21] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.
- [22] Alexey Ignatiev, Nina Narodytska, and João Marques-Silva. Abduction-based explanations for machine learning models. In *AAAI*, pages 1511–1519. AAAI Press, 2019.
- [23] Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognit.*, 65:211–222, 2017.
- [24] Ramaravind Kommiya Mothilal, Amit Sharma, and Chenhao Tan. Explaining machine learning classifiers through diverse counterfactual explanations. In *FAT**, pages 607–617. ACM, 2020.
- [25] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *Int. J. Comput. Vis.*, 128(2):336–359, 2020.
- [26] Vojtech Franc, Daniel Prruvs, and Václav Voráček. Optimal strategies for reject option classifiers. *J. Mach. Learn. Res.*, 24:11:1–11:49, 2023.
- [27] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Boosting with abstention. In *NIPS*, pages 1660–1668, 2016.
- [28] Lorenzo Perini and Jesse Davis. Unsupervised anomaly detection with rejection. In *NeurIPS*, 2023.

- [29] Yiyu Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*, pages 20827–20840. PMLR, 2022.
- [30] Dries Van der Plas, Wannes Meert, Johan Verbraecken, and Jesse Davis. A reject option for automated sleep stage scoring. In *Workshop on Interpretable ML in Healthcare at International Conference on Machine Learning (ICML)*, 2021.
- [31] Luca Stradiotti, Lorenzo Perini, and Jesse Davis. Combining active learning and learning to reject for anomaly detection. In *ECAI*, volume 392 of *Frontiers in Artificial Intelligence and Applications*, pages 2266–2273. IOS Press, 2024.
- [32] Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. In *NIPS*, pages 4878–4887, 2017.
- [33] Zixi Chen, Varshini Subhash, Marton Havasi, Weiwei Pan, and Finale Doshi-Velez. What makes a good explanation?: A harmonized view of properties of explanations. *arXiv preprint arXiv:2211.05667*, 2022.
- [34] Sanjoy Dasgupta, Nave Frost, and Michal Moshkovitz. Framework for evaluating faithfulness of local explanations. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*, pages 4794–4815. PMLR, 2022.
- [35] Ramaravind Kommiya Mothilal, Divyat Mahajan, Chenhao Tan, and Amit Sharma. Towards unifying feature attribution and counterfactual explanations: Different means to the same end. In *AIES*, pages 652–663. ACM, 2021.
- [36] Steve Azzolin, Antonio Longa, Stefano Teso, and Andrea Passerini. Perks and pitfalls of faithfulness in regular, self-explainable and domain invariant gnns. *arXiv preprint arXiv:2406.15156*, 2024.
- [37] Alexandros Kalousis, Julien Prados, and Melanie Hilario. Stability of feature selection algorithms: a study on high-dimensional spaces. *Knowl. Inf. Syst.*, 12(1):95–116, 2007.
- [38] Dylan Slack, Anna Hilgard, Sameer Singh, and Himabindu Lakkaraju. Reliable Post-hoc Explanations: Modeling Uncertainty in Explainability. *NeurIPS*, 2021.
- [39] Rémi Kazmierczak, Steve Azzolin, Eloïse Berthier, Anna Hedström, Patricia Delhomme, Nicolas Bousquet, Goran Frehse, Massimiliano Mancini, Baptiste Caramiaux, Andrea Passerini, et al. Benchmarking xai explanations with human-aligned evaluations. *arXiv preprint arXiv:2411.02470*, 2024.
- [40] Julien Colin, Thomas Fel, Rémi Cadène, and Thomas Serre. What I cannot predict, I do not understand: A human-centered evaluation framework for explainability methods. In *NeurIPS*, 2022.
- [41] Robert R. Hoffman, Shane T. Mueller, Gary Klein, and Jordan Litman. Metrics for explainable AI: challenges and prospects. *CoRR*, abs/1812.04608, 2018.
- [42] Janet Hui-wen Hsiao, Hilary Hei Ting Ngai, Luyu Qiu, Yi Yang, and Caleb Chen Cao. Roadmap of designing cognitive metrics for explainable artificial intelligence (XAI). *CoRR*, abs/2108.01737, 2021.
- [43] Stefano Teso and Kristian Kersting. Explanatory interactive machine learning. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 239–245, 2019.
- [44] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20:273–297, 1995.
- [45] Jianlong Zhou, Amir H. Gandomi, Fang Chen, and Andreas Holzinger. Evaluating the quality of machine learning explanations: A survey on methods and metrics. *Electronics*, 10(5), 2021. ISSN 2079-9292. doi: 10.3390/electronics10050593. URL <https://www.mdpi.com/2079-9292/10/5/593>.

- [46] Vivian Lai and Chenhao Tan. On human predictions with explanations and predictions of machine learning models: A case study on deception detection. In danah boyd and Jamie H. Morgenstern, editors, *Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT* 2019, Atlanta, GA, USA, January 29-31, 2019*, pages 29–38. ACM, 2019. doi: 10.1145/3287560.3287590. URL <https://doi.org/10.1145/3287560.3287590>.
- [47] Ahmed Zaoui, Christophe Denis, and Mohamed Hebiri. Regression with reject option and application to knn. In *NeurIPS*, 2020.
- [48] Umang Bhatt, Adrian Weller, and José M. F. Moura. Evaluating and aggregating feature-based model explanations. In *IJCAI*, pages 3016–3022. ijcai.org, 2020.
- [49] Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *ICLR*, 2018.
- [50] Markelle Kelly, Rachel Longjohn, and Kolby Nottingham. The uci machine learning repository, 2023.
- [51] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [52] Pieter Robberechts, Cem Arslan, and Jesse Davis. Enhancing xG models with freeze frame data. Available at <https://dtai.cs.kuleuven.be/sports/blog/enhancing-xg-models-with-freeze-frame-data/>, 2020tjha.
- [53] Pieter Robberechts and Jesse Davis. *How Data Availability Affects the Ability to Learn Good xG Models*, pages 17–27. 12 2020. ISBN 978-3-030-64911-1. doi: 10.1007/978-3-030-64912-8_2.
- [54] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- [55] Hudl Statsbomb. The 2015/16 Big 5 Leagues Free Data Release: La Liga. Available at <https://statsbomb.com/news/the-2015-16-big-5-leagues-free-data-release-la-liga/>, 2023.
- [56] Alina Pommeranz, Joost Broekens, Pascal Wiggers, Willem-Paul Brinkman, and Catholijn M Jonker. Designing interfaces for explicit preference elicitation: a user-centered investigation of preference representation and elicitation process. *User Modeling and User-Adapted Interaction*, 22:357–397, 2012.
- [57] Loretta J Stalans. Frames, framing effects, and survey responses. *Handbook of survey methodology for the social sciences*, pages 75–90, 2012.
- [58] Astrid Bertrand, Rafik Belloum, James R Eagan, and Winston Maxwell. How cognitive biases affect xai-assisted decision-making: A systematic review. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, pages 78–91, 2022.
- [59] Alka Joshi, Satyendra Kale, Satish Chandel, and Divya K Pal. Likert scale: Explored and explained. *British Journal of Applied Science & Technology*, 7(4):396–403, 2015. doi: 10.9734/BJAST/2015/14975.
- [60] Katherine A Batterton and Kimberly N Hale. The likert scale what it is and how to use it. *Phalanx*, 50(2):32–39, 2017.
- [61] David Madras et al. Predict Responsibly: Improving Fairness and Accuracy by Learning to Defer. *NeurIPS*, 2018.
- [62] Hussein Mozannar and David Sontag. Consistent estimators for learning to defer to an expert. In *ICML*, 2020.

- [63] Maithra Raghu, Katy Blumer, Greg Corrado, Jon Kleinberg, Ziad Obermeyer, and Sendhil Mullainathan. The algorithmic automation problem: Prediction, triage, and human effort. *arXiv:1903.12220*, 2019.
- [64] Nastaran Okati et al. Differentiable learning under triage. *NeurIPS*, 2021.
- [65] Abir De, Paramita Koley, Niloy Ganguly, and Manuel Gomez-Rodriguez. Regression under human assistance. In *AAAI*, 2020.
- [66] Abir De et al. Classification under human assistance. In *AAAI*, 2021.
- [67] Bryan Wilder et al. Learning to complement humans. In *IJCAI*, 2021.
- [68] Gagan Bansal, Besmira Nushi, Ece Kamar, Eric Horvitz, and Daniel S Weld. Is the most accurate ai the best teammate? optimizing ai for teamwork. In *AAAI*, 2021.
- [69] Jessie Liu et al. Incorporating uncertainty in learning to defer algorithms for safe computer-aided diagnosis. *Scientific Reports*, 2022.
- [70] Amina Adadi and Mohammed Berrada. Peeking inside the black-box: A survey on explainable artificial intelligence (XAI). *IEEE Access*, 6:52138–52160, 2018.
- [71] Stan Lipovetsky and Michael Conklin. Analysis of regression in game theory approach. *Applied stochastic models in business and industry*, 17(4):319–330, 2001.
- [72] Erik Štrumbelj and Igor Kononenko. An efficient explanation of individual classifications using game theory. *The Journal of Machine Learning Research*, 11:1–18, 2010.
- [73] Erik Štrumbelj and Igor Kononenko. Explaining prediction models and individual predictions with feature contributions. *Knowledge and information systems*, 41:647–665, 2014.
- [74] Anupam Datta, Shayak Sen, and Yair Zick. Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems. In *2016 IEEE symposium on security and privacy (SP)*, pages 598–617. IEEE, 2016.
- [75] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- [76] Jinqiang Yu, Alexey Ignatiev, and Peter J Stuckey. On formal feature attribution and its approximation. *arXiv preprint arXiv:2307.03380*, 2023.
- [77] André Artelt, Roel Visser, and Barbara Hammer. “I do not know! but why?” – Local Model-agnostic Example-based Explanations of Reject. *Neurocomputing*, 2023.
- [78] David Alvarez-Melis and Tommi S. Jaakkola. On the robustness of interpretability methods. *CoRR*, abs/1806.08049, 2018.
- [79] Prasad Chalasani, Jiefeng Chen, Amrita Roy Chowdhury, Xi Wu, and Somesh Jha. Concise explanations of neural networks using adversarial training. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 1383–1391. PMLR, 2020.
- [80] An-phi Nguyen and María Rodríguez Martínez. On quantitative aspects of model interpretability. *arXiv preprint arXiv:2007.07584*, 2020.
- [81] Stefano Teso, Öznur Alkan, Wolfgang Stammer, and Elizabeth Daly. Leveraging explanations in interactive machine learning: An overview. *Frontiers in Artificial Intelligence*, 2023.
- [82] Pedro Lopes, Eduardo Silva, Cristiana Braga, Tiago Oliveira, and Luís Rosado. Xai systems evaluation: a review of human and computer-centred methods. *Applied Sciences*, 12(19):9423, 2022.
- [83] Giulia Vilone and Luca Longo. Notions of explainability and evaluation approaches for explainable artificial intelligence. *Inf. Fusion*, 76:89–106, 2021.

- [84] Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Sai Suggala, David I. Inouye, and Pradeep Ravikumar. On the (in)fidelity and sensitivity of explanations. In *NeurIPS*, pages 10965–10976, 2019.
- [85] Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3681–3688, 2019.
- [86] Sarah Nogueira and Gavin Brown. Measuring the stability of feature selection. In *ECML/PKDD (2)*, volume 9852 of *Lecture Notes in Computer Science*, pages 442–457. Springer, 2016.
- [87] Saumitra Mishra, Sanghamitra Dutta, Jason Long, and Daniele Magazzeni. A survey on the robustness of feature importance and counterfactual explanations. *CoRR*, abs/2111.00358, 2021.
- [88] Sebastian Bordt, Michèle Finck, Eric Raidl, and Ulrike von Luxburg. Post-hoc explanations fail to achieve their purpose in adversarial contexts. In *FAccT*, pages 891–905. ACM, 2022.
- [89] Naman Bansal, Chirag Agarwal, and Anh Nguyen. SAM: the sensitivity of attribution methods to hyperparameters. In *CVPR*, pages 8670–8680. Computer Vision Foundation / IEEE, 2020.
- [90] David Alvarez-Melis and Tommi S Jaakkola. Towards robust interpretability with self-explaining neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 7786–7795, 2018.
- [91] Laura Rieger, Chandan Singh, William Murdoch, and Bin Yu. Interpretations are useful: penalizing explanations to align neural networks with prior knowledge. In *International Conference on Machine Learning*, pages 8116–8126. PMLR, 2020.
- [92] Cristian Bucila, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *KDD*, pages 535–541. ACM, 2006.
- [93] Nelson Cowan. The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *Behavioral and brain sciences*, 24(1):87–114, 2001.
- [94] Christoph Molnar. *Interpretable machine learning*. Lulu. com, 2020.
- [95] Sidra Naveed, Gunnar Stevens, and Dean Robin-Kern. An overview of the empirical evaluation of explainable ai (xai): A comprehensive guideline for user-centered evaluation in xai. *Applied Sciences*, 14(23):11288, 2024.
- [96] Jürgen Dieber and Sabrina Kirrane. A novel model usability evaluation framework (muse) for explainable artificial intelligence. *Inf. Fusion*, 81:143–153, 2022.
- [97] Hamed Taherdoost. What is the best response scale for survey and questionnaire design; review of different lengths of rating scale/attitude scale/likert scale. *International Journal of Academic Research in Management (IJARM)*, 8, 2019.

A Broader Impact

Rejecting low-quality explanations can be beneficial from at least two perspectives. First, when human involvement is expensive and time-consuming, this reject option serves as an effective mechanism to filter outputs based on human-validated reasoning. Second, since modern decision-making often relies on both predictions and their corresponding explanations, explanation quality becomes critical to prevent harmful decisions.

Our approach contributes to this goal by enhancing trust in the system and supporting human-validated decision-making, ultimately promoting more effective human-AI interaction. Our findings represent an initial step in this direction, showing that our method can reject more low-quality explanations than several existing and adapted learning-to-reject strategies.

B Explanation quality metrics

Explanation quality metrics aim to assess to what extent explanations satisfy the general goal of explaining a decision. These metrics can be broadly categorized into two families [82, 45, 83]: *machine-side* and *human-side* metrics. The former focus exclusively on the relationship between the explainer and the predictor, whereas the latter involve human subjects in evaluating the quality of the explanations.

B.1 Machine-side metrics.

The simplest way to evaluate an explanation is by verifying whether it effectively reveals the predictor’s underlying reasoning. Several metrics have been proposed to assess the relationship between explanations and the predictor. [33] categorize existing machine-side metrics - and provide their mathematical formulations — into three groups: stability, faithfulness, and complexity. We exclude homogeneity from our analysis because it is defined for groups of explanations rather than individual ones.

Stability measures the similarity of explanations under changes to the input instance, the training data or the model hyperparameters [84, 78, 85, 37, 86, 87]. This can be harmful because an attacker can selectively choose explanations based on their (potentially adversarial) interests [15, 88]. Following [89], we define the stability of an explanation as the average similarity across multiple runs of the same explainer, each potentially yielding a different explanation. Formally, given an instance \mathbf{x} and prediction $f(\mathbf{x})$ with associated explanation \mathbf{z} , *stability* is defined as:

$$\text{stab}(\mathbf{z}) = \mathbb{E}_{\mathbf{z}' \sim \mathcal{Z}} [\text{Sim}(\mathbf{z}, \mathbf{z}')] \quad (3)$$

where *Sim* is a similarity metric and \mathcal{Z} denotes the space of possible explanations for the given prediction. In practice, we compute stability using the Pearson correlation coefficient as the similarity metric and average it across ten independently generated explanations.

Faithfulness measures how accurately an explanation captures the true underlying behavior of the predictor [48, 90, 91, 80, 34, 39]. Given an explanation \mathbf{z} , we define the sets of relevant features $\mathbf{z}_{\mathcal{R}} = \{i < d : |\mathbf{z}_i| > 0\}$ and irrelevant features $\mathbf{z}_{\mathcal{I}} = \{i < d : |\mathbf{z}_i| = 0\}$. Intuitively, an explanation is faithful if perturbing irrelevant features causes little to no change in the predictor’s output, while perturbing relevant features induces significant changes. Building on [36], we define *faithfulness* (faith) as the harmonic mean of *sufficiency* (suf) and *necessity* (nec), which estimate the sensitivity of the prediction to perturbations in irrelevant and relevant features, respectively. Formally, given a instance-prediction pair $(\mathbf{x}, f(\mathbf{x}))$ with associated explanation \mathbf{z} , and the predictor to be explained f , *sufficiency* and *necessity* are defined as:

$$\text{suf}_{d, p_{\mathcal{I}}}(\mathbf{z}) = \mathbb{E}_{\mathbf{x}' \sim p_{\mathcal{I}}} [\Delta_f(\mathbf{x}, \mathbf{x}')], \quad \text{nec}_{d, p_{\mathcal{R}}}(\mathbf{z}) = \mathbb{E}_{\mathbf{x}' \sim p_{\mathcal{R}}} [\Delta_f(\mathbf{x}, \mathbf{x}')] \quad (4)$$

where Δ_f measures prediction change between \mathbf{x} and its perturbed version \mathbf{x}' , and $p_{\mathcal{R}}$ and $p_{\mathcal{I}}$ are interventional distributions that specify how to perturb relevant and irrelevant features, respectively. Equations 4 are then normalized to $[0, 1]$ range, the higher the better, via a non-linear transformation i.e., respectively $\exp(-\text{suf}_{d, p_{\mathcal{I}}})$ and $1 - \exp(-\text{nec}_{d, p_{\mathcal{R}}})$. Operationally, for a given instance-explanation pair (\mathbf{x}, \mathbf{z}) sampling from $p_{\mathcal{R}}$ ($p_{\mathcal{I}}$) involves perturbing the features in $\mathbf{z}_{\mathcal{R}}$ ($\mathbf{z}_{\mathcal{I}}$) following [92], while keeping the remaining features fixed. Additionally, the prediction change Δ_f is computed

either as the absolute difference in positive class probability for classification tasks, *i.e.*, $|P(Y = 1|\mathbf{x}) - P(Y = 1|\mathbf{x}')|$, or the absolute prediction difference in regression, *i.e.*, $|f(\mathbf{x}) - f(\mathbf{x}')|$.

Complexity refers to the cognitive burden associated with parsing an explanation [48, 79, 80]. In general, a less complex explanation is easier for a human to understand, making complexity a common proxy for understandability [93, 94]. Following [48], given an instance \mathbf{x} with prediction $f(\mathbf{x})$ and explanation \mathbf{z} , we formally define *complexity* as:

$$\text{compl}_{d, p_{\mathcal{X}}} = \mathbb{E}[-\ln(\bar{\mathbf{z}})] = -\sum_{i=1}^d \bar{z}_i \ln(\bar{z}_i) \quad (5)$$

where \bar{z}_i is the fractional contribution of feature i , *i.e.*, the ratio of its absolute relevance score $|z_i|$ to the sum of all the absolute relevance scores $\sum_{j=1}^d |z_j|$.

B.2 Human-side metrics

Despite the literature recognizing the importance of human-centered evaluations [39, 83], only a few metrics have been proposed to evaluate explanations from perspective of a human [95]. This gap stems from the inherently subjective nature of human evaluations, which typically makes it challenging to provide a precise mathematical formulation for a metric [33]. Moreover, there is no consensus in the literature regarding standard criteria for human-side evaluation metrics [45].

PASTA uses a model to score each explanation based on how this is perceived by humans [39]. The authors first construct a dataset in which users rated several explanations according to four key desiderata: faithfulness, robustness, complexity, and objectivity. Then, the *PASTA-metric* is trained on these ratings to derive a metric value for new explanations. Specifically, this model consists of two main components: an embedding network that leverages a foundation model to generate feature embeddings from the explanations, and a scoring network that employs a linear layer to predict the human ratings based on these embeddings. PASTA is the closest competitor to our work in that it also aims to assess explanations based on human feedback. However, there are three substantial differences with our approach. First, PASTA is designed for image data and relies on an embedding network to create embeddings from this high-dimensional space, whereas we focus on tabular data and learn directly from feature-importance explanations. Second, PASTA does not include a rejection mechanism and always returns a score regardless of quality, while we explicitly aim to develop a reject option based on explanation quality. Third, PASTA seeks to create a dataset-agnostic metric and thus annotates 25 explanations per dataset to encourage generalization. In contrast, we aim to train a dataset-specific rejector and therefore collect 1050 annotations for a single dataset.

Other human-side metrics. *Understandability* measures whether an explanation is easy to comprehend for the human [82]. The rationale behind this metric is to examine whether the explanations facilitate the user’s understanding of the model’s decisions [96]. *Plausibility* is high if \mathbf{z} matches the ground-truth explanation \mathbf{z}^* , assuming the latter exists and is unique. Depending on the model’s behavior and structure of the underlying learning problem, the model’s reasoning may or may not reflect the ground-truth explanation \mathbf{z}^* . Our approach implicitly addresses both metrics. The user’s rating depends on how understandable the explanation is, *i.e.*, users tend to assign low scores to explanations they find difficult to interpret. Furthermore, the per-feature feedback we collect encourages users to identify features that substantially deviate from their expectations, thereby aligning the underlying ground truth.

C Experiments: extended details and results

C.1 Dataset characteristics and predictors’ performance

Table 3 presents the characteristics of the eight datasets used in the empirical evaluation, along with the performance of the oracle \mathcal{O} and the predictor f . We report the balanced accuracy (*BACC*) for classification tasks; for regression tasks, we report the mean squared error (*MSE*) after normalizing the target variable to the $[0, 1]$ -range. Specifically, both predictors are trained on a training set \mathcal{T} and evaluated on a test set \mathcal{D} . The size of \mathcal{D} is limited because obtaining human-judgment labels on explanation quality is expensive [39]. Additionally, the table reports the proportion of low-quality explanations γ in \mathcal{D} for each dataset, as determined using the procedure described in Section 4.1.

Table 3: **Datasets’ characteristics and predictors’ performance.** This table reports the datasets’ characteristics (*i.e.*, size of the training set, number of features, size of the test set, proportion of low-quality explanations) and the oracle \mathcal{O} and f ’s performance on the eight benchmark datasets used in the experiments.

dataset	$\#(\mathcal{T})$	d	$\#(\mathcal{D})$	$BACC_{\mathcal{O}} \uparrow$	$BACC_f \uparrow$	γ
compas	10000	12	2000	0.770	0.690	0.05
creditcard	10000	23	2000	0.660	0.608	0.12
adult	10000	12	2000	0.756	0.757	0.02
churn	1000	13	1850	0.838	0.696	0.15
dataset	$\#(\mathcal{T})$	d	$\#(\mathcal{D})$	$MSE_{\mathcal{O}} \downarrow$	$MSE_f \downarrow$	γ
news	10000	58	2000	0.004	0.009	0.48
wine	1000	11	2000	0.014	0.015	0.02
parkinson	1000	19	2000	0.008	0.044	0.46
appliances	10000	13	2000	0.005	0.010	0.32

C.2 Hyperparameter Selection

We optimize all hyperparameters using a grid search on the validation split \mathcal{D}_{val} . Specifically, for ULER we optimize the SVM kernel (linear, polynomial, RBF), the cost of mistakes $C \in \{0.1, 1, 10\}$, the number of augmentations per explanation $k \in \{5, 10, 20\}$ and the noise $\epsilon_0 \in \{0.1, 0.5, 1\}$. For PASTA, we employ the authors’ code for the scoring network and optimize the loss hyperparameters $\alpha \in \{0.1, 1, 10\}$, $\beta \in \{0.001, 0.01, 0.1\}$ and $\gamma \in \{0.01, 0.1, 1\}$. For NovRej_X and NovRej_Z , we optimize the number of neighbors $k_{NN} \in \{1, 5, 10\}$.

C.3 Robustness to the choice of the explainer

In this section, we assess the robustness of our approach to the choice of explanation method. Specifically, we replicate the experimental setup from Section 4.1, but generate all explanations using LIME [20] with its default hyperparameters.

Fig. 4 shows the percentage of low-quality explanations for the accepted and the rejected set as a function of the rejection rate $\rho\%$ averaged across the eight datasets considered. Even when using LIME, ULER outperforms the competitors across all rejection rates. Interestingly, all explanation-aware rejectors perform better in this setting compared to when explanations are generated using KernelSHAP. Upon inspection, we observed that LIME tends to produce more correlated explanations for both predictors, making it easier to identify cases where explanations differ. On average, across all datasets and rejection rates, ULER reduces the percentage of low-quality explanations in the accepted set by 18% compared to the best competitor PASTARej, which is very similar to the gain when KernelSHAP is used as explainer.

Finally, Table 4 reports the average AUROC per dataset. Again, ULER achieves the highest AUROC on all datasets, demonstrating superior ability to distinguish low- from high-quality explanations. While NovRej_Z ties with ULER on *compas* and *FaithRej* does so on *churn*, only ULER consistently performs well across all datasets, whereas the baselines exhibit higher variance in performance.

C.4 Robustness to the choice of the oracle

We further assess the robustness of our approach by employing an alternative predictor for the oracle \mathcal{O} . In particular, we replicate the experimental setup from Section 4.1, but in this setting, the oracle \mathcal{O} uses a SVM with an RBF kernel. This allows us to evaluate our method in a scenario where the oracle and the predictor have a similar inductive bias. The average balanced accuracy on the classification datasets is 0.67, while the average MSE on the regression datasets is 0.02.

Fig. 5 shows the percentage of low-quality explanations for the accepted and the rejected set as a function of the rejection rate $\rho\%$ averaged across the eight datasets considered. Even when using a different predictor for \mathcal{O} , ULER outperforms the competitors across all rejection rates. On average, across all datasets and rejection rates, ULER reduces the percentage of low-quality explanations in

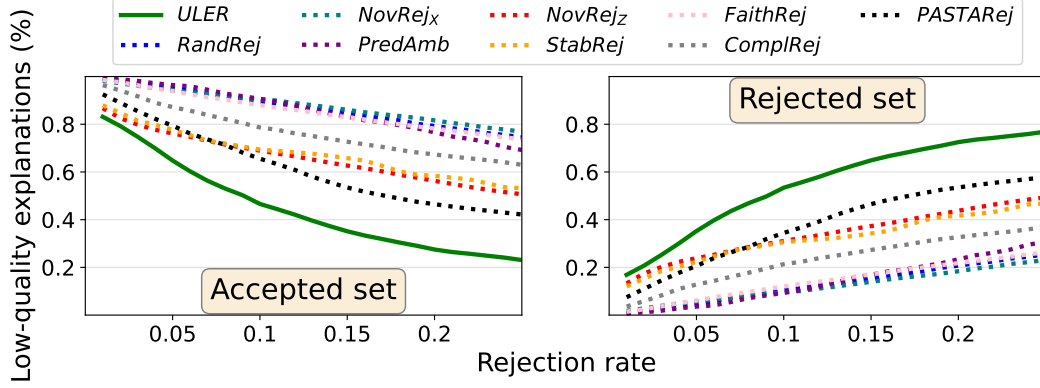


Figure 4: **ULER rejects on average more low-quality explanations than all competitors when LIME is used as explainer.** Average percentage of low quality explanations in the accepted and rejected set for all the considered strategies over the 8 datasets for 25 rejection rates $\rho\%$. ULER outperforms all the competitors for all rejection rates considered, demonstrating its robustness to the choice of the explainer.

Table 4: **ULER outperforms the competitors at separating low-quality from high-quality explanations LIME is used as explainer.** Average AUROC for all the rejection strategies over the 8 datasets and its standard deviation. ULER consistently obtains the best results in all datasets, demonstrating its robustness to the choice of the explainer

	compas	Classification			churn	news	Regression			appliances
		creditcard	adult				wine	parkinson		
ULER	1.00 \pm 0.00	0.81 \pm 0.05	0.90 \pm 0.02	0.74 \pm 0.02	1.00 \pm 0.00	0.57 \pm 0.03	0.99 \pm 0.00	0.82 \pm 0.02		
RandRej	0.52 \pm 0.05	0.50 \pm 0.02	0.53 \pm 0.06	0.49 \pm 0.02	0.50 \pm 0.02	0.51 \pm 0.07	0.49 \pm 0.01	0.50 \pm 0.02		
NovRej _x	0.46 \pm 0.04	0.58 \pm 0.02	0.30 \pm 0.05	0.36 \pm 0.02	0.54 \pm 0.01	0.51 \pm 0.04	0.58 \pm 0.02	0.56 \pm 0.02		
PredAmb	0.56 \pm 0.03	0.46 \pm 0.02	0.71 \pm 0.03	0.85 \pm 0.01	0.50 \pm 0.00	0.50 \pm 0.00	0.50 \pm 0.00	0.50 \pm 0.00		
StabRej	0.91 \pm 0.04	0.76 \pm 0.04	0.79 \pm 0.05	0.53 \pm 0.03	0.63 \pm 0.01	0.46 \pm 0.03	0.63 \pm 0.01	0.53 \pm 0.02		
FaithRej	0.13 \pm 0.05	0.48 \pm 0.02	0.29 \pm 0.03	0.74 \pm 0.03	0.66 \pm 0.01	0.42 \pm 0.03	0.55 \pm 0.02	0.56 \pm 0.03		
ComplRej	0.00 \pm 0.00	0.75 \pm 0.05	0.74 \pm 0.02	0.52 \pm 0.02	0.91 \pm 0.01	0.56 \pm 0.03	0.34 \pm 0.01	0.53 \pm 0.02		
PASTARej	0.30 \pm 0.48	0.78 \pm 0.06	0.65 \pm 0.08	0.62 \pm 0.03	0.99 \pm 0.01	0.49 \pm 0.03	0.96 \pm 0.03	0.79 \pm 0.04		
NovRej _z	1.00 \pm 0.00	0.69 \pm 0.06	0.78 \pm 0.02	0.46 \pm 0.02	0.81 \pm 0.01	0.51 \pm 0.02	0.70 \pm 0.01	0.55 \pm 0.02		

the accepted set by 26% compared to the best competitor PASTARej, which is even higher than the improvement obtained using Random Forest for the oracle \mathcal{O} (see Fig. 2).

Finally, Table 5 reports the average AUROC, measuring each rejector’s ability to distinguish low- from high-quality explanations across datasets. ULER achieves the highest AUROC on all datasets and consistently performs well across datasets, while the baselines show greater variance in performance.

C.5 Ablation study - Training the rejector without augmenting the data

In this section, we evaluate whether the augmentation effectively improves the performance of the rejector. To this end, we compare ULER with an ablated variant, ULER-NOAUG, which does not leverage the feedback-aware augmentation strategy. Specifically, ULER-NOAUG trains the rejector as described in Section 3.1, but uses \mathcal{D} instead of the augmented data \mathcal{D}_{aug} .

Table 6 reports the average AUROC per dataset for both ULER and ULER-NOAUG, assessing their performance in distinguishing low-quality from high-quality explanations. For comparison, we also report the best-performing baseline from Q1 for each dataset, denoted as *best*. ULER consistently outperforms its ablated variant across all considered datasets. While the performance gain in performance is quite small ($\approx 2\%$), it is consistent: ULER always outperforms the variant without augmentation across all datasets. We argue that this improvement is still worth it given the minimal additional cost to obtain the feature-level feedback. Once user-provided quality judgments are collected, obtaining per-feature feedback is inexpensive because users are already focused on identifying features with

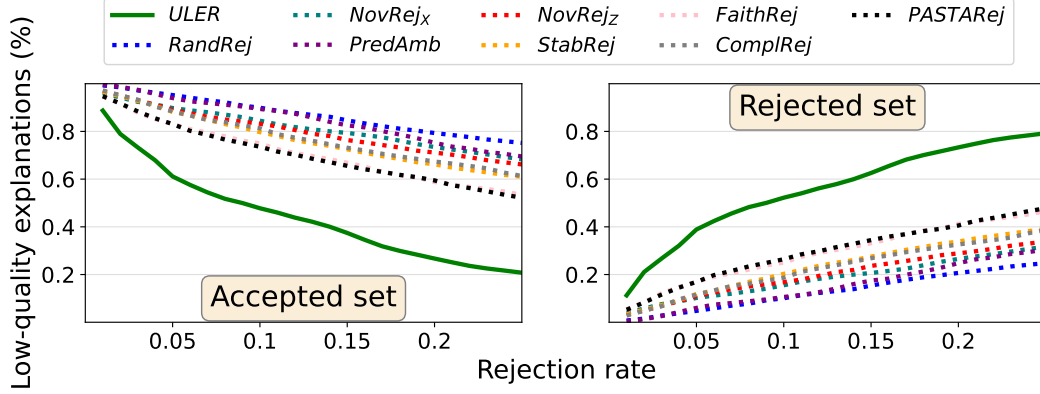


Figure 5: ULER rejects on average more low-quality explanations than all competitors when \mathcal{O} and f have a similar inductive bias. Average percentage of low quality explanations in the accepted and rejected set for all the considered strategies over the 8 datasets for 25 rejection rates $\rho\%$. ULER outperforms all the competitors for all rejection rates considered, demonstrating its robustness to the choice of the oracle predictor.

Table 5: ULER outperforms the competitors at separating low-quality from high-quality explanations when \mathcal{O} and f have a similar inductive bias. Average AUROC for all the rejection strategies over the 8 datasets and its standard deviation. ULER consistently obtains the best results in all datasets, demonstrating its robustness to the choice of the oracle

	compas	Classification			churn	news	Regression		
		creditcard	adult				wine	parkinson	appliances
ULER	0.98 \pm 0.02	0.99 \pm 0.00	0.88 \pm 0.02	0.95 \pm 0.01	0.78 \pm 0.01	0.81 \pm 0.05	0.96 \pm 0.02	0.94 \pm 0.01	
RandRej	0.50 \pm 0.11	0.50 \pm 0.01	0.50 \pm 0.03	0.50 \pm 0.03	0.50 \pm 0.02	0.50 \pm 0.07	0.50 \pm 0.04	0.48 \pm 0.03	
PredAmb	0.63 \pm 0.06	0.57 \pm 0.02	0.56 \pm 0.03	0.49 \pm 0.03	0.50 \pm 0.02	0.53 \pm 0.10	0.50 \pm 0.04	0.49 \pm 0.04	
NovRej _X	0.54 \pm 0.08	0.29 \pm 0.03	0.59 \pm 0.03	0.73 \pm 0.02	0.46 \pm 0.02	0.71 \pm 0.03	0.63 \pm 0.03	0.54 \pm 0.02	
StabRej	0.60 \pm 0.05	0.99 \pm 0.00	0.64 \pm 0.03	0.62 \pm 0.03	0.65 \pm 0.01	0.60 \pm 0.04	0.53 \pm 0.04	0.67 \pm 0.02	
FaithRej	0.85 \pm 0.06	0.61 \pm 0.01	0.65 \pm 0.03	0.48 \pm 0.03	0.72 \pm 0.01	0.61 \pm 0.05	0.66 \pm 0.03	0.85 \pm 0.02	
ComplRej	0.48 \pm 0.05	0.90 \pm 0.01	0.55 \pm 0.02	0.83 \pm 0.01	0.63 \pm 0.01	0.54 \pm 0.05	0.68 \pm 0.02	0.65 \pm 0.01	
PASTARej	0.69 \pm 0.14	0.82 \pm 0.06	0.70 \pm 0.04	0.74 \pm 0.04	0.70 \pm 0.02	0.62 \pm 0.07	0.75 \pm 0.07	0.59 \pm 0.03	
NovRej _Z	0.63 \pm 0.05	0.21 \pm 0.02	0.74 \pm 0.02	0.77 \pm 0.02	0.43 \pm 0.03	0.73 \pm 0.03	0.51 \pm 0.03	0.52 \pm 0.03	

Table 6: ULER shows a small but consistent improvement over its variant without augmentation in separating low-quality from high-quality explanations. Average AUROC for ULER and ULER-NOAUG across the eight datasets. For comparison, we also report the *best* performing baselines. ULER consistently achieves a modest but consistent improvement in AUROC across all datasets, while ULER-NOAUG still outperforms the *best* baseline.

	compas	Classification			churn	news	Regression		
		creditcard	adult				wine	parkinson	appliances
ULER	0.75 \pm 0.04	0.87 \pm 0.02	0.85 \pm 0.04	0.92 \pm 0.01	0.91 \pm 0.01	0.93 \pm 0.03	0.87 \pm 0.01	0.82 \pm 0.01	
ULER-NOAUG	0.72 \pm 0.03	0.85 \pm 0.02	0.83 \pm 0.04	0.91 \pm 0.01	0.90 \pm 0.02	0.92 \pm 0.03	0.86 \pm 0.01	0.80 \pm 0.01	
<i>best</i>	0.69 \pm 0.04	0.82 \pm 0.03	0.71 \pm 0.03	0.87 \pm 0.02	0.76 \pm 0.04	0.76 \pm 0.04	0.61 \pm 0.03	0.61 \pm 0.03	

wrong scores to assess explanation quality. In cases where per-feature feedback is not available, one could skip the augmentation step and simply use ULER-NOAUG, which still consistently outperforms the *best* baseline across all datasets and achieves an average AUROC that is 12% higher than its nearest competitors.

Table 7: ULER **predicts the human-judgments better than all competitors**. Average AUROC and its standard deviation for all the rejection strategies on the user study data.

rejector	AUROC (\pm std)
ULER	0.63 \pm 0.05
RandRej	0.55 \pm 0.07
PredAmb	0.46 \pm 0.06
NovRej _X	0.44 \pm 0.07
StabRej	0.55 \pm 0.11
FaithRej	0.55 \pm 0.08
ComplRej	0.45 \pm 0.09
NovRej _Z	0.45 \pm 0.04

C.6 Q3: Comparison with the other competitors

Additionally, we replicate the same experiments described in Section 4.2 including all competitors in Section 4 to further validate that standard LtR strategies and machine-side metrics cannot reliably reflect user judgments.

Table 7 reports the average AUROC for ULER and the other seven competitors (results for PASTARej are reported in the main paper), measuring their ability to distinguish between high-quality and low-quality explanations. ULER outperforms all competitors, achieving at least an 8% improvement in AUROC and demonstrating more consistent performance, as indicated by the lower standard deviation. We observe that StabRej and FaithRej perform similarly to the random rejector, while ComplRej performs even worse.

Additionally, we found that human annotators identified, on average, 1.8 features with incorrect relevance scores in low-quality explanations, compared to only 0.7 features in high-quality ones. This supports our intuition that low-quality explanations are perceived by users as containing more wrong relevance scores.

D User study

D.1 Data

For this user study, we used the publicly available StatsBomb 360 event stream data [55]. This contextualized event stream data is extracted from broadcast video and contains event stream data, and snapshots of player positioning at the moment of each event. The event stream data describes semantic information about the on-the-ball actions, such as which actions are performed, their start and end location, the outcome of the action, which players performed them, and the time in the match they were performed at.

D.2 Obtaining the explanations

To obtain the explanations, we begin by preprocessing the data [55] to obtain the features needed to train the classifier. From each shot snapshot, we extract the following features: (i) the distance from the ball to the center of the goal, (ii) the angle between the ball and the goalposts, (iii) the distance of the goalkeeper from the goal line, (iv) the distance of the goalkeeper from the midline (*i.e.*, the line that passes through the center of the field and the middle of the goals), and (v) the distance to the closest defender (excluding the goalkeeper). We select only these features for two main reasons: they are easily interpretable from the snapshot (see Fig. 3), and their meanings are non-overlapping, which makes it easier for annotators to disentangle their individual contributions as we found empirically that working with strongly correlated features can complicate human assessment. Using these features, we train an XGBoost ensemble [54] consisting of 50 trees with a maximum depth of 3, as it is standard practice in soccer analytics [52]. The model is trained on shots from the 2015–2016 season across four major top-tier leagues (Germany, Spain, England, and France). We evaluate the classifier on a separated test set of 1050 shots from the Italian top division during the same season, where it achieves an AUROC of 0.81 and a Brier score of 0.067.

We then use the test set to generate the explanations. As for the benchmark datasets, explanations are generated using KernelSHAP [21] with 100 samples and the training set used as background.

D.3 Human annotation process


Participants were recruited using Prolific, a crowd sourcing platform. We applied Prolific’s filters to ensure that participants possessed sufficient soccer expertise. Specifically, we applied filters to recruit subjects that (i) live in countries where soccer is widespread (UK, Germany, France, Spain, Belgium, Italy, Netherlands, or Portugal), and (ii) actively watch and play soccer. All participants were compensated with £3 for an expected completion time of 25 minutes, as estimated from the pilot studies.

After conducting pilot studies to ensure that the task was clear and comprehensible and to verify intra-annotator consistency, we launched the main user study. Participants were first requested to give their consent to participate. Then, they were provided with a link to an external Google Doc containing task instructions, which they could consult at any time during the session. The document provides general introduction for the task setting and objective, the description and illustration of the predictor’s features, and 3 exemplary snapshots. After the task introduction, participants completed three warm-up trials to familiarize themselves with the interface and the task; this was followed by the real annotation session comprising of 30 trials. In each trial, participants were asked three questions: two 5-point Likert-scale questions to separately assess the quality of the prediction and explanation, and one multiple choice question to identify the features with a wrong relevance score. We used two separate questions, presented in distinct sections of the form, to disentangle participants’ agreement with the prediction from their perception of the explanation’s quality and to minimize spurious correlations between their responses. 5-point Likert scales have been chosen as they provide satisfactory reliability and validity [97]. Specifically, in the first question, participants were shown an image containing only the shot snapshot along with the predicted probability of scoring (see Fig. 6) and asked to assess their agreement with the prediction - *"The AI thinks that the probability that the shooter will score is 1%, which is much lower than the average (10%). To what extent do you agree with the AI's prediction?"*, where the comparison *much lower* was dynamically adapted based on the predicted probability. For the second and third questions, participants were shown a different image containing the shot snapshot, the prediction, and the explanation (see Fig. 7). To facilitate interpretation, features relevance are visualized as independent arrows: blue indicates a positive impact on the prediction, while red indicates a negative impact. The second question - *"To what extent is the AI's explanation consistent with how you would explain the predicted probability of scoring?"* - was used to collect the perceived explanation quality. While the third question - *"Which features are being used incorrectly, if any?"* - is used to obtain the feature-level feedback about the features with an incorrect relevance score in the prediction. To ensure high-quality annotations, we included an attention check requiring specific answers for a trial. This allowed us to detect and discard inattentive or randomly answering participants.

D.4 Annotations preprocessing

To ensure high-quality annotations, we applied several filtering steps. First, we excluded participants who failed more than one attention check question, as well as those who consistently provided the same score for every explanation (typically a score of 3), since this means they were not able (or did not bother) to discriminate between explanations. We also removed two participants who did not flag any relevance score as incorrect. Additionally, given the subjective nature of the task (for instance, we saw that showed very low annotator agreement, *e.g.*, 1 vs 5) we removed explanations for which the standard deviation of the explanation quality scores exceeded 1.25. This step helped ensuring that our dataset contains only explanations where annotators’ opinions are reasonably consistent. After applying these filters, 718 explanations remained for our experiments.

The AI thinks that the probability that the shooter will score is 1%, which is much lower than the average (10%). To what extent do you agree with the AI's prediction? *



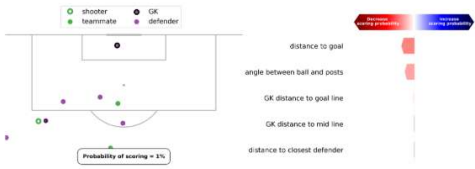
Probability of scoring = 1%

1 2 3 4 5

Completely disagree ☐ ☐ ☐ ☐ ☐ Completely agree

Figure 6: Example of the first image of each trial

To what extent is the AI's explanation consistent with how you would explain the predicted probability of scoring? *



Probability of scoring = 1%

1 2 3 4 5

Not consistent at all ☐ ☐ ☐ ☐ ☐ Fully consistent

Which features are being used incorrectly, if any? *

☐ distance to goal

☐ angle between ball and posts

☐ GK distance to goal line

☐ GK distance to mid line

☐ distance to closest defender

☐ None

Figure 7: Example of the second image of each trial

