# The Disagreement Problem in Explainable Machine Learning: A Practitioner's Perspective

Satyapriya Krishna *[1], Tessa Han *[1], Alex Gu[2], Javin Pombra[1],
Shahin Jabbari[3], Steven Wu[4], and Himabindu Lakkaraju[1]

[1]Harvard University
[3]Drexel University
[4]Carnegie Mellon University
[2]Massachusetts Institute of Technology

February 4, 2022

## Abstract

As various post hoc explanation methods are increasingly being leveraged to explain complex models in high-stakes settings, it becomes critical to develop a deeper understanding of if and when the explanations output by these methods disagree with each other, and how such disagreements are resolved in practice. However, there is little to no research that provides answers to these critical questions. In this work, we introduce and study the *disagreement problem* in explainable machine learning. More specifically, we formalize the notion of disagreement between explanations, analyze how often such disagreements occur in practice, and how do practitioners resolve these disagreements. To this end, we first conduct interviews with data scientists to understand what constitutes disagreement between explanations generated by different methods for the same model prediction, and introduce a novel quantitative framework to formalize this understanding. We then leverage this framework to carry out a rigorous empirical analysis with four real-world datasets, six state-of-the-art post hoc explanation methods, and eight different predictive models, to measure the extent of disagreement between the explanations generated by various popular explanation methods.

In addition, we carry out an online user study with data scientists to understand how they resolve the aforementioned disagreements. Our results indicate that state-of-the-art explanation methods often disagree in terms of the explanations they output. Worse yet, there do not seem to be any principled, well-established approaches that machine learning practitioners employ to resolve these disagreements, which in turn implies that they may be relying on misleading explanations to make critical decisions such as which models to deploy in the real world. To the best of our knowledge, this work is the first to highlight and study the *disagreement problem* in explainable machine learning. Our findings underscore the importance of developing principled evaluation metrics that enable practitioners to effectively compare explanations.

## 1 Introduction

As machine learning (ML) models are increasingly being deployed to make consequential decisions in domains such as healthcare, finance, and policy, there is a growing emphasis on ensuring that these models are readily interpretable to ML practitioners and other domain experts (e.g., doctors, policy makers). In order to assess when and how much to rely on these models, and detect systematic errors and potential biases in them,

---

*These authors contributed equally to this work

1

practitioners often seek to understand the behavior of these models [17]. However, the increasing complexity as well as the proprietary nature of predictive models make it challenging to understand these complex black boxes, and thus motivate the need for tools and techniques that can explain them in a faithful and human interpretable manner. To this end, several techniques have been proposed in recent literature to explain complex models in a *post hoc* fashion [36, 40, 45, 47, 50, 51]. Most of the popular *post hoc explanation methods* focus on explaining individual predictions (i.e., local explanations) of any given model, and can be broadly categorized into *perturbation based* (e.g., LIME, SHAP [36, 42]) and *gradient based* (e.g., Gradient times Input, SmoothGrad, Integrated Gradients, GradCAM [45, 47, 50, 51]) methods.

Owing to their generality, post hoc explanation methods are increasingly being utilized to explain a number of complex models in high stakes domains such as medicine, finance, law, and science [18, 24, 52]. Therefore, it becomes critical to ensure that the explanations generated by these methods are reliable. To this end, prior works [33, 37, 49, 54] proposed various evaluation metrics to quantify how *faithfully* or *accurately* a given explanation mimics the behavior of the underlying model. However, one of the biggest drawbacks of these metrics is that they are not general enough to be applicable to all model classes and real world settings. For example, Liu et al. [33] evaluate fidelity (faithfulness) of post hoc explanations by comparing them with the ground truth (e.g., true feature importances) of the underlying model. However, such ground truth is typically unavailable in most real world applications where post hoc explanations are employed to understand complex black boxes [54]. Hooker et al. [23] proposed "Remove and Retrain" (ROAR), which measures the fidelity of an explanation by retraining the underlying model with and without the features deemed as most important by the explanation. However, post hoc explanations are often employed in settings where there is no access to the underlying model. Prior work also leveraged some of the aforementioned metrics to analyze the behavior of post hoc explanation methods and their vulnerabilities – e.g., Ghorbani et al. [19] and Slack et al. [48] demonstrated that methods such as LIME and SHAP may result in explanations that are not only inconsistent and unstable, but also prone to adversarial attacks and fair washing [8].

While prior research has already taken the first steps towards analyzing the behavior of explanation methods, several critical aspects pertaining to these methods still remain unexplored. For instance, data scientists and ML practitioners do not typically rely on a single explanation method, but instead employ multiple such methods simultaneously to understand the rationale behind individual model predictions [25]. While ML practitioners can obtain a coherent understanding of model behavior if multiple methods generate consistent explanations, this may not always be the case. There may be instances for which explanations generated by various methods disagree with each other – e.g., the top-$k$ most important features output by different methods may differ. When faced with such a *disagreement problem*, practitioners will need to decide which explanation to rely on. The extent to which this disagreement problem occurs in practice is unclear because there is little to no research on understanding how often explanations produced by state-of-the-art methods disagree with each other. Furthermore, if and when the disagreement problem occurs, practitioners need to tackle it carefully as they may end up relying on misleading explanations otherwise, which may in turn lead to catastrophic consequences – e.g., trusting and deploying racially biased models, trusting incorrect model predictions and recommending sub-optimal treatments to patients, etc. [48]. However, the lack of reliable, general purpose evaluation metrics (as discussed in the previous paragraph) which can help ascertain and compare the quality of explanations may pose a serious challenge to addressing the disagreement problem in practice. Given all the above, it is critical to not only understand and quantify how often explanations output by state-of-the-art methods disagree with each other, but also study how such disagreements are currently being resolved by ML practitioners. However, there is no existing work that focuses on these important aspects.

We address the aforementioned gaps by introducing and studying the *disagreement problem* in explainable ML. To the best of our knowledge, this work is the first to highlight the disagreement problem, determine the extent to which it occurs in the real world, and understand how it is being resolved in practice. We make the following key contributions:

a) We first obtain practitioner inputs on what constitutes explanation disagreement, and the extent to which they encounter this problem in their day-to-day workflow. To this end, we conduct semi-structured

interviews[*] with data scientists (N = 25) who regularly work with explainability tools.

b) Using the insights obtained from the aforementioned interviews, we formalize the notion of explanation disagreement, and propose a novel evaluation framework which can quantitatively measure the disagreement between any two explanations that explain the same model prediction.

c) We leverage the aforementioned framework to carry out a rigorous empirical analysis with real-world data to quantify the level of disagreement between popular post hoc explanations. We experiment with four real-world datasets, six state-of-the-art explanation methods, and various popular predictive models (e.g., logistic regression, tree based models, deep neural networks, recurrent neural networks such as LSTM, and convolutional neural networks such as ResNet).

d) Lastly, we study how explanation disagreements are resolved in practice. We carry out an online user study with data scientists (N = 24) where we show them pairs of explanations that disagree with each other, and ask them which explanation (if any) would they rely on and why. At the end of this survey, we also ask the participants to provide a high-level description of the strategies they use to resolve explanation disagreements in their day-to-day workflow.

Results from our empirical analysis, user interviews and studies indicate that state-of-the-art explanation methods often disagree in terms of the explanations they output, and worse yet, there do not seem to be any principled, well established approaches that ML practitioners employ to resolve these disagreements. More specifically, 84% of our interview participants reported encountering the disagreement problem in their day-to-day workflow. Our empirical analysis further confirmed that explanations generated by state-of-the-art methods often disagree with each other, and this phenomenon persists across various model classes and data modalities. Furthermore, 86% of our online user study responses indicated that ML practitioners either employed arbitrary heuristics (e.g., choosing a favorite method) or just simply did not know how to resolve the disagreement problem. Our findings not only shed light on the previously unexplored disagreement problem, but also underscore the importance of developing principled evaluation metrics to effectively compare explanations, and educating practitioners about the same.

## 2  Related Work

Our work builds on the vast literature in explainable ML. We discuss prior works and their connections to this research.

**Inherently Interpretable Models and Post hoc Explanations.** Many approaches learn inherently interpretable models for various tasks including classification and clustering. Examples of such models include decision trees, decision lists [31], decision sets [29], prototype based models [12, 26], and generalized additive models [13, 34]. However, complex models such as deep neural networks often achieve higher accuracy than simpler models [42]; thus, there has been a lot of interest in constructing post hoc explanations to understand their behavior. To this end, several techniques have been proposed to construct *post hoc explanations* of complex models. These techniques differ in their access to the complex model (i.e., black box vs. access to internals), scope of approximation (e.g., global vs. local), search technique (e.g., perturbation based vs. gradient based), and basic units of explanation (e.g., feature importance vs. rule based). For instance, LIME, SHAP, Anchors, BayesLIME and BayesSHAP [35, 39, 40, 49] are *perturbation based local* explanations as they leverage perturbations of individual instances to construct interpretable local approximations (e.g., linear models).On the other hand, methods such as Gradient*Input, SmoothGrad, Integrated Gradients and GradCAM [45, 47, 50, 51] are *gradient based local* explanations as they leverage gradients computed with respect to input dimensions of individual instances to explain model predictions. An alternate class of methods known as *global* explanations attempt to summarize the behavior of black box models as a whole [10, 30]. In contrast, our work focuses on analyzing the disagreements between explanations generated by state-of-the-art methods.

---

[*]All the user interviews and studies in this work were approved by our institution's IRB.

**Analyzing and Evaluating Post hoc Explanations.** Prior research has studied several notions of explanation quality such as fidelity, stability, consistency and sparsity [33, 37, 49, 54]. Several metrics to quantify each of these aspects of explanation quality were also proposed [14, 20, 33, 54]. As discussed in the introduction, most of these metrics are not general enough to cater to all models or real world settings. Follow up works leveraged these properties and metrics to theoretically and empirically analyze the behavior of popular post hoc explanations [6, 7, 9, 15, 16, 19, 32, 48]. More specifically it has been shown that these explanations can be inconsistent or unstable [19, 48], prone to fair washing [8, 28, 48], and can be unfaithful to the model to the extent that their usefulness can be severely compromised [43]. However, none of these works highlight or study the disagreement problem in explainable ML which is the focus of our work.

**Human Factors in Explainability.** Many user studies evaluate how well humans can understand and utilize explanations [17]. Kaur et al. [25] show that data scientists do not have a good understanding of the state-of-the-art interpretability techniques, and are unable to effectively leverage them in debugging ML models. Bhatt et al. [11], conduct a survey to understand the use-cases for local explanations. Hong et al. [22] conduct a similar survey to identify a variety of stakeholders across the model lifecycle, and higlight core goals: improving the model, making decisions, and building trust in the model. Furthermore, Lakkaraju and Bastani [28] study if misleading explanations can fool domain experts into deploying racially biased models. Similarly, Poursabzi-Sangdeh et al. [38] find that supposedly-interpretable models can lead to a decreased ability to detect and correct model mistakes. Lage et al. [27] use insights from rigorous human-subject experiments to inform regularizers used in explanation algorithms. However, none of these works focus on understanding if and how often practitioners face explanation disagreement, and how they resolve it.

# 3   Understanding and Measuring Disagreement between Model Explanations

In this section, we discuss practitioner perspectives on what constitutes disagreement between two explanations, and then formalize the notion of explanation disagreement. To this end, we first describe the study that we carry out with data scientists to understand what constitutes explanation disagreement, and the extent to which they encounter this problem in practice. We then discuss the insights from this study, and leverage these insights to propose a novel framework which can quantitatively measure the disagreement between any two explanations.

## 3.1   Characterizing Explanation Disagreement Using Practitioner Inputs

Here, we describe the study that we conducted with data scientists to characterize explanation disagreement, and then outline our findings and insights from this study.

### 3.1.1   Interviews with Practitioners: Study Details

great,We conducted 30-minute long semi-structured interviews with 25 data scientists who employ explainability techniques to understand model behavior and explain it to their customers and managers. All of these data scientists were employed in for-profit organizations, and worked for various companies in the technology and financial services sectors in the United States. Furthermore, all the participants used state-of-the-art (local) post hoc explanation methods such as LIME, SHAP, and gradient based methods in their day-to-day workflow. 19 of these participants (76%) were male, and 6 of them (24%) were female. 16 participants (64%) had more than 2 years of experience working with explainability techniques, and the remaining 9 (36%) had about 8 to 12 months of experience. Our interviews included, but were not limited to the following questions: Q1) *How often do you use multiple explanation methods to understand the same model prediction?* Q2) *What constitutes disagreement between two explanations that explain the same model prediction?* Q3) *How often do you encounter disagreements between explanations output by different methods for the same model prediction?*

### 3.1.2 Findings and Insights

Our study revealed a wealth of information about how data scientists utilize explanation methods and their perspectives on disagreement between explanations. 22 out of the 25 participants (88%) said that they almost always use multiple explanation methods to understand the same model prediction. Furthermore, 21 out of the 25 participants (84%) mentioned that they have often run into some form of disagreement between explanations output by different methods for the same prediction. They also elaborated on when they think two explanations disagree:

**Top features are different:** Most of the popular post hoc explanation methods (e.g., LIME, SHAP, Gradient based methods) return a feature importance value associated with each feature. These values indicate which features contribute the most either positively or negatively (i.e., the top features) to the prediction. 21 out of the 25 participants (84%) in our study mentioned that such a set of top features is *"the most critical piece of information"* that they rely on in their day-to-day workflow. They also noted that they typically look at the top 5 to 10 features provided by an explanation for each prediction. When two explanations have different sets of top features, they consider it to be a disagreement.

**Ordering among top features is different:** 18 out of 25 participants (72%) in our study indicated that they also consider the ordering among the top features very carefully in their workflow. Therefore, they consider a mismatch in the ordering of the top features provided by two different explanations to be a disagreement.

**Direction of top feature contributions is different:** 19 out of 25 participants (76%) mentioned that the *sign* or *direction* of the feature contribution (is the feature contributing positively or negatively to the predicted class?) is another critical piece of information. Any mismatch in the signs of the top features between two explanations is a sign of disagreement. As remarked by one of the participants, *"I saw an explanation indicating that a top feature bankruptcy contributes positively to a particular loan denial, and another explanation saying that it contributes negatively. That is a clear disagreement. The model prediction can be trusted with the former explanation, but not with the latter."*.

**Relative ordering of certain features is different:** 16 of our study participants (64%) indicated that they also look at relative ordering between certain features of interest; and if explanations provide contradicting information about this aspect, then it is considered a disagreement. For example, one of the participants remarked, *"I often check if salary is more important than credit score in loan approvals. If one explanation says salary is more important than credit score, and another says credit score is more important than salary; then it is a disagreement."*

A very striking finding from our study is that participants typically characterize explanation disagreement based on factors such as mismatch in top features, feature ordering, and directions of feature contributions, but not on the feature importance values output by different explanation methods. 24 out of 25 participants (96%) in our study opine that feature importance values output by different explanation methods are not directly comparable. They also note that this is due to the fact that while LIME outputs coefficients of a linear model as feature importance values, SHAP outputs Shapley values as feature attributions which sum to the probability of the predicted class. So, they don't try to base explanation disagreement on these numbers not being equal or similar. One of our participants succinctly summarized practitioners' perspective on this explanation disagreement problem – *"The values generated by different explanation methods are clearly different. So, I would not characterize disagreement based on that. But, I would at least want the explanations they output to give me consistent insights. The explanations should agree on what are the most important features, the ordering among them and so on for me to derive consistent insights. But, they don't!"*

## 3.2 Formalizing the Notion of Explanation Disagreement

Our study indicates that ML practitioners consider the following key aspects when they think about explanation disagreement: a) the extent to which explanations differ in the top-$k$ features, the signs (or directions of contribution) and the ordering of these top-$k$ features, and b) the extent to which explanations differ in the relative ordering of certain features of interest. To capture these intuitions about explanation disagreement, we propose six different metrics, namely, *feature agreement*, *rank agreement*, *sign agreement*, *signed rank agreement*, *rank correlation*, and *pairwise rank agreement*. While the first four metrics capture disagreement w.r.t. the top-$k$ features of the explanations, the last two metrics capture disagreement w.r.t. a selected set of features which could be provided as input by an end user.

### 3.2.1 Measuring Disagreement w.r.t. Top-k Features

We now define four metrics, which capture specific aspects of explanation disagreement w.r.t. the top-$k$ features.[*] For all metrics, lower values indicate higher disagreement.

**Feature Agreement:**   ML practitioners in our study (Section 3.1) clearly indicated that a key notion of disagreement between a pair of explanations is that they output different top-$k$ features. To capture this notion, we introduce the feature agreement metric which computes the fraction of common features between the sets of top-$k$ features of two explanations. Given two explanations $E_a$ and $E_b$, the feature agreement metric can be formally defined as:

$$FeatureAgreement(E_a, E_b, k) = \frac{|top\_features(E_a, k) \cap top\_features(E_b, k)|}{k}$$

where $top\_features(E, k)$ returns the set of top-$k$ features (based on the magnitude of the feature importance values) of the explanation $E$. If the sets of top-$k$ features of explanations $E_a$ and $E_b$ match, then $FeatureAgreement(E_a, E_b, k) = 1$.

**Rank Agreement:**   Practitioners in our study also indicated that if the ordering of the top-$k$ features is different for two explanations (even if the feature sets are the same), then they consider it to be a disagreement. To capture this notion, we introduce the rank agreement metric which computes the fraction of features that are not only common between the sets of top-$k$ features of two explanations, but also have the same position in the respective rank orders. Rank agreement is a stricter metric than feature agreement since it also considers the ordering of the top-$k$ features. Given two explanations $E_a$ and $E_b$, the rank agreement metric ($RankAgreement(E_a, E_b, k)$) can be formally defined as:

$$\frac{|\bigcup_{s \in S} \{s \mid s \in top\_features(E_a, k) \land s \in top\_features(E_b, k) \land rank(E_a, s) = rank(E_b, s)\}|}{k}$$

where $S$ is the complete set of features in the data, $top\_features(E, k)$ is defined as above, and $rank(E, s)$ returns the position (or the rank) of the feature $s$ according to the explanation $E$. If the rank-ordered lists of top-$k$ features of explanations $E_a$ and $E_b$ match, then $RankAgreement(E_a, E_b, k) = 1$.

**Sign Agreement:**   In our study, practitioners also mentioned that they consider two explanations to disagree if the feature attribution signs or the directions of feature contribution (does a feature contribute positively or negatively to the prediction?) do not align for the top-$k$ features. To capture this notion, we introduce the sign agreement metric which computes the fraction of features that are not only common between the sets of top-$k$ features of two explanations, but also share the same sign (direction of contribution)

---

[*]The top-$k$ features of an explanation are typically computed only based on the magnitude of the feature importance values and not the signs.

in both explanations. Sign agreement is a stricter metric than feature agreement since it also considers signs (directions of contributions) of the top-$k$ features. More formally, $SignAgreement(E_a, E_b, k)$ is defined as:

$$\frac{|\bigcup_{s \in S}\{s \mid s \in top\_features(E_a, k) \wedge s \in top\_features(E_b, k) \wedge sign(E_a, s) = sign(E_b, s)\}|}{k}$$

where $sign(E, s)$ returns the sign (direction of contribution) of the feature $s$ according to the explanation $E$.

**Signed Rank Agreement:**  This metric fuses together all the above notions, and computes the fraction of features that are not only common between the sets of top-$k$ features of two explanations, but also share the same feature attribution sign (direction of contribution) and position (rank) in both explanations. Signed rank agreement is the strictest compared to all the aforementioned metrics since it considers both the ordering and the signs (directions of contributions) of the top-$k$ features. More formally, $SignedRankAgreement(E_a, E_b, k)$ is formulated as:

$$\begin{aligned} |\bigcup_{s \in S}\{\{s \mid s \in top\_features(E_a, k) \wedge s \in top\_features(E_b, k) \\ \wedge\ sign(E_a, s) = sign(E_b, s) \wedge rank(E_a, s) = rank(E_b, s)\}| \\ \hline k \end{aligned}$$

where $top\_features$, $sign$, $rank$ are all as defined above. $SignedRankAgreement(E_a, E_b, k) = 1$ if the top-$k$ features of two explanations match on all aspects (i.e., features, feature attribution signs, rank ordering) barring the exact feature importance values.

### 3.2.2 Measuring Disagreement w.r.t. Features of Interest

Practitioners also indicated that they consider two explanations to be different if the relative ordering of features of interest (e.g., salary and credit score discussed in Section 3.1) differ between the two explanations. To formalize this notion, we introduce the two metrics below.

**Rank Correlation:**  We adopt a standard rank correlation metric (i.e., Spearman's rank correlation coefficient) to measure the agreement between feature rankings provided by two explanations for a selected set of features. In practice, this selected set of features corresponds to features that are of interest to end users, and can be provided as input by end users. Given two explanations $E_a$ and $E_b$, rank correlation can be computed as:

$$RankCorrelation(E_a, E_b, F) = r_s(Ranking(E_a, F), Ranking(E_b, F))$$

where $F$ is a selected set of features potentially input by an end user, $r_s$ computes Spearman's rank correlation coefficient, and $Ranking(E, F)$ assigns ranks to features in $F$ based on explanation $E$. Lower values indicate higher disagreement.

**Pairwise Rank Agreement:**  Pairwise rank agreement takes as input a set of features that are of interest to the user, and captures if the relative ordering of every pair of features in that set is the same for both the explanations i.e., if feature A is more important than B according to one explanation, then the same should be true for the other explanation. More specifically, this metric computes the fraction of feature pairs for which the relative ordering is the same between two explanations. More formally:

$$PairwiseRankAgreement(E_a, E_b, F) = \frac{\sum_{i,j \text{ for } i<j} \mathbb{1}[RelativeRanking(E_a, f_i, f_j) = RelativeRanking(E_b, f_i, f_j)]}{\binom{|F|}{2}}$$

where $F = \{f_1, f_2 \cdots\}$ is a selected set of features input by an end user, $RelativeRanking(E, f_i, f_j)$ is an indicator function which returns 1 if feature $f_i$ is more important than feature $f_j$ according to explanation $E$, and 0 otherwise.

# 4 Empirical Analysis of Explanation Disagreement

We leverage the metrics outlined in Section 3 and carry out a comprehensive empirical analysis with six state-of-the-art explanation methods and four real-world datasets to study the explanation disagreement problem. In this section, we describe the datasets that we use (Section 4.1), our experimental setup (Section 4.2), and key findings (Section 4.3).

## 4.1 Datasets

To carry out our empirical analysis, we leverage four well known datasets spanning three different data modalities (tabular, text, and images). For **tabular** data, we use the Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) dataset [2] and the German Credit dataset [3]. This dataset comprises of seven features capturing information about the demographics, criminal history, and prison time of 4,937 defendants. Each defendant in the data is labeled either as high-risk or low-risk for recidivism based on the COMPAS algorithm's risk score. The German Credit dataset contains twenty features capturing the demographics, credit history, bank account balance, loan information, and employment information of 1,000 loan applicants. The class label here is a loan applicant's credit risk (high or low). For **text** data, we use Antonio Gulli (AG)'s corpus of news articles (AG_News) [1]. The dataset contains 127,600 sentences (collected from 1,000,000+ articles from 2,000+ sources with a vocabulary size of 95,000+ words). The class label is the topic of the article from which a sentence was obtained (World, Sports, Business, or Science/Technology). For **image** data, we use the ImageNet$-1k$ [4, 44] object recognition dataset. It contains 1,381,167 images belonging to 1000 object categories. We experiment with images from PASCAL VOC 2012 [5] which provides segmentation maps that can be directly used as super-pixels for the explanation methods.

## 4.2 Experimental Setup

We train a variety of black box models on the data. In case of tabular data, we train four models: logistic regression, densely-connected feed-forward neural network, random forest, and gradient boosted trees. In case of text data, we train a widely-used vanilla LSTM-based text classifier on AG_News [53] corpus. For image data, we use the pre-trained ResNet-18 [21] for ImageNet.

Next, we apply six state-of-the-art post hoc explanation methods to explain the black box models' predictions for a set of test data points. We apply two perturbation-based explanation methods (LIME [41] and KernelSHAP [36]), and four gradient-based explanation methods (Vanilla Gradient [47], Gradient*Input [46], Integrated Gradients [51], and SmoothGrad [50]). In case of explanation methods with a sample size hyperparameter, we either run the explanation method to convergence (i.e., select a sample size such that an increase in the number of samples does not significantly change the explanations) or use a sample size that is much higher than the sample size recommended by previous work.

We then evaluate the (dis)agreement between the explanation methods using the metrics described in Section 3.2. For tabular and text data, we apply rank correlation and pairwise rank agreement across all features; and feature agreement, rank agreement, sign agreement, signed rank agreement across top-$k$ features for varying values of $k$. For image data, metrics that operate on the top-$k$ features are more applicable to super-pixels. Thus, we apply the six disagreement metrics on explanations output by LIME and KernelSHAP (which leverage super pixels), and calculate rank correlation (across all pixels as features) between the explanations output by gradient-based methods. See Appendix A for details.

## 4.3 Results and Insights

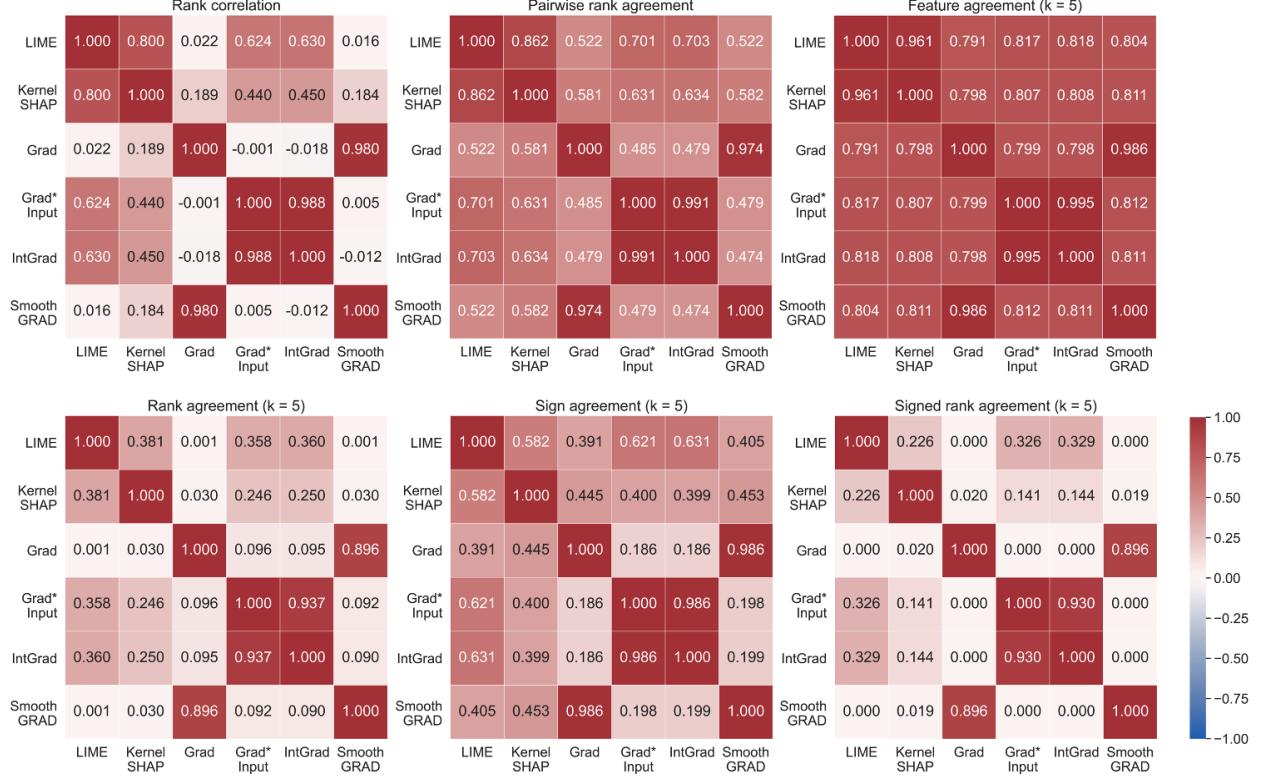We discuss the results of our empirical analysis for each of the three data modalities.

Figure 1: Disagreement between explanation methods for neural network model trained on COMPAS dataset measured by six metrics: rank correlation and pairwise rank agreement across all features, and feature, rank, sign, and signed rank agreement across top $k = 5$ features. Heatmaps show the average metric value over test set data points for each pair of explanation methods, with lighter colors indicating stronger disagreement. Across all six heatmaps, the standard error ranges between 0 and 0.009.

### 4.3.1 Tabular Data

Figure 1 shows the disagreement between various pairs of explanation methods for the neural network model trained on COMPAS dataset. We computed the six metrics outlined in Section 3.2 where we used $k = 5$ (out of 7 features) for metrics that focus on the top features. Each cell in the heatmap shows the metric value averaged over the test data points for each pair of explanation methods with lighter colors indicating more disagreement. We see that explanation methods tend to exhibit slightly higher values on pairwise rank agreement and feature agreement metrics, and relatively lower values on other metrics (indicating more disagreement).

We next study the effect of the number of top features on the degree of disagreement. Figure 2 shows the disagreement of explanation methods for the neural network model trained on COMPAS dataset. We computed rank agreement (top row) and signed rank agreement (bottom row) at top-$k$ features for increasing values of $k$. We see that as the number of top-$k$ features increases, rank agreement and signed rank agreement decrease. This indicates that, as $k$ increases, top-$k$ features of a pair of explanation methods are less likely to contain shared features with the same rank (as measured by rank agreement) or shared features with the same rank and sign (as measured by signed rank agreement). These patterns are consistent across other models trained on the COMPAS dataset. See Appendix B.1.

In addition, across all metrics, values of $k$, and models, the specific explanation method pairs of Grad-SmoothGrad and Grad*Input-IntGrad consistently exhibit strong agreement while the pairs of Grad-IntGrad, Grad-Grad*Input, SmoothGrad-Grad*Input and SmoothGRAD-IntGrad consistently exhibit strong disagree-
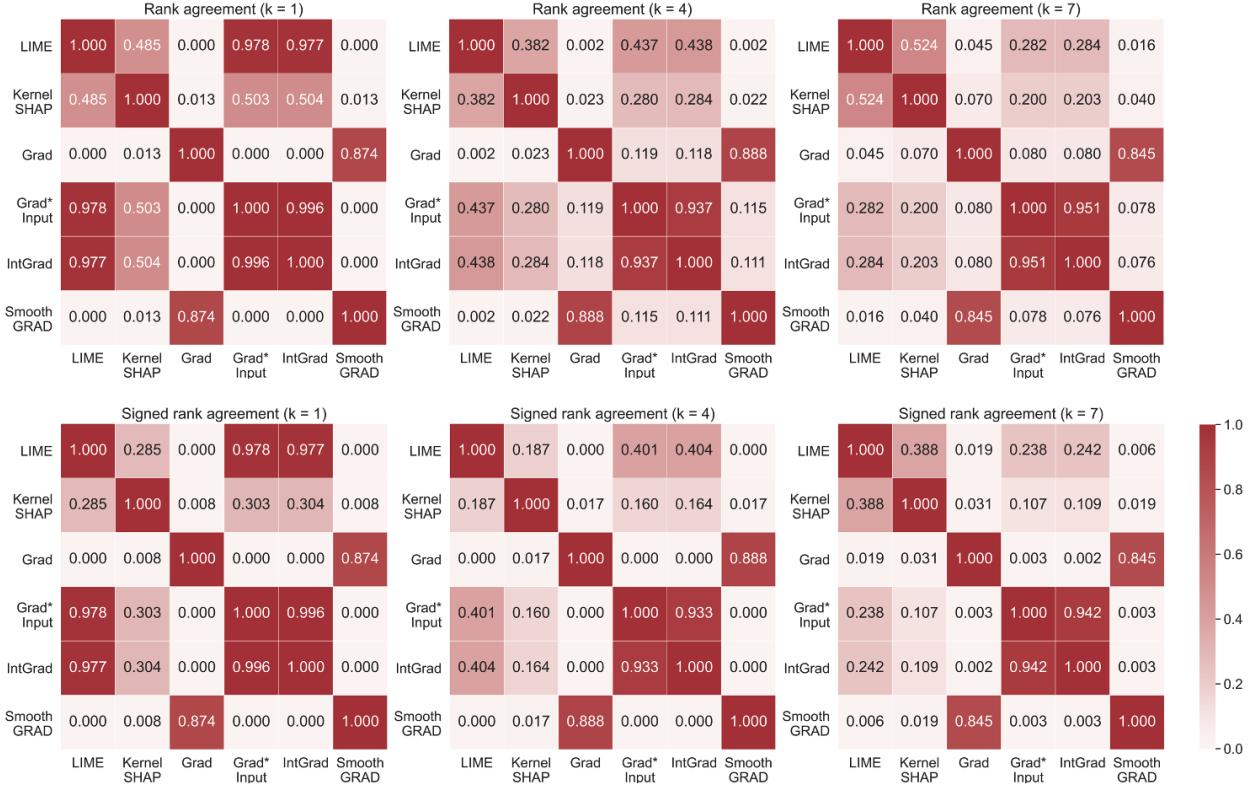
Figure 2: Disagreement between explanation methods for neural network model trained on COMPAS dataset measured by rank agreement (top row) and signed rank agreement (bottom row) at top-$k$ features for increasing values of $k$. Each cell in the heatmap shows the metric value averaged over test set data points for each pair of explanation methods, with lighter colors indicating stronger disagreement. Across all six heatmaps, the standard error ranges between 0 and 0.003.

ment. This suggests a dichotomy among gradient-based explanation methods, i.e., certain gradient-based explanation methods are consistent with one another while others are inconsistent with one another. See Appendix B.1 for more details.

Furthermore, there are varying degrees of disagreement among pairs of explanation methods. For example, for the neural network model trained on the COMPAS dataset, rank correlation displays a wide range of values across explanation method pairs, with 10 out of 15 explanation method pairs even exhibiting negative rank correlation when explaining multiple data points. This is shown in the left panel of Figure 3, which displays the rank correlation over all the features among all pairs of explanation methods for neural network model trained on COMPAS dataset.

All the patterns discussed above are also generally reflected in the German Credit dataset (Appendix B.2). However, explanation methods tend to display stronger disagreement for the German Credit dataset than for the COMPAS dataset. For example, rank agreement and signed rank agreement are lower for the German Credit dataset than for the COMPAS dataset at top 25%, 50%, 75%, and 100% of features for both logistic regression and neural network models (Appendix B.1 and B.2). One possible reason is that the German Credit dataset has a larger set of features than the COMPAS dataset, resulting in a larger number of possible ranking and sign combinations assigned by a given explanation method and making it less likely for two explanation methods to produce consistent explanations.

Lastly, explanation methods display trends associated with model complexity. For example, the disagreement between explanation methods is similar or stronger for the neural network model than for the logistic

regression model across metrics and values of $k$, for both COMPAS and German Credit datasets (Appendix B.1 and B.2). In addition, explanation methods show similar levels of disagreement for the random forest and gradient-boosted tree models. These trends suggest that disagreement among explanation methods may increase with model complexity. As the complexity of the black box model increases, it may be more difficult to accurately approximate the black box model with a simpler model (LIME's strategy, for example) and more difficult to disentangle the contribution of each feature to the model's prediction. Thus, the higher the model complexity, the more difficult it may be for different explanation methods to generate the true explanation and the more likely it may be for different explanation methods to generate differently false explanations, leading to stronger disagreement among explanation methods.
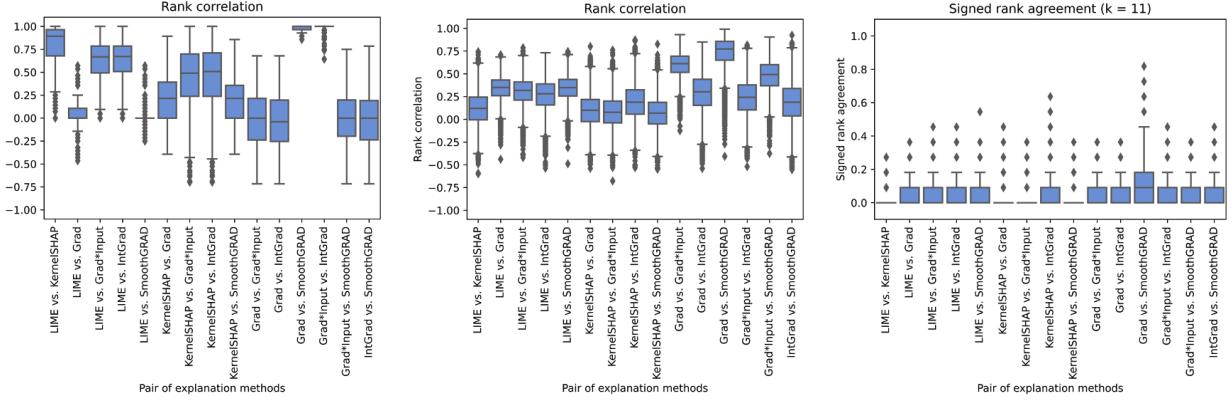


Figure 3: Distribution of rank correlation over all features for neural network model trained on COMPAS (left), and rank correlation across all features (middle) and signed rank agreement across top-11 features (right) for neural network model trained on AG_News.

### 4.3.2 Text Data

In the case of text data, we deal with a high-dimensional feature space where words are features. We plot the six metrics for $k = 11$ which is 25% of the average text length of a sentence (data point) in the dataset (Figure 4). As can be seen, we observe severe disagreements across all the six disagreement metrics. Rank agreement and signed rank agreement are the lowest between explanations with values under 0.1 for most cases indicating disagreement in over 90% of the top-$k$ features. Trends are quite similar for rank correlation and feature agreement with better agreement between gradient based explanation methods, such as a feature agreement of 0.61 for Grad*Input and Gradient method.

In addition to the overall disagreements, we also notice specific patterns of agreement between a group of explanation methods. Based on middle and right panels of Figure 3, we notice that there is a high rank correlation between pairs of gradient based explanations. Although Integrated Gradients has the lowest correlation with the rest of the gradient methods, still this correlation is significantly higher compared to its correlation with KernelSHAP and LIME. We also notice that disagreement is lower between LIME and other explanation methods compared to KernelSHAP and other methods (e.g., rank correlation of 0.4-0.6 for LIME as opposed to 0.2-0.4 for KernelSHAP. See Appendix B for other metrics). Finally, we see a higher disagreement in explanation methods for text compared to tabular data which suggests that disagreement may worsen as the number of features increases. We also observe a similar agreement pattern between LIME and other methods which was also observed earlier in our experiments with tabular datasets, hence, indicating that LIME explanations are most aligned with other post hoc explanation methods.
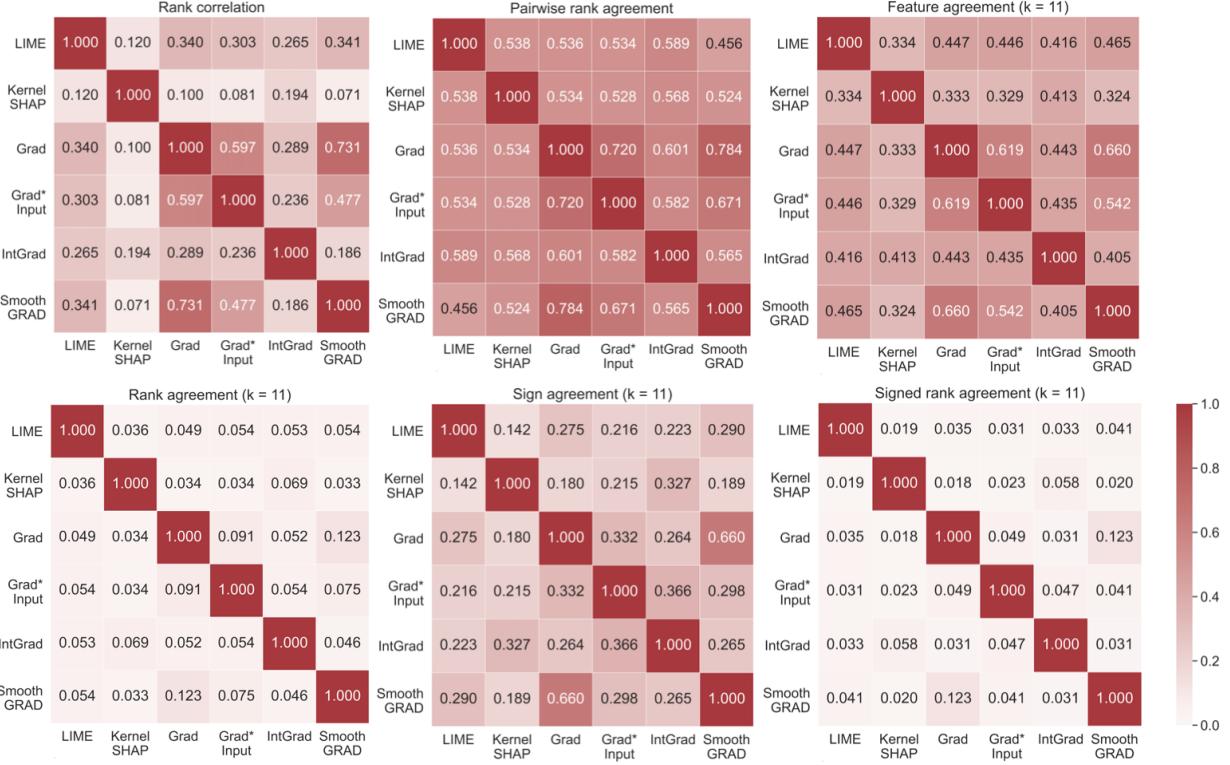
**Rank correlation**

| | LIME | Kernel SHAP | Grad | Grad* Input | IntGrad | Smooth GRAD |
|---|---|---|---|---|---|---|
| LIME | 1.000 | 0.120 | 0.340 | 0.303 | 0.265 | 0.341 |
| Kernel SHAP | 0.120 | 1.000 | 0.100 | 0.081 | 0.194 | 0.071 |
| Grad | 0.340 | 0.100 | 1.000 | 0.597 | 0.289 | 0.731 |
| Grad* Input | 0.303 | 0.081 | 0.597 | 1.000 | 0.236 | 0.477 |
| IntGrad | 0.265 | 0.194 | 0.289 | 0.236 | 1.000 | 0.186 |
| Smooth GRAD | 0.341 | 0.071 | 0.731 | 0.477 | 0.186 | 1.000 |

**Pairwise rank agreement**

| | LIME | Kernel SHAP | Grad | Grad* Input | IntGrad | Smooth GRAD |
|---|---|---|---|---|---|---|
| LIME | 1.000 | 0.538 | 0.536 | 0.534 | 0.589 | 0.456 |
| Kernel SHAP | 0.538 | 1.000 | 0.534 | 0.528 | 0.568 | 0.524 |
| Grad | 0.536 | 0.534 | 1.000 | 0.720 | 0.601 | 0.784 |
| Grad* Input | 0.534 | 0.528 | 0.720 | 1.000 | 0.582 | 0.671 |
| IntGrad | 0.589 | 0.568 | 0.601 | 0.582 | 1.000 | 0.565 |
| Smooth GRAD | 0.456 | 0.524 | 0.784 | 0.671 | 0.565 | 1.000 |

**Feature agreement (k = 11)**

| | LIME | Kernel SHAP | Grad | Grad* Input | IntGrad | Smooth GRAD |
|---|---|---|---|---|---|---|
| LIME | 1.000 | 0.334 | 0.447 | 0.446 | 0.416 | 0.465 |
| Kernel SHAP | 0.334 | 1.000 | 0.333 | 0.329 | 0.413 | 0.324 |
| Grad | 0.447 | 0.333 | 1.000 | 0.619 | 0.443 | 0.660 |
| Grad* Input | 0.446 | 0.329 | 0.619 | 1.000 | 0.435 | 0.542 |
| IntGrad | 0.416 | 0.413 | 0.443 | 0.435 | 1.000 | 0.405 |
| Smooth GRAD | 0.465 | 0.324 | 0.660 | 0.542 | 0.405 | 1.000 |

**Rank agreement (k = 11)**

| | LIME | Kernel SHAP | Grad | Grad* Input | IntGrad | Smooth GRAD |
|---|---|---|---|---|---|---|
| LIME | 1.000 | 0.036 | 0.049 | 0.054 | 0.053 | 0.054 |
| Kernel SHAP | 0.036 | 1.000 | 0.034 | 0.034 | 0.069 | 0.033 |
| Grad | 0.049 | 0.034 | 1.000 | 0.091 | 0.052 | 0.123 |
| Grad* Input | 0.054 | 0.034 | 0.091 | 1.000 | 0.054 | 0.075 |
| IntGrad | 0.053 | 0.069 | 0.052 | 0.054 | 1.000 | 0.046 |
| Smooth GRAD | 0.054 | 0.033 | 0.123 | 0.075 | 0.046 | 1.000 |

**Sign agreement (k = 11)**

| | LIME | Kernel SHAP | Grad | Grad* Input | IntGrad | Smooth GRAD |
|---|---|---|---|---|---|---|
| LIME | 1.000 | 0.142 | 0.275 | 0.216 | 0.223 | 0.290 |
| Kernel SHAP | 0.142 | 1.000 | 0.180 | 0.215 | 0.327 | 0.189 |
| Grad | 0.275 | 0.180 | 1.000 | 0.332 | 0.264 | 0.660 |
| Grad* Input | 0.216 | 0.215 | 0.332 | 1.000 | 0.366 | 0.298 |
| IntGrad | 0.223 | 0.327 | 0.264 | 0.366 | 1.000 | 0.265 |
| Smooth GRAD | 0.290 | 0.189 | 0.660 | 0.298 | 0.265 | 1.000 |

**Signed rank agreement (k = 11)**

| | LIME | Kernel SHAP | Grad | Grad* Input | IntGrad | Smooth GRAD |
|---|---|---|---|---|---|---|
| LIME | 1.000 | 0.019 | 0.035 | 0.031 | 0.033 | 0.041 |
| Kernel SHAP | 0.019 | 1.000 | 0.018 | 0.023 | 0.058 | 0.020 |
| Grad | 0.035 | 0.018 | 1.000 | 0.049 | 0.031 | 0.123 |
| Grad* Input | 0.031 | 0.023 | 0.049 | 1.000 | 0.047 | 0.041 |
| IntGrad | 0.033 | 0.058 | 0.031 | 0.047 | 1.000 | 0.031 |
| Smooth GRAD | 0.041 | 0.020 | 0.123 | 0.041 | 0.031 | 1.000 |

Figure 4: Disagreement between explanation methods for the LSTM model trained on the AG_News dataset using $k = 11$ features for metrics operating on top-$k$ features, and all features for other metrics. Each heatmap shows the metric value averaged over test data for each pair of explanation methods. Lighter colors indicate more disagreement. Standard error ranges from 0.0 to 0.0025 for all six metrics.

### 4.3.3 Image Data

While LIME and KernelSHAP consider super pixels of images as features, gradient based methods consider pixels as features. Furthermore, the notion of top-$k$ features and the metrics we define on top-$k$ features are not semantically meaningful when we consider pixels as features. Given this, we compute all the six metrics to capture disagreement between explanations output by LIME and KernelSHAP with super pixels as features. We also compute rank correlation on all the pixels (features) to capture disagreement between explanations output by gradient based methods.

Unlike earlier trends with tabular and text data, we see higher agreement between KernelSHAP and LIME on all the six metrics: rank correlation of 0.8977, pairwise rank agreement of 0.9302, feature agreement of 0.9535, rank agreement of 0.8478, sign agreement of 0.9218 and signed rank agreement of 0.8193. However, the trends are quite opposite when we compute rank correlation at pixel-level for gradient-based methods (See Appendix B). For instance, rank correlation between Integrated Gradients and SmoothGrad is 0.001 (indicating high disagreement). The disagreement is similarly quite high in case of other pairs of gradient based methods. This suggests that disagreement could potentially vary significantly based on the granularity of image representation.

# 5 Resolving the disagreement problem in practice: a qualitative study

In order to understand how practitioners resolve the disagreement problem, we conduct a qualitative user study targeted towards explainability practitioners. We now describe our user study design and discuss our findings.

## 5.1 User Study Design

In total, 25 participants participated in our study, 13 from academia and 12 from industry. Participants from academia were graduate students, and postdoctoral researchers, while participants from industry were data scientists and ML engineers from three different firms. 20 of these participants indicated that they have used explainability methods in their work in a variety of ways, including doing research, helping clients explain their models, and debugging their own models. Following the setup in Section 4, we asked participants to compare the output of five pairs of explainability methods on the predictions made by the neural network we trained on the COMPAS dataset. We chose the COMPAS dataset because it only has 7 features, making it easy for participants to understand the explanations.

First, the participants are shown an information page explaining the COMPAS risk score binary prediction setting and various explainability algorithms. We indicate that we trained a neural network to predict the COMPAS risk score (low or high) from the seven COMPAS features. We also give a brief description of each of these seven features to the participant and tell them to assume that the criminal defendant's risk of recidivism is correctly predicted to be high risk. In this information page, we also briefly introduce and summarize the six explainability algorithms we use in the study (LIME, KernelSHAP, Gradient, Gradient*Input, SmoothGrad, and Integrated Gradients). Finally, we provide links to the papers describing each of the algorithms. We include a screenshot of this information page in Appendix C.1.

Next, each of the participants is shown a series of 5 prompts, a sample of which is shown in Figure 5. Each prompt presents two explanations of our neural network model's prediction corresponding to a particular data point generated using two different explanation methods (e.g., LIME and KernelSHAP in Figure 5). For each of the 15 pairs of explanation methods, we chose a different data point from COMPAS to run the two methods on, giving us a set of 15 prompts. These prompts were picked to showcase various levels of agreement. We display the full set of $k = 7$ COMPAS features, showing the feature importance of each feature. Red and blue bars indicate that the feature contributes negatively and positively respectively to the predicted class. The participants were first asked the question *"To what extent do you think the two explanations shown above agree or disagree with each other?"* and given four choices: *completely agree, mostly agree, mostly disagree*, and *completely disagree*. If the participant indicated any level of disagreement (any of the latter 3 choices), we then asked *"Since you believe that the above explanations disagree (to some extent), which explanation would you rely on?"* and presented with three choices: the two explainability methods shown and *"it depends"*. They were then asked to explain their response. The users were allowed to take as much time as they wanted to complete the study.

## 5.2 Results and Insights

We now discuss the results and findings from our user study in Sections 5.2.1-5.2.4.

### 5.2.1 Do practitioners observe disagreements?

We aggregated the responses to the first question in each prompt, *"To what extent do you think the explanations shown above agree or disagree with each other?"*. Overall, 4%, 28%, 50%, and 18% of responses indicate *completely agree*, *mostly agree*, *mostly disagree*, and *completely disagree*, respectively, highlighting that there is significant disagreement among our prompts. See Appendix C.4 for more details.

Figure 5: The user interface for a prompt. The user is shown two explanations for a COMPAS data point, showing the feature importance value of each of the 7 features. Red and blue indicate negative and positive feature values, respectively. See the text for more details.

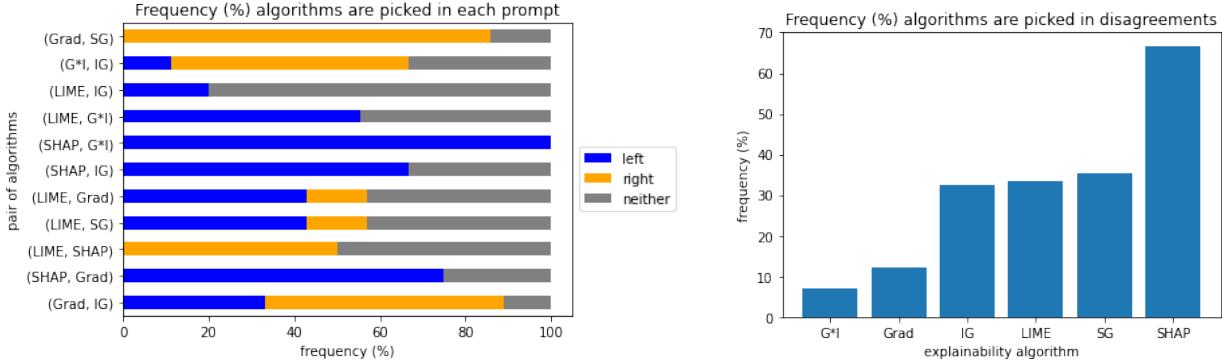### 5.2.2   Are certain explanations favored over others?

Next, since different algorithms have different levels of popularity, we analyze if certain algorithms are chosen more often in disagreements. Figure 6a shows the distribution of how participants resolved disagreements for each prompt (dropping prompts with 4 or less responses). We first emphasize that there is high variability in how participants chose to resolve disagreements, showing a lack of consensus for the majority of prompts. However, when participants do decide to choose an algorithm rather than abstaining, they often choose the same algorithm. For example, in the Gradient vs SmoothGrad (top row in Figure 6a), participants either chose SmoothGrad over Gradient or chose neither. We also aggregate these choices over all prompts, and in Figure 6b, we plot how often each of the six explainability algorithms is chosen, finding that indeed, certain algorithms were favored over others. While KernelSHAP was chosen 66.7% of the time when there are disagreements, Gradient × Input was only chosen 7.0% of the time. We include a further explanation of why participants chose each of the explanations in Appendix C.5, including quotes from participants that supported each algorithm.

### 5.2.3   How do practitioners resolve disagreements?

Throughout all six explainability techniques, we find three unifying themes that dictated why participants chose one explanation over the other. We give a high-level description of these themes below, highlighting direct quotes from participants in Table 1.

*1. One method is inherently better than the other because of their theory or publication time (33%):* Participants often indicated preference towards a particular method without referencing the shown explanation citing features such as the paper's publication time (more recent papers are better), the theory behind the method, and the method's stability.

*2. One of the generated explanations matches intuition better (32%):*   Participants frequently said that one method's explanation aligned with their intuition better, citing the absolute and relative values of specific

(a) The frequency with which each of the explanations in a pair is selected upon disagreement. The blue, gold, and grey bars show the percentage of participants (X axis) that picked the left, right, and neither algorithm when presented with the pair of algorithms shown on the Y axis.

(b) The frequency with which each of the explanations was chosen when there is a disagreement. X axis indicates the explainability algorithms and Y axis indicates the frequency.

Figure 6: Sub-figures highlight which algorithms participants chose when the explanations they were shown disagreed. In (a), we show how participants resolved each particular prompt. In (b), we show the overall frequencies with which each explanation method was selected.

features as evidence.

*3. LIME and SHAP are better because COMPAS is tabular data (23%):* Participants said that they mainly used LIME and SHAP for tabular data and commonly cited this as their sole reason.

Table 1: Themes summarizing how participants decided between explanations when faced with disagreement along with quotes.

| Theme Highlighted | Sample Quotes |
|---|---|
| **1. One method's paper/theory suggests that it's inherently better (33%).** | • *"I have no reason to believe the gradient holds anywhere other than very locally."* <br> • *"[IG is] more rigorous [than SmoothGrad] based on the paper and axioms"* <br> • *"gradient explanations are more unstable"* |
| **2. One explanation matches intuition better (32%).** | • *"seems unlikely that all features contributed to a positive classification"* <br> • *"features such as priors_count and length of stay [are] important for determining"* <br> • *"Gradient*Input only consider[s] sensitive features (age, race) as impactful which could be a sign of a biased underlying data distribution"* |
| **3. LIME/SHAP are better for tabular data (23%).** | • *"I use LIME for structured data"* <br><br> • *"SHAP is more commonly used [than Gradient] for tabular data"* |

### 5.2.4 Experiencing and resolving disagreements in day to day work:

After answering all 5 prompts, participants were asked a set of questions to help us understand their experience with the disagreement problem in their day to day work. First, to filter out participants who didn't use explainability methods, we asked: *"Have you used explainability methods in your work before?"*. Of the 25 participants, 5 of them indicated that they had not. We asked the other 20 participants further questions to

better understand their experience with the disagreement problem. The full set of questions can also be found in Appendix C.3. Having understood what participants look for to determine disagreement in Section 3.1, we next aim to better understand two crucial questions related to the disagreement problem: *(Q1): Do you observe disagreements between explanations output by state of the art methods in your day to day workflow?* and *(Q2): How do you resolve such disagreements in your day to day workflow?*.

One of the 20 participants declined to respond to (Q1) and (Q2) because they were not a practitioner. Out of the other 19, 14 participants (74%) responded *"yes"* to (Q1), indicating that they did in fact encounter explanation disagreement in practice. Of the remaining 5 who said they did not, 3 said they had not really paid attention to the issue. Through (Q2), we aimed to uncover how participants dealt with the disagreement problem when it arose in practice. Of the 14 participants answering yes to (Q1), their responses to (Q2) can be grouped into 3 categories. 50% had personal heuristics for choosing which algorithms to use (*"data scientists picking their favorite algorithm"*, *"rules of thumb based on results in papers"*). These heuristics varied among participants and included ease of implementation, groundedness of theory, recency of publication, ease of understanding, and documentation of packages. 36% didn't indicate any way to resolve these disagreements, but rather showed confusion and uncertainty (*"no clear answer to me"*). Many of the responses indicated the desire for the research community to make progress and help (*"I hope research community can provide some guidance"*). Therefore, we hope that these responses motivate and inspire future work in this direction. The remaining 14% proposed to use other metrics such as fidelity (*"try and use some metric to measure fidelity"*). See Appendix C.7.

# References

[1] Antonio gulli corpus of news articles. `http://groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html`. Accessed: 2021-01-20.

[2] Propublica article on compas. `https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm`. Accessed: 2021-01-20.

[3] German credit dataset. `https://archive.ics.uci.edu/ml/datasets/statlog+(german+credit+data)`. Accessed: 2021-01-20.

[4] Imagenet. `https://www.image-net.org/challenges/LSVRC/index.php`. Accessed: 2021-01-20.

[5] Pascal voc 2012. `http://host.robots.ox.ac.uk/pascal/VOC/index.html/`. Accessed: 2021-01-20.

[6] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages 9505–9515, 2018.

[7] S. Agarwal, S. Jabbari, C. Agarwal, S. Upadhyay, Z. S. Wu, and H. Lakkaraju. Towards the unification and robustness of perturbation and gradient based explanations. 2021.

[8] U. Aivodji, H. Arai, O. Fortineau, S. Gambs, S. Hara, and A. Tapp. Fairwashing: the risk of rationalization. In *International Conference on Machine Learning*, pages 161–170, 2019.

[9] D. Alvarez-Melis and T. Jaakkola. On the robustness of interpretability methods. *CoRR*, abs/1806.08049, 2018.

[10] O. Bastani, C. Kim, and H. Bastani. Interpretability via model extraction. *CoRR, abs/1706.09773*, 2017.

[11] U. Bhatt, A. Xiang, S. Sharma, A. Weller, A. Taly, Y. Jia, J. Ghosh, R. Puri, J. M. Moura, and P. Eckersley. Explainable machine learning in deployment. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 648–657, 2020.

[12] J. Bien and R. Tibshirani. Classification by set cover: The prototype vector machine. *CoRR, abs/0908.2284*, 2009.

[13] R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1721–1730, 2015.

[14] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso. Machine learning interpretability: A survey on methods and metrics. *Electronics*, 8(8):832, 2019.

[15] P. Chalasani, J. Chen, A. R. Chowdhury, X. Wu, and S. Jha. Concise explanations of neural networks using adversarial training. In *International Conference on Machine Learning*, pages 1383–1391, 2020.

[16] A.-K. Dombrowski, M. Alber, C. Anders, M. Ackermann, K.-R. Müller, and P. Kessel. Explanations can be manipulated and geometry is to blame. *CoRR, abs/1906.07983*, 2019.

[17] F. Doshi-Velez and B. Kim. Towards a rigorous science of interpretable machine learning. *CoRR, abs/1702.08608*, 2017.

[18] R. Elshawi, M. H. Al-Mallah, and S. Sakr. On the interpretability of machine learning-based model for predicting hypertension. *BMC medical informatics and decision making*, 19(1):146, 2019.

[19] A. Ghorbani, A. Abid, and J. Zou. Interpretation of neural networks is fragile. In *AAAI Conference on Artificial Intelligence*, volume 33, pages 3681–3688, 2019.

[20] L. H. Gilpin, D. Bau, B. Z. Yuan, A. Bajwa, M. Specter, and L. Kagal. Explaining explanations: An overview of interpretability of machine learning. In *2018 IEEE 5th International Conference on data science and advanced analytics (DSAA)*, pages 80–89. IEEE, 2018.

[21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[22] S. R. Hong, J. Hullman, and E. Bertini. Human factors in model interpretability: Industry practices, challenges, and needs. *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW1):1–26, 2020.

[23] S. Hooker, D. Erhan, P.-J. Kindermans, and B. Kim. Evaluating feature importance estimates. 2018.

[24] M. Ibrahim, M. Louie, C. Modarres, and J. Paisley. Global explanations of neural networks: Mapping the landscape of predictions. *CoRR, abs/1902.02384*, 2019.

[25] H. Kaur, H. Nori, S. Jenkins, R. Caruana, H. Wallach, and J. Wortman Vaughan. 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*, pages 1–14, 2020.

[26] B. Kim, C. Rudin, and J. A. Shah. The bayesian case model: A generative approach for case-based reasoning and prototype classification. In *Advances in Neural Information Processing Systems*, pages 1952–1960, 2014.

[27] I. Lage, E. Chen, J. He, M. Narayanan, B. Kim, S. Gershman, and F. Doshi-Velez. An evaluation of the human-interpretability of explanation. *CoRR, abs/1902.00006*, 2019.

[28] H. Lakkaraju and O. Bastani. " 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*, pages 79–85, 2020.

[29] H. Lakkaraju, S. H. Bach, and J. Leskovec. Interpretable decision sets: A joint framework for description and prediction. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1675–1684, 2016.

[30] H. Lakkaraju, E. Kamar, R. Caruana, and J. Leskovec. Faithful and customizable explanations of black box models. In *AAAI Conference on Artificial Intelligence, Ethics, and Society*, pages 131–138, 2019.

[31] B. Letham, C. Rudin, T. McCormick, and D. Madigan. Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model. *Annals of Applied Statistics*, 2015.

[32] A. Levine, S. Singla, and S. Feizi. Certifiably robust interpretation in deep learning. *CoRR, abs/1905.12105*, 2019.

[33] Y. Liu, S. Khandagale, C. White, and W. Neiswanger. Synthetic benchmarks for scientific research in explainable machine learning. 2021.

[34] Y. Lou, R. Caruana, and J. Gehrke. Intelligible models for classification and regression. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 150–158, 2012.

[35] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Neural Information Processing Systems (NIPS)*, pages 4765–4774. Curran Associates, Inc., 2017.

[36] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, pages 4765–4774, 2017.

[37] V. Petsiuk, A. Das, and K. Saenko. Rise: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*, 2018.

[38] F. Poursabzi-Sangdeh, D. G. Goldstein, J. M. Hofman, J. W. Vaughan, and H. Wallach. Manipulating and measuring model interpretability. *CoRR, abs/1802.07810*, 2018.

[39] M. Ribeiro, S. Singh, and C. Guestrin. Anchors: High-precision model-agnostic explanations. In *AAAI Conference on Artificial Intelligence*, pages 1527–1535, 2018.

[40] M. T. Ribeiro, S. Singh, and C. Guestrin. Model-agnostic interpretability of machine learning. In *ICML Workshop on Human Interpretability in Machine Learning (WHI)*, 2016.

[41] M. T. Ribeiro, S. Singh, and C. Guestrin. "why should I trust you?": Explaining the predictions of any classifier. In *Demo at the Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2016.

[42] M. T. Ribeiro, S. Singh, and C. Guestrin. "Why should I trust you?" Explaining the predictions of any classifier. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144, 2016.

[43] C. Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206, 2019.

[44] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.

[45] R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *IEEE International Conference on Computer Vision*, pages 618–626, 2017.

[46] A. Shrikumar, P. Greenside, and A. Kundaje. Learning important features through propagating activation differences. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 3145–3153. PMLR, 06–11 Aug 2017. URL https://proceedings.mlr.press/v70/shrikumar17a.html.

[47] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *International Conference on Learning Representations*, 2014.

[48] D. Slack, S. Hilgard, E. Jia, S. Singh, and H. Lakkaraju. How can we fool LIME and SHAP? adversarial attacks on post hoc explanation methods. In *AAAI/ACM Conference on AI, Ethics, and Society*, pages 180–186, 2020.

[49] D. Slack, A. Hilgard, S. Singh, and H. Lakkaraju. Reliable post hoc explanations: Modeling uncertainty in explainability. *Advances in Neural Information Processing Systems*, 34, 2021.

[50] D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg. Smoothgrad: Removing noise by adding noise. *CoRR*, abs/1706.03825, 2017.

[51] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pages 3319–3328, 2017.

[52] L. S. Whitmore, A. George, and C. M. Hudson. Mapping chemical performance on molecular structures using locally interpretable explanations. *CoRR, abs/1611.07443*, 2016.

[53] X. Zhang, J. Zhao, and Y. LeCun. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28:649–657, 2015.

[54] J. Zhou, A. H. Gandomi, F. Chen, and A. Holzinger. Evaluating the quality of machine learning explanations: A survey on methods and metrics. *Electronics*, 10(5):593, 2021.

# A   Experimental Setup

## A.1   Black Box Models: Training and Performance

For tabular data, we train four models: a logistic regression model, a gradient-boosted tree model (50 estimators), a random forest model (50 estimators), and a densely-connected feed-forward neural network (with 4 hidden layers with relu activation consisting of 50, 100, 100, and 50 neurons, respectively). For the COMPAS dataset, we train the four models based on a 70%-30% train-test split of the dataset, using features to predict COMPAS risk score group. The test accuracies of the four models are 0.84, 0.83, 0.82, and 0.84, respectively. For the German credit dataset, we train the same four models based on a 80%-20% train-test split of the dataset, using features to predict credit risk group. The test accuracies of the four models are 0.74, 0.69, 0.75, and 0.70, respectively.

For text data, we trained a widely-used LSTM-based text classifier, based on 120,000 training samples and 7,600 test samples, to predict the news category of the article from which a sentence was obtained. The model performs with 90.67% accuracy. The architecture comprises of an embedding layer of dimension 300, followed by an LSTM layer of hidden size 256 connected to a four-dimensional output layer.

For image data, we use the pre-trained ResNet-18 model [21] and analyze explanations generated for predictions made to classify images to one of the 1000 classes. This model performs 69.758 % and 89.078 % on Accuracy@1 and Accuracy@5 metrics[*], respectively.

## A.2   Explanation Methods

For tabular data, the perturbation-based explanation methods (LIME and KernelSHAP) were applied to explain all four models while the gradient-based explanation methods (Vanilla Gradients, Integrated Gradients, Gradient*Input, and SmoothGRAD) were applied to explain the logistic regression and neural network models to explain samples from the test set (1,482 samples for the COMPAS dataset and 200 samples for the German Credit dataset). Because gradients are not computed for tree-based models, the gradient-based explanation methods were not applied to the random forest and gradient-boosted tree models. When applying explanation methods with a sample size hyperparameter (LIME, KernelSHAP, Integrated Gradients, SmoothGRAD), we

---

[*] https://pytorch.org/vision/stable/models.html

performed a convergence check and selected the sample size at which an increase in the number of samples does not significantly change the explanations. Change in explanations is measured by the L2 distance of feature attributions at the current versus previous sample size. For both COMPAS and German Credit datasets, we used the following number of perturbations/samples/steps for the following explanation methods: LIME (3,000), Integrated Gradients (1,500), SmoothGRAD (1,500). For the COMPAS dataset, when applying KernelSHAP, since the number of features is small, we used a sample size large enough to cover the entire coalition space ($2^7 = 128$ samples), thereby calculating exact Shapley values. For the German Credit dataset, when applying KernelSHAP, we used 3,000 samples, based on the convergence analysis.

For text data, we applied all six explanation methods on the LSTM-based classifier to explain predictions for 7,600 samples in the test set. For LIME and KernelSHAP, we follow the convergence analysis described above and find that attributions do not change significantly beyond 500 perturbations; hence, we use 500 perturbations for LIME and KernelSHAP. Integrated Gradients explanations were generated using 500 steps which is higher than the recommended number of steps mentioned in [51]. SmoothGRAD explanations were generated using 500 samples to get the most confident attribution which is significantly higher the recommended number of 50 samples [50].

For image data, we applied all six explanation methods on the ResNet-18 model [21] to explain predictions for the PASCAL VOC 2012 test set of 1,449 samples. Integrated Gradients explanations were generated using 400 steps, significantly higher than the recommendation of 300 [51], to obtain a stable and confident attribution map. Similarly, SmoothGRAD explanations were generated using a sample size of 200 which is also higher than the recommended sample size of 50 [50]. For LIME and KernelSHAP, we chose 100 perturbations to train the surrogate model as we did not notice any significant changes in attributions beyond 50 perturbations. KernelSHAP and LIME were used to compute attributions of super-pixels annotated in PASCAL VOC 2012 segmentation maps. Due to a larger feature space in images compared to the previous tabular and text datasets, disagreement metrics based on top-$k$ features may not provide a clear picture. Hence, we use Rank Correlation and cosine distance between attribution maps generated by a pair of explanation methods as the disagreement metric. Higher cosine distance between attribution maps indicate larger disagreement between explanation methods.

# B    Results from Empirical Analysis of Disagreement Problem

## B.1    COMPAS Dataset

### B.1.1    Figure description: metrics measuring agreement among a set of selected features

Disagreement of explanation methods as measured by rank correlation (left column) and pairwise rank agreement (right column) over test set data points. Both metrics are calculated across all features. Heatmaps show the average metric value and boxplots show the distribution of metric values for each pair of explanation methods. In heatmaps, lighter colors indicate stronger disagreement. Minimum and maximum standard errors are indicated below each heatmap.

### B.1.2    Figure description: metrics measuring agreement among top-$k$ features

Disagreement of explanation methods as measured by rank agreement, feature agreement, sign agreement, and signed rank agreement (each row is one metric). By definition, when $k$ equals the full set of features, feature agreement equals one. Heatmaps show the average metric value for each pair of explanation methods, with lighter colors indicating stronger disagreement. Minimum and maximum standard errors are indicated below each heatmap.

**Neural Network**



Figure 7: Disagreement of explanation methods for neural network model trained on COMPAS dataset. Figure description in Appendix B.1.1.

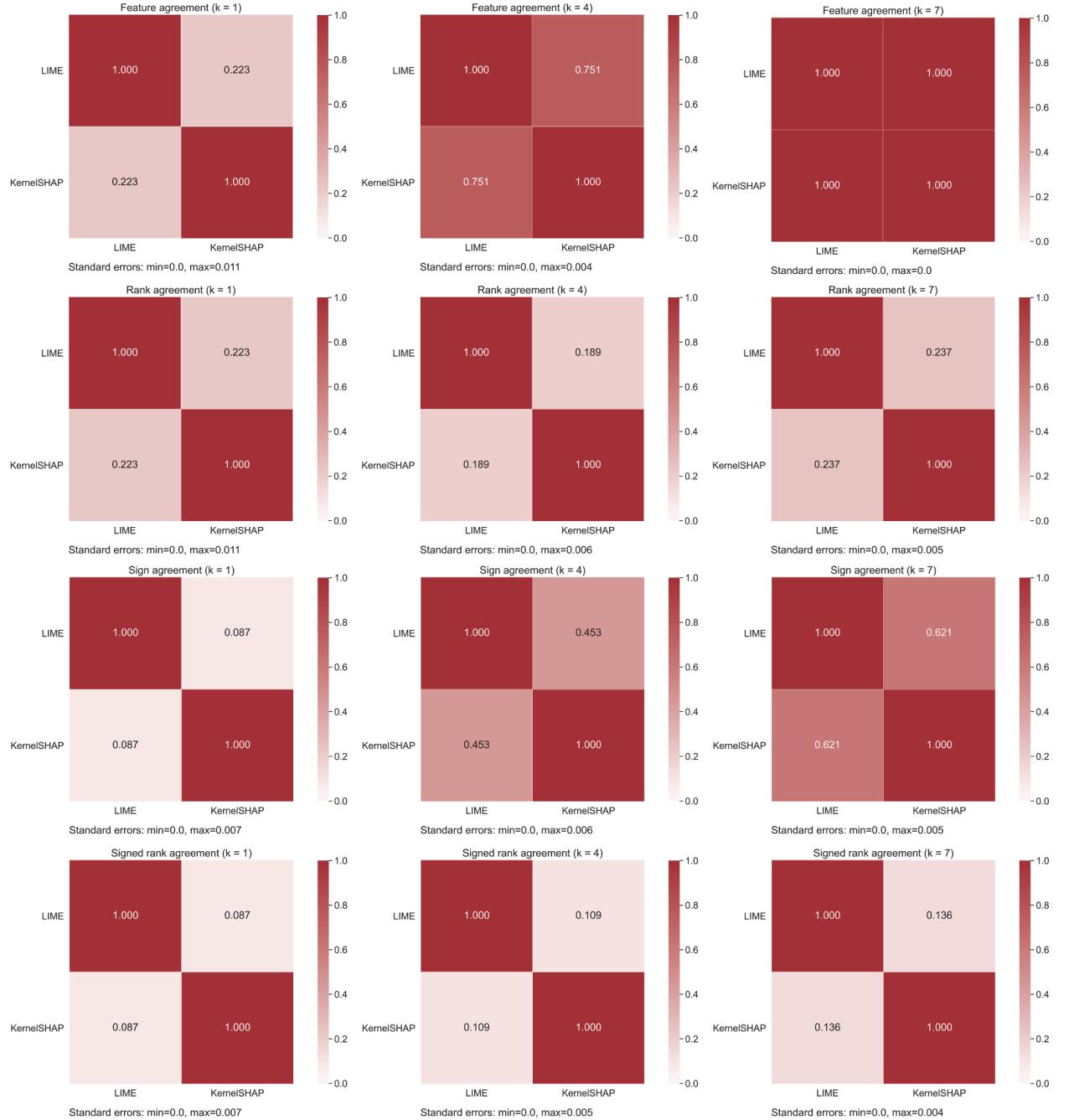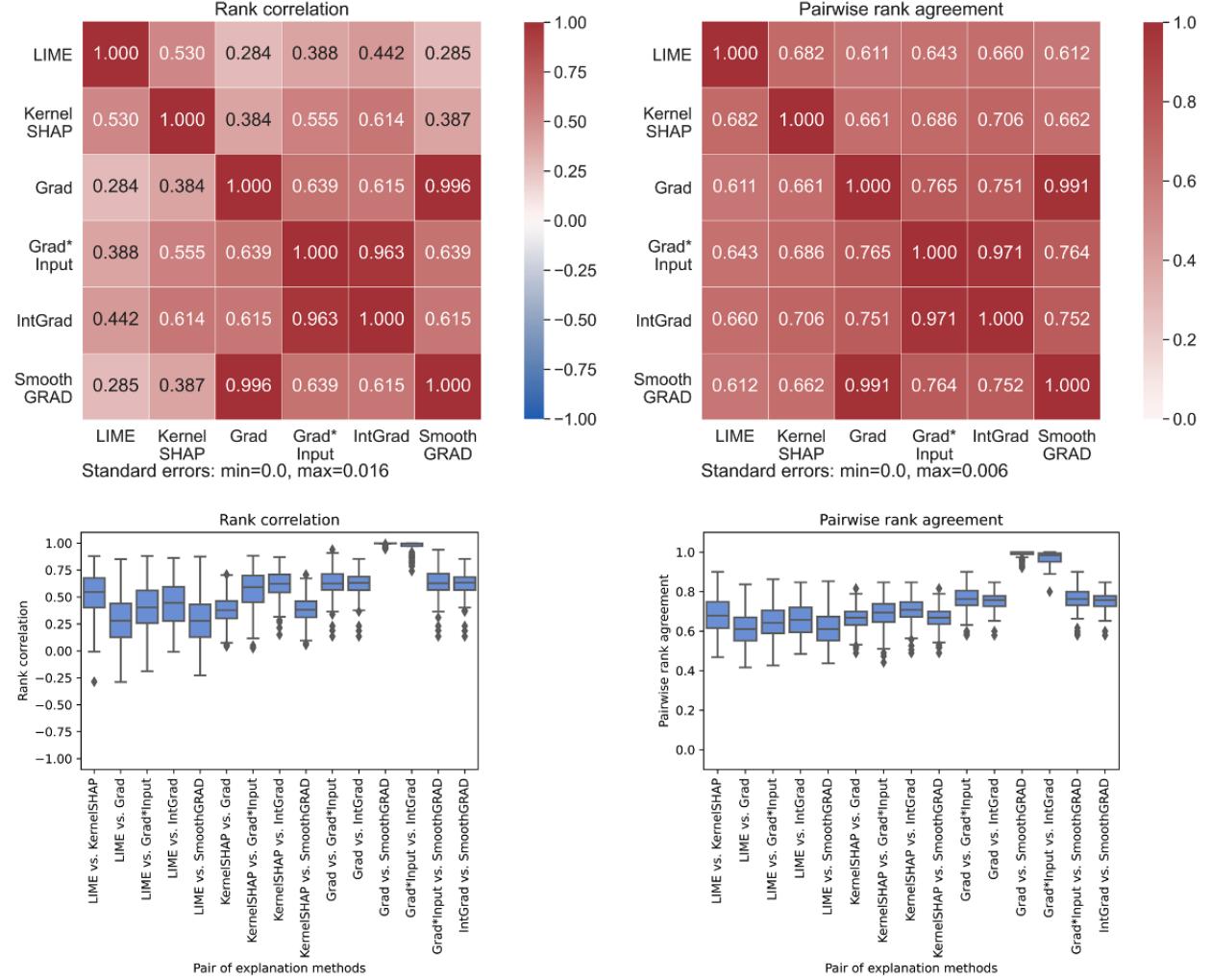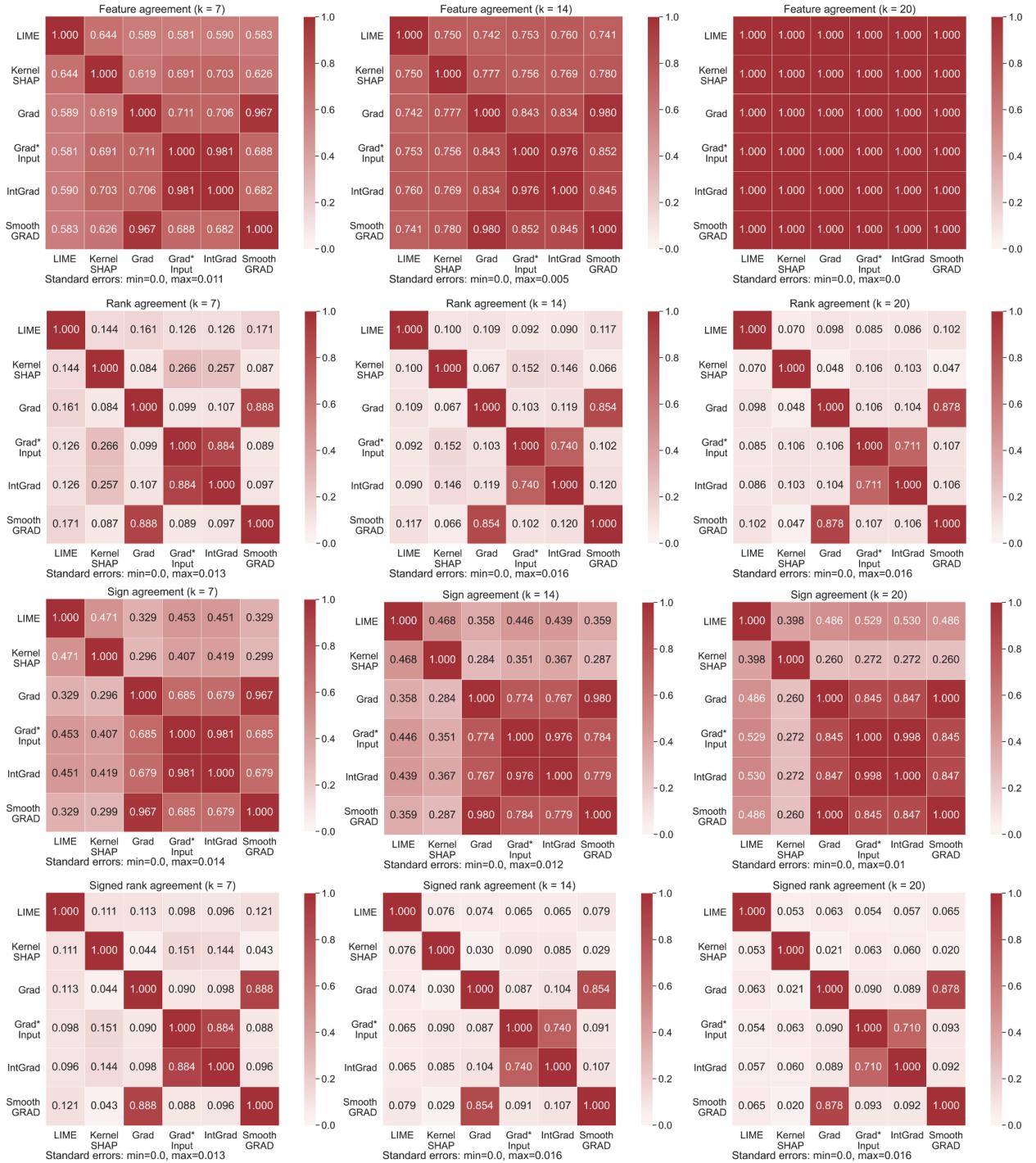Figure 8: Disagreement of explanation methods for neural network model trained on COMPAS dataset. Figure description in Appendix B.1.2.
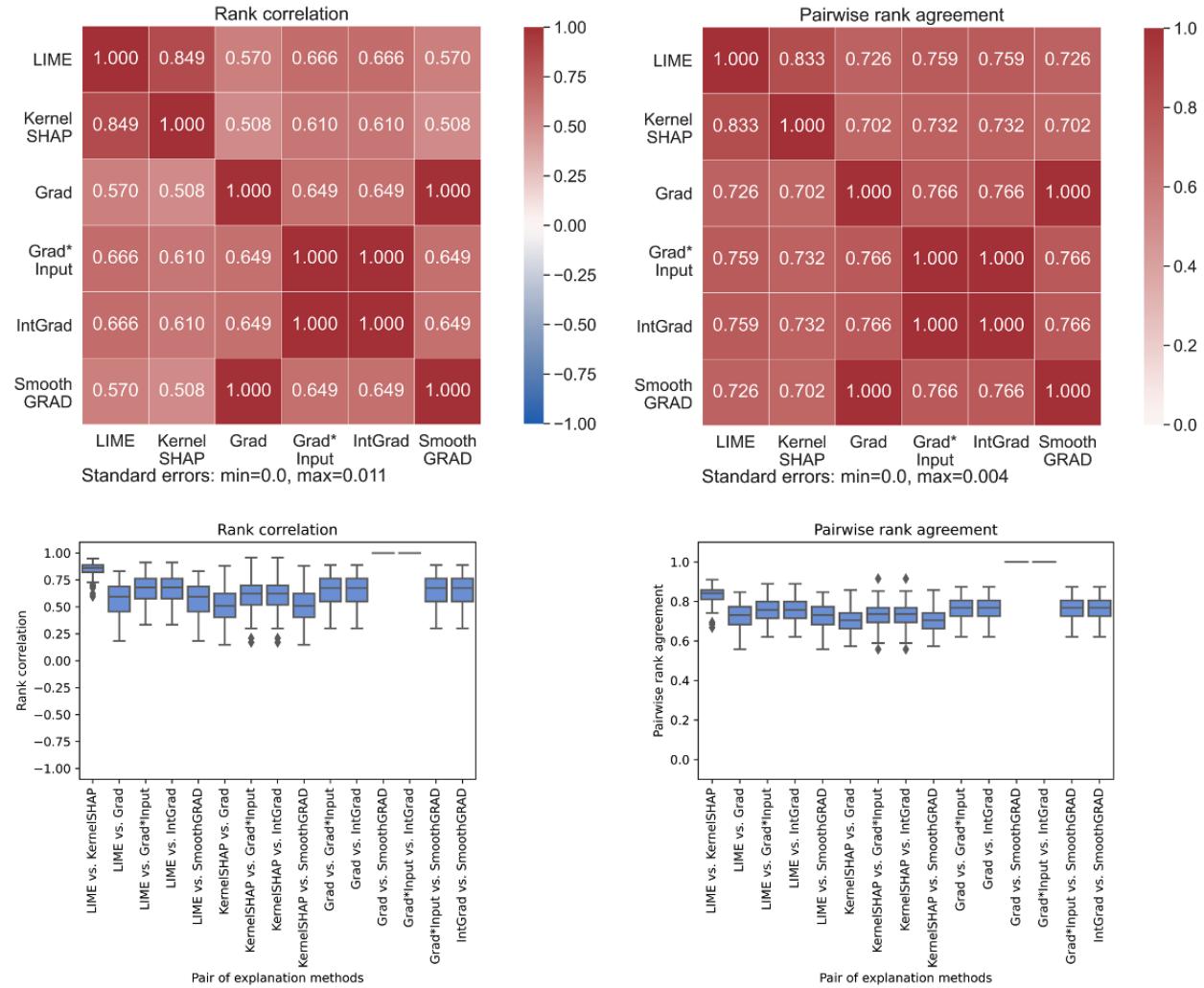
**Logistic Regression**



Figure 9: Disagreement of explanation methods for logistic regression model trained on COMPAS dataset. Figure description in Appendix B.1.1.

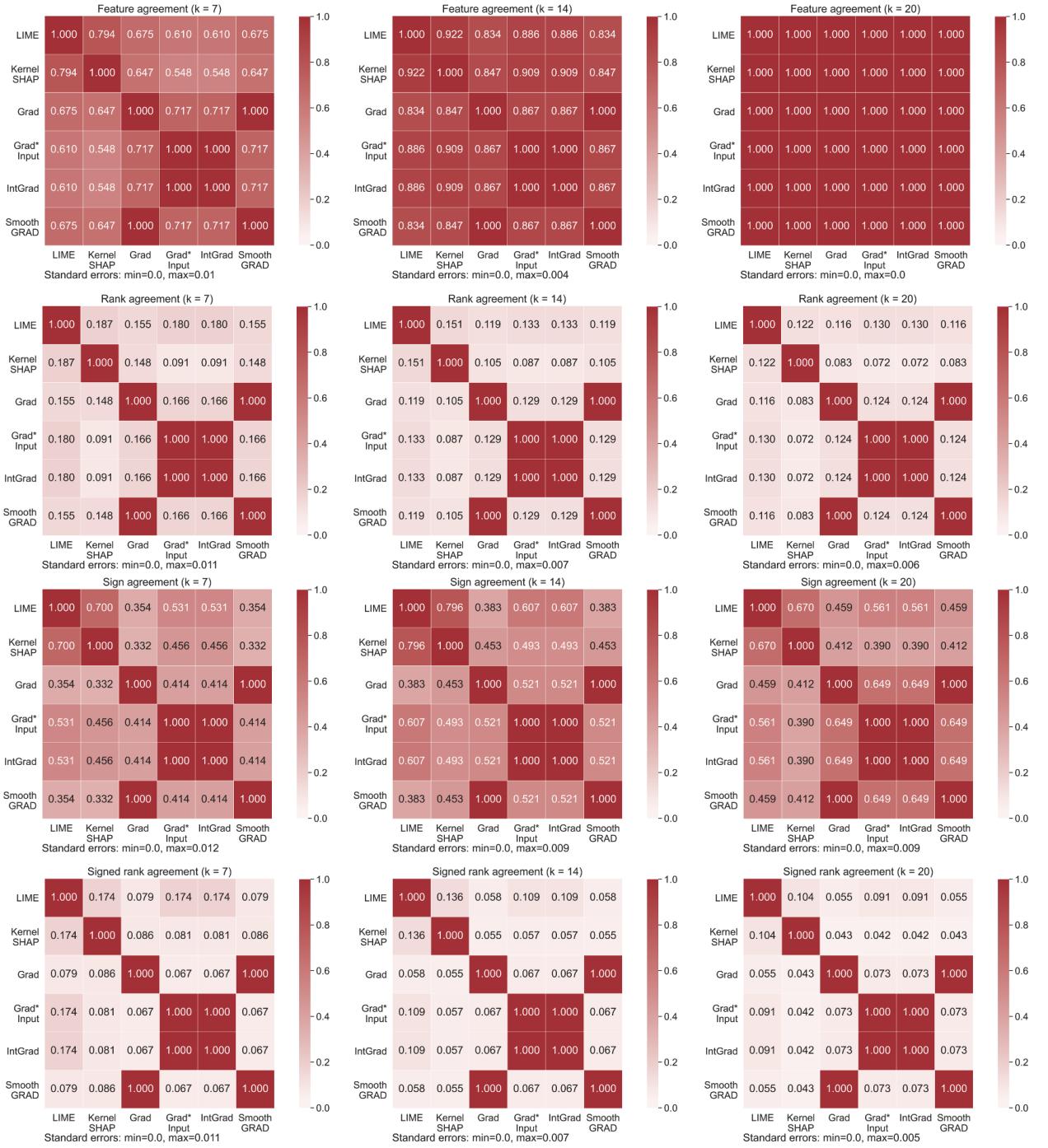Figure 10: Disagreement of explanation methods for logistic regression model trained on COMPAS dataset. Figure description in Appendix B.1.2.

**Random Forest**


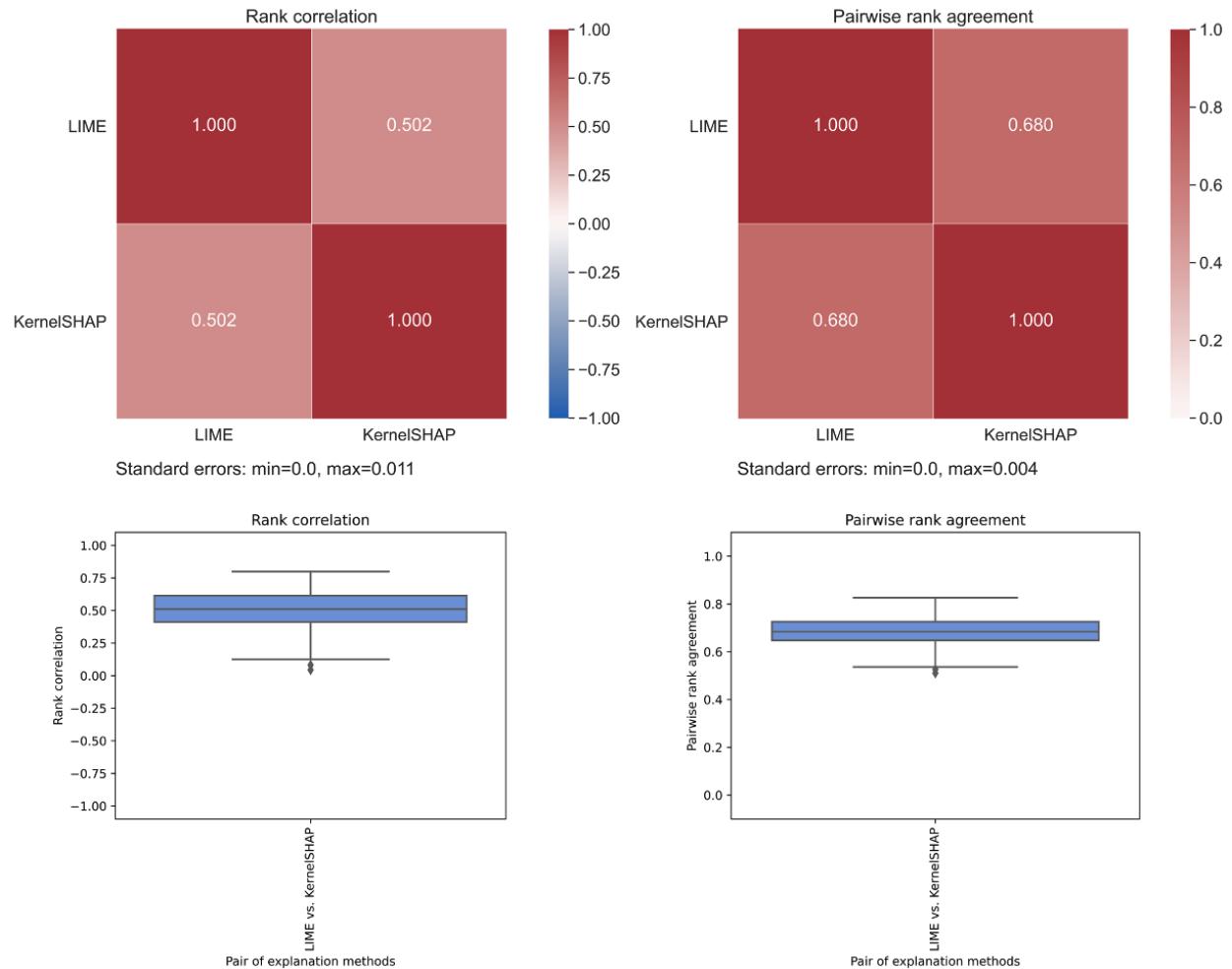
Figure 11: Disagreement of explanation methods for random forest model trained on COMPAS dataset. Figure description in Appendix B.1.1.
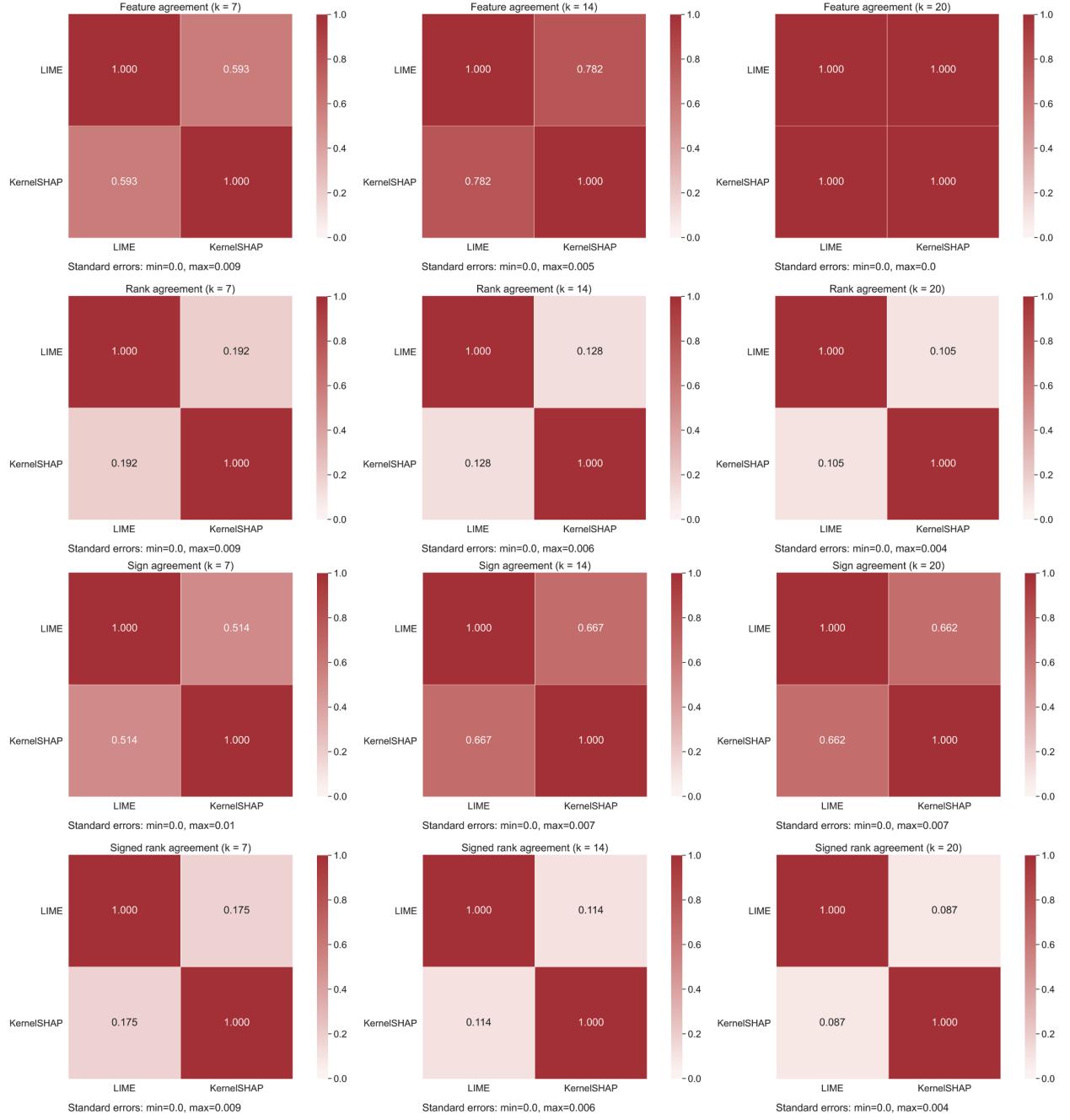
Figure 12: Disagreement of explanation methods for random forest model trained on COMPAS dataset. Figure description in Appendix B.1.2.

**Gradient-Boosted Tree**



Figure 13: Disagreement of explanation methods for gradient-boosted tree model trained on COMPAS dataset. Figure description in Appendix B.1.1.

Figure 14: Disagreement of explanation methods for gradient-boosted tree model trained on COMPAS dataset. Figure description in Appendix B.1.2.

## B.2 German Credit Dataset

**Neural Network**


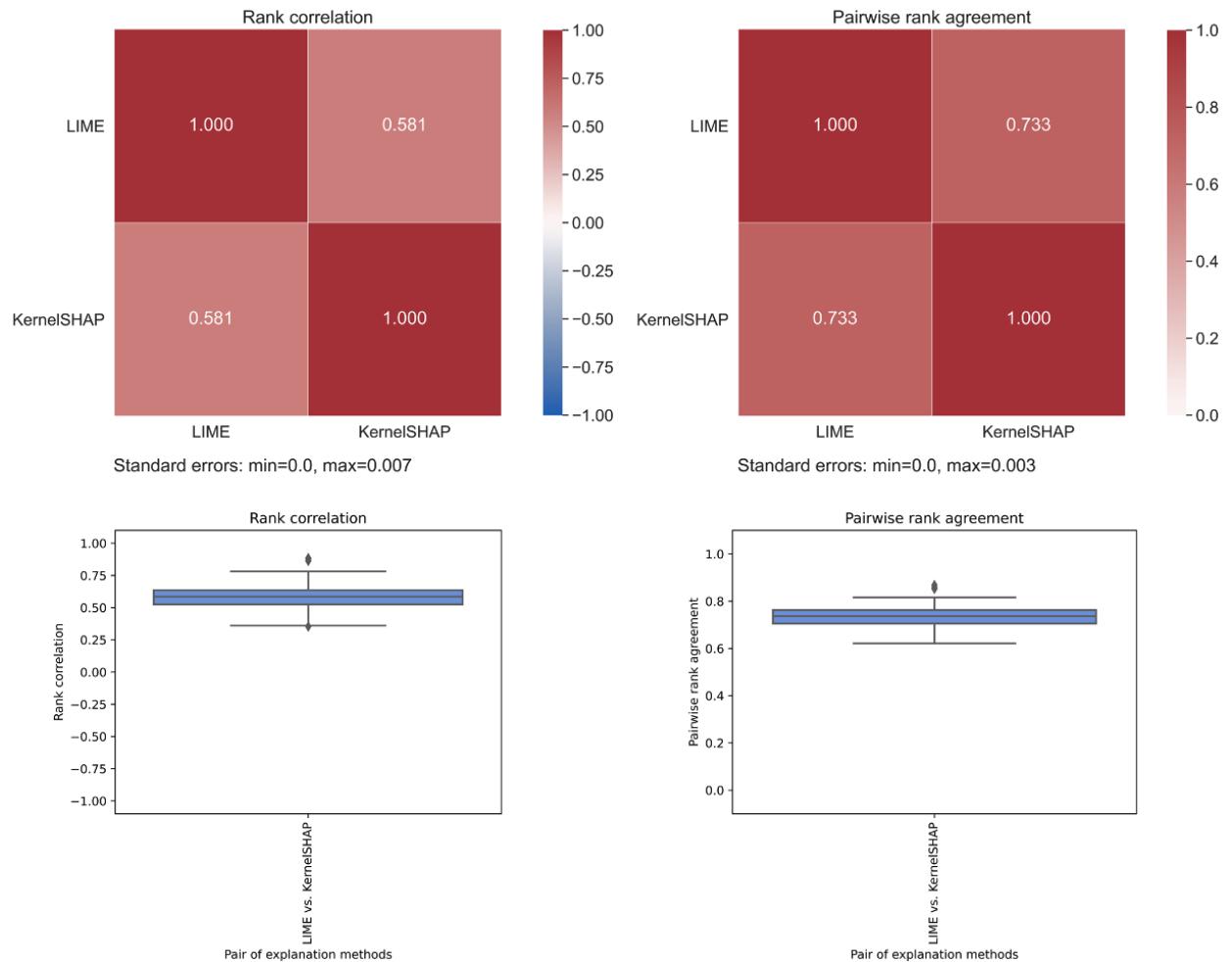
Figure 15: Disagreement of explanation methods for neural network model trained on German Credit dataset. Figure description in Appendix B.1.1.

Figure 16: Disagreement of explanation methods for neural network model trained on German Credit dataset. Figure description in Appendix B.1.2.

**Logistic Regression**



Figure 17: Disagreement of explanation methods for logistic regression model trained on German Credit dataset. Figure description in Appendix B.1.1.
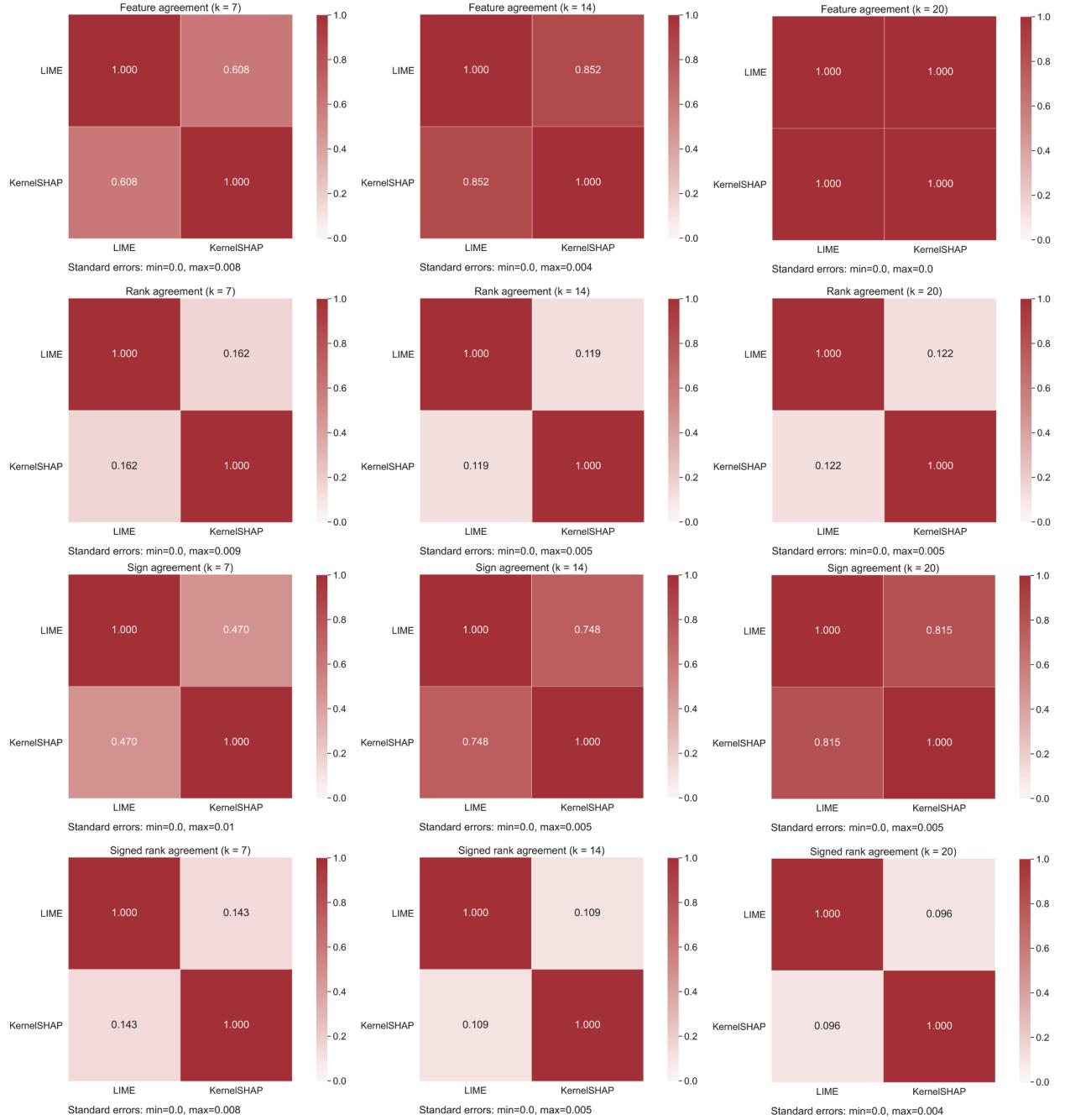
Figure 18: Disagreement of explanation methods for logistic regression model trained on German Credit dataset. Figure description in Appendix B.1.2.

**Random Forest**



Figure 19: Disagreement of explanation methods for random forest model trained on German Credit dataset. Figure description in Appendix B.1.1.

Figure 20: Disagreement of explanation methods for random forest model trained on German Credit dataset. Figure description in Appendix B.1.2.

**Gradient-Boosted Tree**



Figure 21: Disagreement of explanation methods for gradient-boosted tree model trained on German Credit dataset. Figure description in Appendix B.1.1.
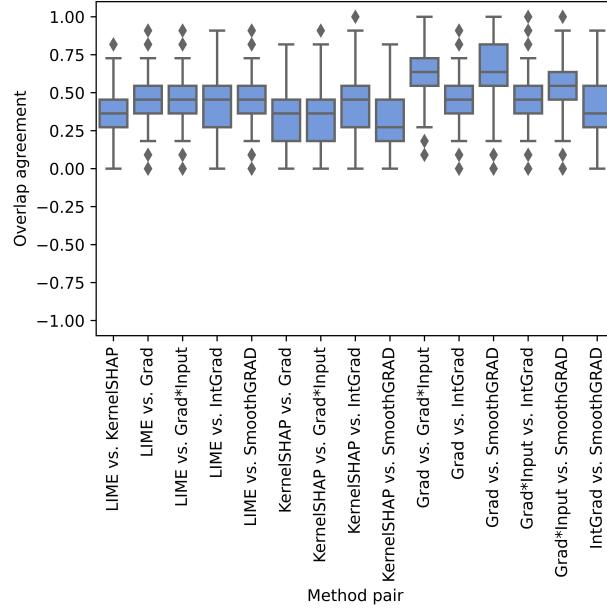
Figure 22: Disagreement of explanation methods for gradient-boosted tree model trained on German Credit dataset. Figure description in Appendix B.1.2.

## B.3 AG_News Dataset



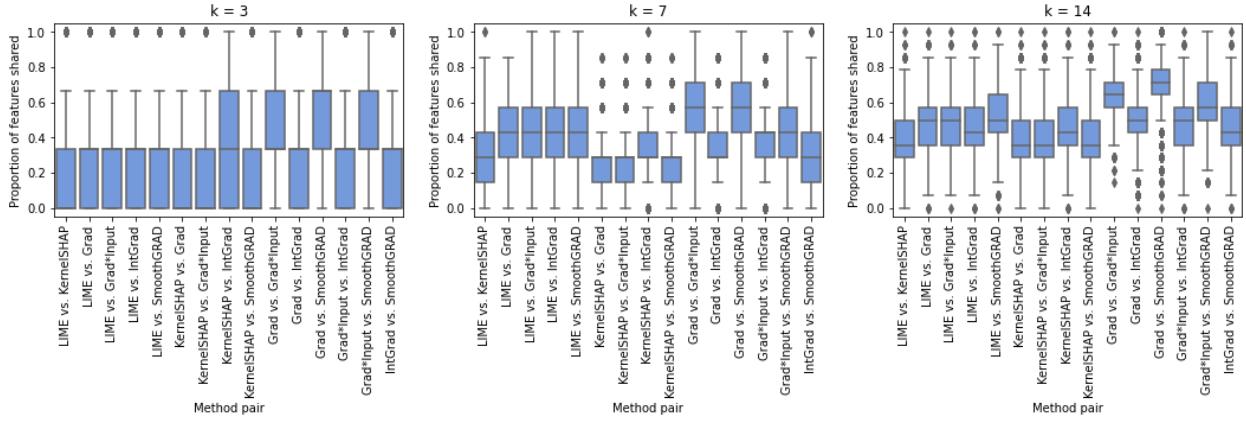Figure 23: Box plot for feature agreement for pair of explanations for AG_News



Figure 24: Box plot for feature agreement for pair of explanations for AG_News

## B.4 ImageNet Dataset

| Metrics | ResNet-18 |
| --- | --- |
| **Rank correlation** | 0.8977 |
| **Pairwise rank agreement** | 0.9302 |
| **Feature agreement** | 0.9535 |
| **Rank agreement** | 0.8478 |
| **Sign agreement** | 0.9218 |
| **Signed rank agreement** | 0.8193 |

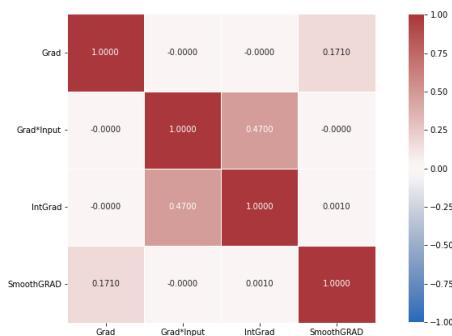Table 2: Disagreement in ImageNet between LIME and KernelSHAP



Figure 25: Rank Correlation for explanations computed at pixel level for gradient based explanation methods

# C  Omitted Details from Section 5

## C.1  Screenshots of UI

In Figures 26 and 27, we present screenshots of the UI that participants are presented with before beginning the study. The purpose of this introduction page is to familiarize the participants with the COMPAS prediction setting, the six explainability methods we use, and the explainability plots we show in each of the prompts.

## Introduction

COMPAS is a popular commercial algorithm used by judges for determining a criminal defendant's likelihood of reoffending (recidivism).

The COMPAS dataset consists of 7 features:

- **age**
- **two_year_recid**: whether the defendant recidivated within 2 years of the original crime
- **priors_count**: number of prior crimes committed
- **length_of_stay**: length the defendant stayed in jail
- **c_charge_degree**: one of Misdemeanor, Felony
- **sex**: one of Male, Female
- **race**: one of African-American, Asian, Caucasian, Hispanic, Native American, or Other

For this study, we trained a neural network on the COMPAS dataset to **predict a criminal defendant's COMPAS risk score (low or high), corresponding to whether he/she would commit a crime after two years past the date of the original crime**. Since it is important to understand our model's predictions (explainability), we also ran six popular explainability algorithms on various input points.
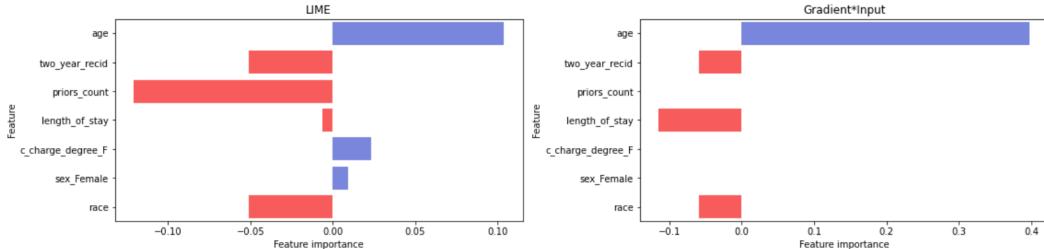
The explainability algorithms we use are listed here. **You do not need to understand them past what we have described here.**

- **LIME**: an explanation based on a locally linear approximation of the model at that input
- **KernelSHAP**: a combination of LIME and Shapley Values, which identify the contribution of each feature based on interactions with other features
- **Gradient**: the gradient of the model at the input
- **Gradient*Input**: the dot product of the input features and the gradient explanation
- **SmoothGrad**: weighted average of the gradient at points around the input
- **Integrated Gradients**: a modification of the gradient method to satisfy two axioms, *sensitivity* and *implementation invariance*

Figure 26: This is a screenshot of the first half of the introductory page, describing our COMPAS risk score prediction setting and briefly summarizing the six explainability algorithms used (with links to their corresponding papers for the interested participant).

## Your Task

On each of the next 5 pages, you will see the result of two explainability methods on the same input sample from COMPAS, as shown below. **Assume that the criminal defendant's risk of recidivism was correctly predicted to be high.** The explanation of the prediction will then be shown to you, as in the figure below.



The y-axis lists each of the 7 COMPAS features, and the x-axis shows the importance of that feature. Positive importance values are shown in blue, while negative importance values are shown in red. A **high positive importance** for a feature means that the feature contributed greatly to the correct prediction, while a **high negative importance** means that the feature negatively contributed (was misleading) to the prediction. Note the different x-axis scales resulting from different methods. You will be asked to compare the two explanations.

Figure 27: This is a screenshot of the second half of the introductory page, describing the concrete task and an explanation of what is shown in the explainability plots.

## C.2 Prompts Used

In this section, we share the 15 prompts that we showed users. Each prompt highlights a pair of different explainability algorithms on a COMPAS data point. For each pair, we chose the data point from the entire COMPAS set that maximized the rank correlation between the explanations.

## C.3 User Study Questions

In each of the five prompts, we asked participants the following questions, which we refer to as *Set 1*. Questions 3-4 were only shown if the user selected *Mostly agree*, *Mostly disagree*, or *Completely disagree* to Question (1).

1. To what extent do you think the two explanations shown above agree or disagree with each other? (choice between *Completely agree, Mostly agree, Mostly disagree, Completely disagree*)

2. Please explain why you chose the above answer.

3. Since you believe that the above explanations disagree (to some extent), which explanation would you rely on? (choice between *Algorithm 1 explanation, Algorithm 2 explanation, It depends*)

4. Please explain why you chose the above answer.

After answering all five prompts, the user was then asked the following set of questions, which we refer to as *Set 2*. Questions 4-9 were only shown if the user selected *Yes* to Question 3.
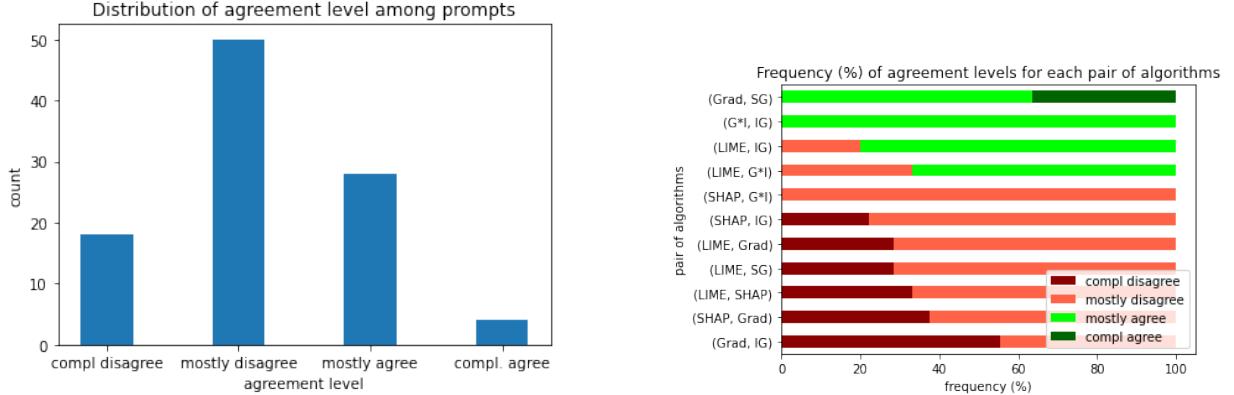
1. (Optional) What is your name?

2. What is your occupation? (eg: PhD student, software engineer, etc.)

3. Have you used explainability methods in your work before? (*Yes/No*)

4. What do you use explainability methods for?

5. Which data modalities do you run explainability algorithms on in your day to day workflow? (eg: tabular data, images, language, audio, etc.)

6. Which explainability methods do you use in your day to day workflow? (eg: LIME, KernelSHAP, SmoothGrad, etc.)

7. Which methods do you prefer, and why?

8. Do you observe disagreements between explanations output by state of the art methods in your day to day workflow?

9. How do you resolve such disagreements in your day to day workflow?

## C.4 Further analysis of overall agreement levels

In this section, we present further plots analyzing responses to questions (1) in Set 1. As shown in Figure 29a, only 32% of responses were *Mostly Agree/Completely Agree* and 68% were *Mostly Disagree/Completely Disagree*, indicating that participants experienced the disagreement problem. We also grouped the responses by prompt, shown in Figure 29b, highlighting that different pairs of algorithms can have different levels of disagreement. We removed prompts with less than 4 total responses. We see that there are varying levels of disagreements among prompts. For example, all participants who were shown the Gradient vs. SmoothGrad prompt believed they agreed to some extent, while all participants who were shown the Gradient vs. Integrated Gradients prompt believed they disagreed to some extent.

Figure 28: Images showing the 15 prompts we used. Each prompt shows the explanation of the same input point with two different interpretability algorithms.

(a) This figure shows the distribution of responses in aggregation over all prompts. The x-axis shows the four possible responses, and the y-axis shows the number of times that response was chosen. Observe that in 68% of cases, participants indicated that the prompts mostly or completely disagreed.

(b) This figure shows the distribution of responses, sorted by prompt. The y-axis shows the pair of explainability algorithms shown in the prompt, and the x-axis shows the frequency that each response was chosen.

Figure 29: These figures show the distribution of answers to Question (1) in Set (1) from Section C.4 in aggregation over all participants.

## C.5 Further analysis of reasons participants chose specific algorithms

In this section, we analyze the responses to Set 1, Question (3) in Section C.3. We saw, in 5.2.2, that algorithms such as KernelSHAP were favored over other algorithms. In Table 3, we list the top reasons the four most frequently chosen algorithms were preferred, showcasing direct quotes from participants.

## C.6 Analysis of reasons participants chose neither algorithm

In this section, we analyze the responses to Set 1, Question (4) in Section C.3, focusing on when participants selected *"It depends"* in Question (3), which was chosen in 38% of cases. Again, we present an overarching summary of the reasons participants made this decision in Table 4.

## C.7 Further analysis of concluding questionnaire

In this section, we extend the analysis presented in 5.2.3, analyzing the responses to questions in Set 2 of Section C.3. As stated in 5.2.3, we received a total of 20 positive responses to Question (3), but one declined to answer Questions (4) through (9). Therefore, we analyze the remaining 19 responses.

In Question (4), we found that study participants use explainability methods for a variety of reasons such as understanding models, debugging models, help explain models to clients, research. In Question (5), we found that 16 of 19 participants employed explanations for tabular data, 6 of 19 participants for text and language data, 11 of 19 participants for image data, and 1 of 19 for audio data. In Question (6), we found that 14 of 19 participants used LIME, 14 of 19 participants used SHAP, and 13 of 19 participants used some sort of gradient-based methods. Participants also indicated using methods like GradCAM, dimensionality reduction, MAPLE, and rule-based methods. In Question (7), 9 of 19 participants stated that they preferred both LIME and SHAP, with another 3 of 19 participants stating LIME only. We showcase some intriguing answers from Question (7) below:

- "LIME and SHAP seem to be the most universally applicable and I can understand."

- "Methods with underlying theoretical justifications such as KernelSHAP and Integrated Gradients"

Table 3: Reasons participants chose the top four most favored explainability algorithms (KernelSHAP, SmoothGrad, LIME, and Integrated Gradients) over others when explanations disagreed.

| Algorithm | Reasons that algorithm was chosen in disagreement |
|---|---|
| **KernelSHAP** | • [36%] SHAP is better for tabular data (*"SHAP is more commonly used [than Gradient] for tabular data"*)<br>• [25%] SHAP is more familiar (*"More information present + more familiarity"*)<br>• [14%] SHAP is a better algorithm overall (*"SHAP seems more methodical than LIME"*, *"SHAP is a more rigorous approach [than LIME] in theory"*) |
| **SmoothGrad** | • [33%] SmoothGrad paper is newer or better (*"SmoothGrad is apparently more robust"*, *"SmoothGrad is often considered improved verison of grad"*)<br>• [58%] Reasons based on the explainability map shown (*"directionality of the attributions … [agree] with intuition"*, *"gradient has unstability problems [, so] smoothgrad"*) |
| **LIME** | • [54%] LIME is better for tabular data (*"I use LIME for structured data."*)<br>• [15%] LIME is more familiar/easier to interpret (*"I am more familiar with LIME"*, *"LIME is easy to interpret"*) |
| **Integrated Gradients** | • [86%] Integrated Gradients paper is better (*"IG came after gradients and paper shows improvements"*, *"integrated gradients paper showed improvements [over Gradient × Input]"* |

Table 4: Reasons people answered *"It depends"* after being asked to choose between disagreements

| Rationale | Representative Quote |
|---|---|
| **1. Need more information** | • *"need to see the final prediction of the model and the feature values"* |
| **2. Pick neither explanation** | • *"No compelling reason to choose one over the other. Both don't align with intuition."* |
| **3. Unsure/Don't know** | • *"I'm not sure which of the two methods is more trustworthy"* |
| **4. Would consult an expert** | • *"I would ask a domain expert for his/her opinion"* |
| **5. Combine explanations** | • *"I would combine both – note that age might be doing weird things, but that length of stay and race both contribute to a negative prediction"* |
| **6. Depends on use case** | • *"The two methods have different interpretations - it depends on if I'm more interested in comparing my explanation to some baseline individual state versus just interested in understanding the immediate local behavior"* |

- "lime and shap ... [easy to implement] and can work with black box"

- "shap and lime because ... [they are] easy to understand and have standard implementations"

- "LIME, because everything else isn't necessarily capturing what I actually want to know about the local behavior"

Finally, we provide additional quotes highlighting the responses to Questions (8) and (9), which were briefly analyzed in Section 5.2.3. These are shown in Table 5.

Table 5: Representative quotes highlighting themes of how participants address the disagreement problem in their day to day work

| Category of Response | Samples Quotes |
|---|---|
| **1. Make arbitrary decisions (50%).** | • *"Such disagreements are resolved by data scientists picking their favorite algorithm"* <br> • *"I try to use rules of thumb based on results in research papers and/or easy to understand outputs."* <br> • *"I favor lime and shap because there is well documented packages on github"* |
| **2. Unsure/Don't know/Don't resolve (36%)** | • *"there is no clear answer to me. I hope research community can provide some guidance"* <br> • *"unfortunately there is no good answer at my end ... I hope you can help me with finding an answer"* |
| **3. Use other metrics (fidelity) (14%).** | • *"By quantitative assessment of feature importance methods that assess specific properties like faithfulness"* <br> • *"I might try and use some metric to measure fidelity."* |