

An Empirical Evaluation of the Rashomon Effect in Explainable Machine Learning

Sebastian Müller (✉)^{1,4}[0000-0002-0778-9695], Vanessa
Toborek^{1,4}[0009-0009-8372-8251], Katharina Beckh^{3,4}[0000-0002-7824-6647],
Matthias Jakobs^{2,4}[0000-0003-4607-8957], Christian
Bauckhage^{1,3,4}[0000-0001-6615-2128], and Pascal Welke⁵[0000-0002-2123-3781]

¹ University of Bonn, Bonn, Germany

² TU Dortmund University, Dortmund, Germany

³ Fraunhofer IAIS, Sankt Augustin, Germany

⁴ Lamarr Institute, Germany

⁵ TU Wien, Vienna, Austria

Abstract. The Rashomon Effect describes the following phenomenon: for a given dataset there may exist many models with equally good performance but with different solution strategies. The Rashomon Effect has implications for Explainable Machine Learning, especially for the comparability of explanations. We provide a unified view on three different comparison scenarios and conduct a quantitative evaluation across different datasets, models, attribution methods, and metrics. We find that hyperparameter-tuning plays a role and that metric selection matters. Our results provide empirical support for previously anecdotal evidence and exhibit challenges for both scientists and practitioners.

Keywords: Explainable ML · Interpretable ML · Attribution Methods
· Rashomon Effect · Disagreement Problem

1 Introduction

We demonstrate the impact of the Rashomon Effect when analyzing ML models. The Rashomon Effect [8] describes the phenomenon that there may exist many models within a hypothesis class which solve a dataset equally well. The set of these models is referred to as the Rashomon Set [12, 37]. From a data-centric perspective this phenomenon is also called Predictive Multiplicity [23], meaning that there exist many strategies to solve a task on a dataset. Other works use Rashomon Sets to analyze and describe data [12, 30]. Somewhat surprisingly, the Rashomon Effect has not yet found wider attention in the Explainable Machine Learning (XML) literature. Although a few works have observed the effect it was only anecdotally or without referring to its proper name [14, 20, 35].

Accepted for presentation at European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML/PKDD 2023)

XML has recently become a very active area of research and numerous explanation methods exist [1, 9, 24]. Many approaches explain black-box models in a post-hoc manner by providing attribution scores [22, 27] which assign each input dimension a numerical value that represents this feature’s importance with respect to the model decision. Attribution scores are used to answer questions such as “What feature was the most important in this input sample?” and have been used to uncover spurious correlations in the data [29] and biased behavior of models [21]. However, attribution scores are sometimes ambiguous and their interpretation depends on the application context. It is hard to decide at what magnitude a feature is still important, particularly, if magnitudes of attribution scores can be sorted into an evenly descending order. It follows that the task of comparing different attribution methods is a difficult problem. Several works touch upon the problem of explanation comparison [5, 7, 19, 26, 35] from different perspectives.

Our main contribution is an empirical analysis of one novel and two existing perspectives, 1) demonstrating model-specific sensitivity regarding the hyperparameter choice for explanation methods, 2) comparison of different explanations from the same attribution method on differently initialized but otherwise identical model architectures [5, 35] and 3) the disagreement between different explanations applied to the same architecture and parameterization [19, 26]. We place these three perspectives into a unified framework to investigate how the Rashomon Effect manifests itself in each situation. Our evaluation is conducted on four datasets of entirely different nature, analyzing differences in models explained by five popular attribution methods using both naive and established human-centered similarity measures.¹ Our results highlight the need to fine-tune the hyperparameters of XML methods on a per-model basis. We do find empirical support for the disagreement problem, meaning practitioners cannot expect consistent explanations across methods. Further, the high solution diversity across models hinders the use of XML as an epistemic tool.

Next, Section 2 discusses how we connect different parts of the literature for our analysis. Section 3 describes the experimental setup in detail. Sections 4.1, 4.2 and 4.3 present results and discuss the three perspectives we analyze. Section 4.4 summarizes our main findings. Section 5 concludes.

2 Comparing Attribution Scores

Given a classifier and a datum, an attribution method assigns each input dimension a numerical value that represents this feature’s importance with respect to the model decision. Hence, an attribution scoring depends on three variables: 1) the model, 2) the input sample, and 3) the attribution method. This distinction enables us to systematically investigate the consequences of the Rashomon Effect on established and novel perspectives in XML in one framework. This framework is the first to bring the different perspectives into a unifying picture

¹ Our code is available at github.com/lamarr-xai-group/RashomonEffect

Table 1: We investigate the Rashomon Effect in Explainable Machine Learning for a set of models and a set of attribution methods. Three interesting scenarios arise for a fixed input-sample from a given dataset.

Same Model	Same Sample	Same Attr Method	Scenario	Examples
1	1	1	Numerical Stability	–
0	1	1	Solution Diversity	[12] [17] [35]
1	1	0	Disagreement Problem	[11] [19] [26]

which we present in Table 1. Our main mode of comparison is centered around comparing pairs of attribution scores. Hence, we assume that the scores belong to the same sample from the same dataset. We investigate model- or attribution method-dependent effects and do not consider the scenario where the data is the same but both models and methods are different.

Numerical Stability (111): In Section 4.1 we discuss the scenario where the same model and same attribution method are applied to the same sample. This perspective is relevant to non-deterministic explanation methods that can be controlled by hyperparameters. We investigate whether there are model-specific differences regarding optimal parameter choice and find that the hyperparameter choice is significantly dependent on both the investigated model and dataset. This suggests that blindly applying non-optimal hyperparameters can lead to erroneous explanations and thus wrong takeaways in an application scenario. This need for rigorous hyperparameter tuning is mostly overlooked in the literature.

Solution Diversity (011): We can compare how similar or dissimilar two models are w.r.t. their solution strategies by comparing explanations that were computed for each of them using the same attribution method. Comparing any two models not only by one, but by the average difference on explanations over several samples, will only be able to measure a difference, if two models consistently behave differently. This is a coarse, but sufficiently sensitive measure. Using this measure as a basis, we provide a large quantitative view of the Rashomon Effect itself, recently also observed in [35]. In Section 4.2, we extend existing results by comparing substantially more models on additional data domains and investigate how diverse the strategies of the models within a Rashomon Set are. We observe very high diversity in most cases and discuss practical implications for machine learning (ML) as an epistemic tool [28, 39].

Disagreement Problem (110): Aiming to find the “right” explanation, prior work compared different attribution methods applied to the same model on the same sample. It was found that explanations of different attribution methods often differ significantly, which is now known as the Disagreement Problem [5, 11, 13, 15, 16, 19, 26]. So far, the Disagreement Problem was only reported on individual or a very small number of models. It has not been sufficiently explored whether the disagreement actually is model-dependent, i.e., whether any pair of attribution methods is consistently less similar than other pairs across models.

We investigate this question in Section 4.3. We provide quantitative support for anecdotal observations from the literature and add practically relevant insights.

A fundamental question is which metric should be used to compare two attribution scores. One possible approach is to use feature (dis-)agreement, i.e. the overlap of the top-k “most important” features, which ML practitioners indicated as a key measure for disagreement [19]. Along the same lines, ranking correlation measures are used, such as Kendall’s τ [26]. Another option is to base the comparison on typical distance measures, such as cosine similarity [7] or Euclidean distance [11,25]. In previous studies, only one metric or metric type has been considered. In this work, we provide a comparison of both Euclidean and (dis-)agreement based measures.

3 Experimental Framework

Before we report our results we introduce the experimental setup. To emphasize the extent of the Rashomon Effect we will remove randomness from the training process with the exception of model initialization.

3.1 Datasets

For the comparison we chose four publicly available datasets. *AG News* [40], a benchmark dataset for text classification with an average sentence length of 43 words. Three tabular datasets containing only real valued variables: *Dry Bean* [18], a 16-dimensional multi-class dataset with 7 classes of dry beans, *Breast Cancer Wisconsin (Diagnostic)* [36], a classical dataset posing a binary classification problem over 30 features, and *Ionosphere* [31], a binary classification problem over 34 features based on radar signal returns. The amount of data available with each dataset differs greatly. A random subset X_{ref} was held out from each dataset during training and later used for the computation of explanations. X_{ref} contains 300 samples for AG News, 1050 for Dry Bean, 114 for Breast Cancer and 71 for Ionosphere.

3.2 Models: Architecture, Training and Selection

For the tabular datasets we use small, fully connected Feed-Forward Neural Networks with ReLU activation functions. Models for Dry Bean, Breast Cancer and Ionosphere use 3x16, 16 and 8 neurons, respectively. For the AG News dataset we use a Bi-LSTM model with 128 dimensions for each direction and a fully connected output layer. We learn a 128 dimensional word embedding from scratch. We use the softmax function as output activation in all models.

We trained 100 models on each tabular dataset and 20 models on AG News. We fixed all random aspects of the model training except for the initialization of the network parameters. Each model observed exactly the same amount of data in the exact same order. All differences in model behavior will thus only stem from the initialization. To build the final Rashomon Set for each dataset,

Table 2: Mean accuracy and mean pairwise Jensen-Shannon-Distance (JSD) of all models over X_{ref} . All models were selected to lie within 5% accuracy of the best model. According to both metrics, all models perform nearly indistinguishably. JSD is bounded to $[0, 1]$.

	AG News	Dry Bean	Breast Cancer	Ionosphere
Mean accuracy on X_{ref}	0.91 ± 0.01	0.89 ± 0.01	0.95 ± 0.01	0.86 ± 0.01
Mean JSD on X_{ref}	0.0315 ± 0.004	0.0207 ± 0.004	0.0019 ± 0.001	0.0103 ± 0.007

we choose all models with at most 5% difference in accuracy to the best model. With the exception of the Ionosphere dataset, nearly all models are selected. We present average model accuracy and average pairwise output similarity computed with the Jensen-Shannon-Distance over X_{ref} in Table 2. All models achieve a high accuracy and are nearly indistinguishable by their output distributions.

3.3 Attribution Methods

We compare five attribution methods. From the family of gradient based methods we use Vanilla Grad (VG) [32], Smooth Grad (SG) [33], and Integrated Gradient (IG) [34]. From the family of perturbation based methods we include KernelSHAP (KS) [22] and LIME (LI) [27] for which we use the implementations provided by Captum². For IG, KS, and LI we use zero-baselines. SG samples with a noise ratio of 10%. Hyperparameters that further impact approximation behavior will be discussed in Section 4.1.

3.4 Model Dissimilarity Measures based on Attribution Scores

We use the following formula to express the scenarios in Table 1:

$$\mathcal{D}(f_a, f_b, X, \phi_1, \phi_2, d) = \frac{1}{|X|} \sum_{x \in X} d(\phi_1(f_a, x), \phi_2(f_b, x)) \quad (1)$$

where $f_a, f_b \in R$ are classifier functions from our Rashomon Set, $X \subseteq X_{\text{ref}} : \{x | x \in X_{\text{ref}} \wedge \arg \max f_a(x) = \arg \max f_b(x)\}$ is a subset of the reference set where both classifiers agree on the label, $\phi_1, \phi_2 \in \Phi = \{\text{VG}, \text{SG}, \text{IG}, \text{KS}, \text{LI}\}$ are the aforementioned attribution methods and $d \in D = \{\text{Feature Disagreement}, \text{Sign Disagreement}, \text{Euclid}, \text{Euclid-abs}\}$ are dissimilarity measures on attribution scores that we introduce now.

Feature Disagreement considers only the k top features (indices of k features of highest magnitude) from each of the two explanations and computes then the fraction of common features between them. *Sign Disagreement* is a more strict

² See project page at github.com/pytorch/captum

version of Feature Disagreement. It applies Feature Disagreement and then sub-selects only the top features that also have the same sign in both explanations. *Euclid* and *Euclid-abs* are the Euclidean distance and the Euclidean distance over absolute values of two attribution scores. Analogously to [19], for the disagreement measures we set $k = 11$ for AG News, $k = 4$ for Dry Bean, and $k = 8$ for both Breast Cancer and Ionosphere.

4 Examining the Rashomon Effect

We now present and discuss the experiments on numerical stability (Section 4.1), the Rashomon Effect itself (Section 4.2) and the Disagreement Problem (Section 4.3). Each section provides its own discussion.

4.1 Numerical Stability and the Rashomon Effect (111)

In this section we investigate the setting $\mathcal{D}(f_a, f_a, X = X_{\text{ref}}, \phi_1, \phi_1, \text{Euclid})$ to analyze the numerical stability of all ϕ_* w.r.t. differences of individual f_* .

Attribution methods often require to choose hyperparameters that control approximation behavior. For IG this is the number of steps used to approximate the integral. For SG, KS, and LI one has control over the number of samples evaluated during computation. This allows to adjust the computation time but if the parameter is too small, the resulting explanations may differ between two computations. We investigate this approximation stability across many models: Do explanations converge at the same hyperparameter for all models and if not, how large are the differences between individual models?

On the AG News dataset for SG we evaluate sampling hyperparameters $p \in [25, 50, 75, 100, 150]$, for IG, KS, and LI we evaluate $p \in [25, 50, 100, 150, 300]$. On the tabular datasets we compute the approximation stability for $p \in [25, 50, 75, 100, 125]$ for all methods. We quantify numerical stability in the following way: We compute ten SG, KS, and LI explanations for each sample in X_{ref} for each p . Next, we compute the pairwise Euclidean distances between all ten explanations. To obtain a stability score for one model, the average is taken across all samples in X_{ref} . As a final stability score we report the mean and standard deviation of this score across all models. IG depends deterministically on the number of steps in the integral, hence, we do not compute ten explanations per sample. Instead, we compute the pairwise distance between explanations for the same point obtained by p_i and p_{i+1} . To assess model dependent differences regarding the optimal choice of p , we compute for each model the smallest p_i in the set of parameters, where p_{i+1} did not improve the average stability by a factor of two.

Results for all datasets are presented in Table 3. The rows that start with SG, IG, KS, and LI report numerical stability for each method. The last row (#) reports the accumulated number of models whose explanations are stable at $\leq p$. The aggregated counts correspond to the attribution methods in the order as they appear in the rows: SG, IG, KS, LI. Unsurprisingly, numerical stability improves across all models with increasing p . At the same time, models clearly

Table 3: Explanation stability for sampling parameter p . We report mean \pm std across all models and samples for each attribution method. The values for IG describe the difference between using p_{i+1} instead of p_i . The last row (#) accumulates the number of models that converged at $\leq p$ for SG/IG/KS/LI.

	25	50	75	100	150
SG	0.0062 ± 0.0028	0.0044 ± 0.0020	0.0036 ± 0.0016	0.0039 ± 0.0029	0.0027 ± 0.0013
	25	50	100	150	300
IG	0.0734 ± 0.2740	0.0532 ± 0.2201	0.0412 ± 0.1389	0.0329 ± 0.1380	—
KS	$6.48e4 \pm 1.3e5$	$3.34e5 \pm 4.94e5$	$6.39e6 \pm 1.15e7$	1.7121 ± 0.1655	0.9515 ± 0.0695
LI	0.0470 ± 0.0117	0.0330 ± 0.0079	0.0220 ± 0.0055	0.0175 ± 0.0045	0.0119 ± 0.0032
#	-/1/-/-	10/2/-/1	18/18/-/19	19/18/20/19	20/20/20/20

(a) AG News. Total number of models is 20. The set of evaluated parameters is different from other datasets and different for SG from other methods.

	25	50	75	100	125
SG	0.0005 ± 0.0003	0.0004 ± 0.0002	0.0003 ± 0.0002	0.0003 ± 0.0001	0.0003 ± 0.0001
IG	0.0002 ± 0.0001	0.0001 ± 0.0001	0.0001 ± 0.0000	0.0000 ± 0.0000	—
KS	0.6087 ± 0.1724	0.2936 ± 0.0576	0.2232 ± 0.0422	0.1872 ± 0.0350	0.1645 ± 0.0306
LI	0.1561 ± 0.0413	0.0956 ± 0.0272	0.0723 ± 0.0219	0.0596 ± 0.0190	0.0515 ± 0.0170
#	-/-/-/-	39/-/73/96	62/99/99/99	87/99/99/99	99/99/99/99

(b) Beans. Total number of models is 99.

	25	50	75	100	125
SG	0.0253 ± 0.0140	0.0181 ± 0.0099	0.0148 ± 0.0081	0.0128 ± 0.0070	0.0114 ± 0.0063
IG	0.0030 ± 0.0012	0.0013 ± 0.0006	0.0008 ± 0.0004	0.0006 ± 0.0003	—
KS	$8.28e4 \pm 1.67e5$	0.4523 ± 0.0735	0.3036 ± 0.0402	0.2462 ± 0.0312	0.2126 ± 0.0265
LI	0.0506 ± 0.0152	0.0345 ± 0.0104	0.0279 ± 0.0085	0.0241 ± 0.0073	0.0215 ± 0.0065
#	-/-/-/-	28/-/-/65	58/100/95/86	91/100/100/100	100/100/100/100

(c) Breastcancer. Total number of models is 100.

	25	50	75	100	125
SG	0.0298 ± 0.0135	0.0209 ± 0.0095	0.0172 ± 0.0077	0.0148 ± 0.0066	0.0133 ± 0.0060
IG	0.0036 ± 0.0017	0.0017 ± 0.0008	0.0011 ± 0.0005	0.0009 ± 0.0004	—
KS	$1.02e5 \pm 2.06e5$	$1.81e3 \pm 3.67e3$	0.1990 ± 0.0279	0.1561 ± 0.0205	0.1325 ± 0.0169
LI	0.0995 ± 0.0239	0.0678 ± 0.0166	0.0547 ± 0.0136	0.0472 ± 0.0119	0.0423 ± 0.0107
#	-/-/-/-	13/-/-/34	28/51/42/43	47/51/51/50	51/51/51/51

(d) Ionosphere. Total number of models is 51.

respond differently to an increase in p . For SG the spread spans four values of p on each dataset and the selected values for p can differ by a factor of up to three. For IG there is no spread on the tabular datasets, but it has the largest spread compared to any other method on AG News. For KS and LI the models mostly split between two consecutive values. Note that KS displays conspicuously large numerical instability for smaller p , even on the smaller tabular datasets. Default parameters for KS and LI are set to 25 and 50 in Captum, which is insufficient for a large number of models. For the remainder of the paper we use explanations computed with the following p for all datasets: SG 100, IG 200, KS and LI 300.

Our results show that, for a rigorous workflow, hyperparameters need to be tuned not only based on the dataset but, in fact, for each model individually. Hence, choosing sensible default parameters is difficult. Providing implementations without default values or with very large values might be an option, though impeding user-friendliness. This learning also impacts any down-stream use of explanations such as benchmarking methods to assess the fidelity of an attribution method [3, 4, 10, 11, 17, 38] or explanation methods that build atop attributions to extract rules as explanations [2]. In those contexts, numerical stability is a pre-requisite to obtain reliable results.

4.2 Solution Diversity or: The Rashomon effect as seen with Different Dissimilarity Measures (011)

In this section we investigate how the Rashomon Effect manifests under different metrics over different attribution methods. In Table 2 we saw that the output behavior of the models is extremely similar. We are now interested to see how diverse the Rashomon Sets appear if we use the explanation based dissimilarity measure defined above. For each dataset and all dissimilarity measures $d \in D$ we evaluate $\mathcal{D}(f_a, f_b, X, \phi_1, \phi_1, d)$ for all pairs of $f_a, f_b \in R$ with $f_a \neq f_b$. Because all attribution scores are specific to the predicted class, we restrict $X \subseteq X_{\text{ref}} : \{x | x \in X_{\text{ref}} \wedge \arg \max f_a(x) = \arg \max f_b(x)\}$.

The distances produced by Feature Disagreement and Sign Disagreement are naturally bounded to the $[0, 1]$ interval. The gradient based attribution scores lie in a bounded range because we compute the gradient through the softmax output. Attribution scores for KS and LI produced distances larger than 1 with the two Euclidean metrics on all datasets. In those cases we normalize the Euclidean distances to the range $[0, 1]$ by dividing by the maximal distance observed. Fig. 1 visualizes the pairwise distances of all models as histograms. The x-axis discretizes dissimilarities, the farther to the right the more dissimilar. The y-axis is the number of distances in each bin.

Euclid and Euclid-abs overlap significantly in all cases except for IG and SG on the Ionosphere dataset. The Disagreement measures diverge on AG News and Dry Bean. Naturally, Sign Disagreement produces larger dissimilarity scores than Feature Disagreement.

In most of the cases, the means of the disagreement based measures and the Euclidean based measures lie relatively far apart. The exceptions are IG, KS, and LI on AG News, as well as KS and LI on Dry Bean. This means that one

metric always measures significantly more differences than the other, but what metric that is depends both on the dataset and method.

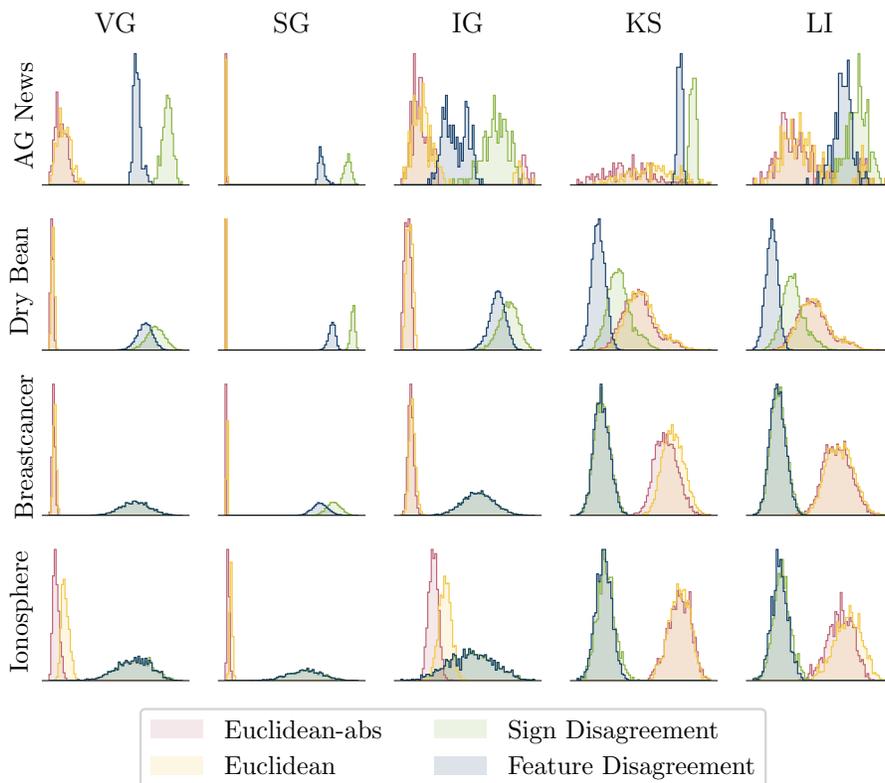


Fig. 1: Histograms over pairwise distances of all models according to Formula 1. Disagreement metrics computed with $k = 11, 4, 8, 8$ for AG News, Beans, Breast Cancer, and Ionosphere, respectively. In the bottom rows both disagreement metrics overlap nearly exactly.

We see that with most attribution methods and metrics the Rashomon Sets produce a large variety of distances across all models. This has strong implications for use cases where ML models, specifically (Deep) Neural Networks, are used as epistemic tools to develop hypotheses about the data generation process as it is becoming frequent practice in several disciplines [28, 39]. The variance in our results illustrates that the number of viable solution strategies is extensively large, hence, discovering all possibilities is highly improbable in cases where training a large number of models is infeasible. Methods such as ROAR [17] (despite being developed for a different purpose) could be useful to iteratively narrow down the search space but may still fail to uncover all pos-

Table 4: Kendall rank correlation coefficient (τ) between rankings of attribution method pairs. On average we observe a strong or very strong correlation, but the standard deviation indicates that for some models the set of methods that (dis)agree are very different compared to other models.

	AG News	Dry Bean	Breastcancer	Ionosphere
Feature Disagreement	0.67 ± 0.14	0.65 ± 0.26	0.66 ± 0.31	0.63 ± 0.28
Sign Disagreement	0.63 ± 0.10	0.57 ± 0.38	0.64 ± 0.37	0.73 ± 0.24
Euclidean	0.77 ± 0.29	0.69 ± 0.15	0.94 ± 0.12	0.95 ± 0.11
Euclidean-abs	0.79 ± 0.20	0.85 ± 0.15	0.81 ± 0.18	0.84 ± 0.14

sible correlations. The Rashomon Effect also has implications in user-centered scenarios. In cases where users interact with model explanations and expect a certain behavioral consistency over time, the deployment of a new model, even if performance itself is very similar, would pose a risk to user trust. Depending on the explanation method, the data domain, and the model, computing explanations can be very costly. Storing explanations for later re-use as a way to mitigate costs only works if the model stays the same.

In nearly all cases the human-oriented agreement metrics provide a very different picture than the Euclidean distances. Without additional knowledge about the suitability of a metric in a given context, practitioners should not rely on either disagreement or Euclidean measure alone. Use cases like [7], that use explanations to produce training signals for models, could benefit from exploring both kinds of metrics separately or from mixing them in a curriculum.

4.3 The Rashomon Effect and the Disagreement Problem (110)

In this section we investigate the Rashomon Effect on the Disagreement Problem. For all datasets and measures $d \in D$ we compare $\mathcal{D}(f_a, f_a, X, \phi_1, \phi_2, d)$ over all pairs $(\phi_1, \phi_2) \in \Phi \times \Phi$ with $\phi_1 \neq \phi_2$. As before, X is the set of all samples in X_{ref} where the predictions of both models agree.

Existing literature on the Disagreement Problem compares disagreement of method pairs for individual or very few models and only with the disagreement measures [5, 11, 19, 26]. These works report no consistent ranking between method pairs, especially when the data complexity increases.

We now analyze whether we find quantitative support for those observations. Additionally, we extend the analysis of the Disagreement Problem to include results based on the Euclidean distances.

For each individual model we rank the ten possible method pairs from most agreeing to most disagreeing. We calculate Kendall’s rank correlation coefficient τ for all model pairs with a sufficiently small p-value (< 0.05). For the remaining τ the mean and standard deviation across all models are reported in Table 4.

Two levels of correlation can be observed: 1) Stronger correlation $\gtrsim 0.8$ for Euclid on AG News, Euclid-abs on Dry Bean as well as both Euclidean based

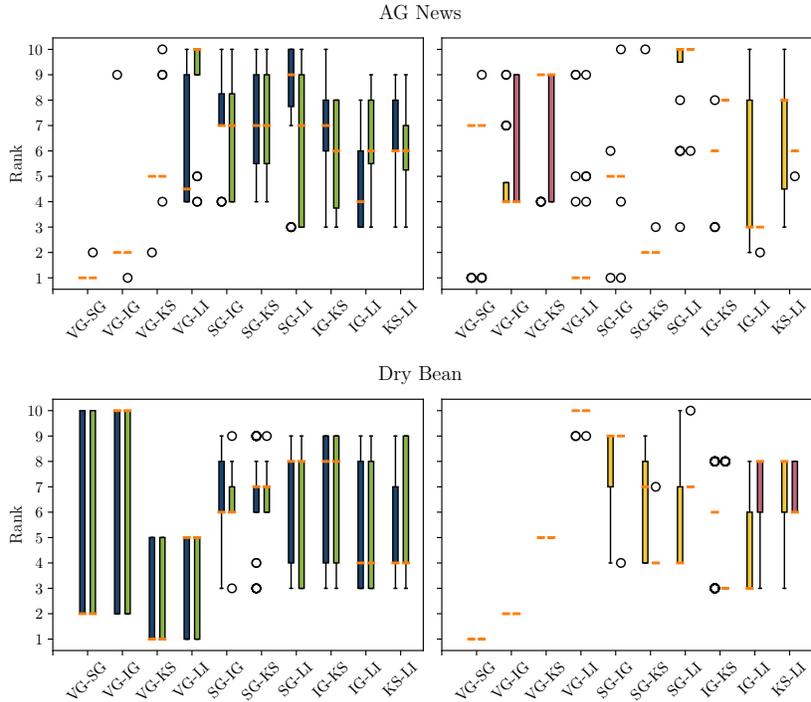


Fig. 2: Box plots of rankings over which pair of attribution methods disagrees most or least for individual models on AG News (top) and Dry Bean (bottom). Higher rank means larger disagreement. Plots on the left: Feature Disagreement Sign Disagreement, plots on the right: Euclid-abs Euclid; Orange lines in each boxplot indicate the median.

metrics on Breast Cancer and Ionosphere. 2) A moderate correlation ≈ 0.65 for Feature Disagreement on all datasets with a lower standard deviation on AG News. Sign Disagreement also falls in this range on all datasets but Dry Bean, with a notably lower standard deviation on AG News compared to other datasets. The lowest correlation (0.57) is produced by Sign Disagreement on Dry Bean, showing the largest standard deviation (0.38) at the same time. The large standard deviations suggest that a fair amount of models produces very different rankings, particularly in the case of the disagreement based rankings.

Are lower correlations structural? I.e. is it always specific method pairs that tend to swap ranks? We visualize the rank that each pairing occupies for every model in the box plots in Fig. 2 and Fig. 3. The y-axis shows the rank, higher rank meaning stronger disagreement relative to the other methods. Plots on the left pair both disagreement based rankings (blue/ green) while plots on the right show results for Euclidean based rankings (yellow/ red).

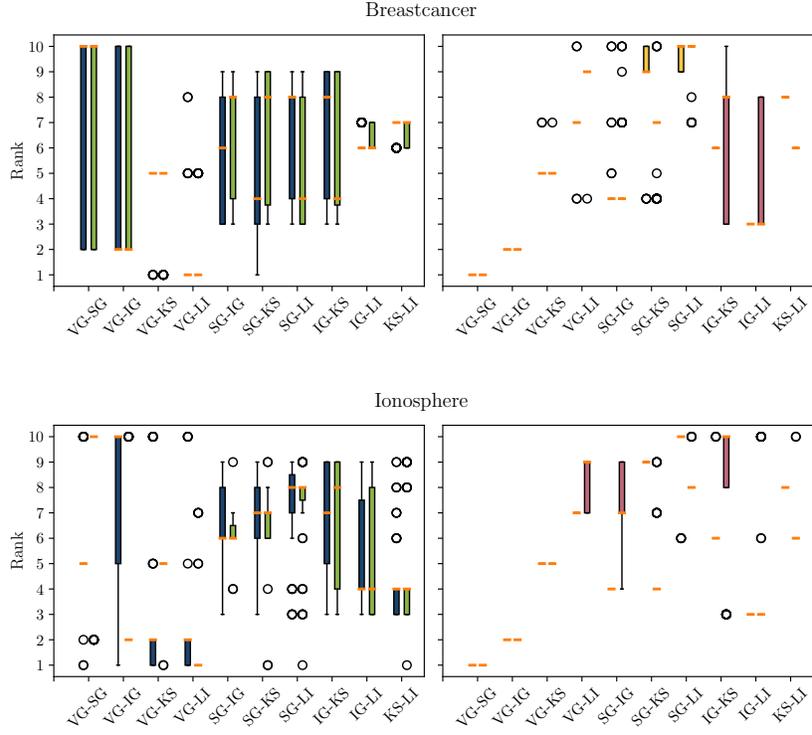


Fig. 3: Box plots of rankings over which pair of attribution methods disagrees most or least for individual models on Breast Cancer (top) and Ionosphere (bottom). Higher rank means larger disagreement. Plots on the left: Feature Disagreement Sign Disagreement, plots on the right: Euclid-abs Euclid; Orange lines in each boxplot indicate the median.

Generally, we can make the following observations about the Euclidean metrics: 1) For most explanation pairs, both Euclidean metrics show little to no variance within each dataset, signifying agreement on the ranking of the respective explainability pair, but no consistent ranking across all datasets. 2) All three tabular datasets agree for VG-SG being on rank one and VG-IG being on rank two, both with no variance.

Looking at the results for the disagreement metrics for each dataset in detail, we can see the following: AG News (Fig. 2) shows very stable rankings for both Disagreement metrics for VG- $\{SG, IG, KS\}$. For Dry Bean (Fig. 2) across both disagreement metrics, the median lies 8/20 times exactly on one of the quartiles which show no whisker. This means that 50% of the models agree on the respective ranking. This is interesting because at the same time VG- $\{SG, IG\}$ span nearly the whole ranking, meaning that all rankings in the fourth

quartile assign the maximum rank. The pairs SG- $\{IG, KS\}$ seem to swap places but are otherwise rather consistently placed in the lower middle of the ranking. Breast Cancer (Fig. 3) shows stable rankings for VG- $\{KS, LI\}$ with both disagreement metrics. More interestingly, for VG- $\{SG, IG\}$ the medians lie again on “whiskerless”-quartiles and the ranking agrees with the one on Dry Bean (rank 2 for VG-SG and rank 10 for VG-IG). In contrast to Dry Bean, here it is IG-LI and KS-LI that place comparably stable towards the middle of the ranking. On Ionosphere (Fig. 3 bottom) the plot shows smaller boxes compared to the other tasks. Taking outliers into account, multiple pairings span the whole ranking for disagreement based rankings. Ignoring outliers, there are five stable rankings for VG- $\{IG, SG, KS, LI\}$, four of which are achieved with Sign Disagreement.

We summarize our observations: We did not see a consistent ranking across all datasets and metrics. Our results for the disagreement based metrics support the observation from the literature that there is no consistent ranking among method pairs. However, we do not observe that results on the more complex AG News appear less correlated than for smaller tabular tasks. Our evaluation of Euclidean based rankings shows them to be notably more stable than their disagreement counterparts.

Interestingly, we cannot identify a single pair of methods that produces high disagreement across all tasks and metrics consistently, but there are pairs of methods for each dataset that consistently take mid-range rankings. Practitioners that seek diverse explanations would be recommended to start their search with comparing VG-KS, SG- $\{IG, KS\}$, and KS-LI.

4.4 Summary

In the first scenario in Section 4.1 we evaluated how sensitive individual models are to hyperparameter choices for non-deterministic attribution methods. Expectedly, a higher sampling rate always improves the numerical stability of the approximations. However, we found stark differences between the individual models, some requiring larger parameter values by a factor of up to twelve. This has direct implications for scientists and developers using XML methods, as it means that prior knowledge is not necessarily transferable between two models. Choosing default values is de facto impossible. Especially scenarios where parameters have to be chosen as small as possible require rigorous testing.

After verifying the numerical stability of our explanations, in Section 4.2 we assessed how the Rashomon Effect manifests itself on different datasets, depending on the different attribution methods and dissimilarity measures. We illustrated the solution diversity under different dissimilarity measures. We found that gradient based attribution methods in conjunction with Euclidean metrics showed smaller distances and low variance on the simpler tabular datasets. Disagreement based dissimilarity measures produced high distances and variances in nearly all cases. The distances are notably higher for Sign Disagreement compared to Feature Disagreement in half of the cases. We saw a large spectrum of distances for perturbation based methods in all cases. Our observation of large magnitudes and high variances in the distances has implications for ML

as an epistemic tool. It illustrates how large the space of possible viable solution strategies is, indicating the need to develop informed search strategies in the future [6], especially in complex or resource constrained scenarios. Also, the histograms of Euclidean and disagreement-based measures rarely show overlap, meaning practitioners will have to make context-specific choices on what type of metric to use. In cases where model behavior is explained to users, deploying a model update can lead to irritations as the explanations will likely change drastically between any two models; using the computationally intensive KS or LI seems to give the best chances to maintain somewhat consistent explanations. Conversely to the use-case of ML as an epistemic tool, a possible direction of future work is the inverse search problem of finding a better performing model that functions most similarly.

Investigating the Rashomon Effect on the Disagreement Problem in Section 4.3 revealed stark differences between results from the disagreement measures and the Euclidean distances. Neither of the two metric types produced rankings consistent across all datasets. Within tasks, the Euclidean metrics produced very stable rankings while the disagreement measures only occasionally produced a stable rank for a few pairs. Thus, our work provides quantitative support to the observations in [11, 19, 26] that are based on a small number of models, only. However, contrary to the literature, our results do not look more stable for smaller models on tabular datasets than for the Bi-LSTM model on AG News.

5 Conclusion

We have quantitatively shown how the Rashomon Effect impacts the application and interpretation of XML techniques and argue that it has to be taken into account by the XML community in the future. Along the three variables 1) the model, 2) the datum, and 3) the attribution method, we presented a structured investigation of the Rashomon Effect from three perspectives within XML.

Our quantitative analysis on numerical stability showed models to have individual sensitivity to hyperparameters of explanation methods. We have shown that choosing the most efficient setting requires careful tuning not only to a specific task or architecture, but in fact to every model instance individually, in order to guarantee stable explanations. Assessing the Rashomon Effect itself by measuring the diversity of solution strategies, we found that the solution space appears extensive, especially under the disagreement metrics. This poses challenges to applications of ML as an epistemic tool, as well as use cases where models are offered to consumers that expect consistent behavior. Our study of the Disagreement Problem provides quantitative support for previously anecdotal evidence. No consistent ranking persists across all datasets and the only option for practitioners that seek diverse explanations is trial and error. However, for each dataset individually we were able to identify a pair of methods that consistently take mid-range ranks. Using those rankings to systematically compare methods might yield insight into differences regarding what parts of model behavior each method is sensitive to.

Acknowledgments This research has been funded by the Federal Ministry of Education and Research of Germany and the state of North-Rhine Westphalia as part of the Lamarr-Institute for Machine Learning and Artificial Intelligence Lamarr22B. Part of PWs work has been funded by the Vienna Science and Technology Fund (WWTF) project ICT22-059.

Ethical Statement

In critical contexts, where persons are directly or indirectly impacted by a model, and where explanations are used to verify that model behavior is compliant with a given standard, proper use of explanation methods is of utmost importance. Hyperparameter choices have to be validated for each model individually. For model testing and validation procedures to be reliable they have to integrate this knowledge. Our work demonstrated that it is unreasonable to expect an explanation computed for one model, to be valid for another model, however similar their performance otherwise may be. Re-using explanations from one model to give as an explanation of behavior for another model is not possible and has to be avoided in critical scenarios.

References

1. Adadi, A., Berrada, M.: Peeking inside the black-box: A survey on explainable artificial intelligence (XAI). *IEEE Access* **6**, 52138–52160 (2018)
2. Alkhatib, A., Boström, H., Vazirgiannis, M.: Explaining Predictions by Characteristic Rules. In: *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML/PKDD) (2022)*
3. Alvarez-Melis, D., Jaakkola, T.S.: On the Robustness of Interpretability Methods. In: *Workshop on Human Interpretability in Machine Learning (WHI@ICML) (2018)*
4. Ancona, M., Ceolini, E., Öztireli, C., Gross, M.: Towards better understanding of gradient-based attribution methods for deep neural networks. In: *International Conference on Learning Representations, (ICLR) (2018)*
5. Atanasova, P., Simonsen, J.G., Lioma, C., Augenstein, I.: A Diagnostic Study of Explainability Techniques for Text Classification. In: *Conference on Empirical Methods in Natural Language Processing (EMNLP) (2020)*
6. Beckh, K., Müller, S., Jakobs, M., Toborek, V., Tan, H., Fischer, R., Welke, P., Houben, S., von Rueden, L.: Harnessing prior knowledge for explainable machine learning: An overview. In: *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*. pp. 450–463 (2023). <https://doi.org/10.1109/SaTML54575.2023.00038>
7. Bogun, A., Kostadinov, D., Borth, D.: Saliency Diversified Deep Ensemble for Robustness to Adversaries. In: *AAAI-22 Workshop on Adversarial Machine Learning and Beyond (2021)*
8. Breiman, L.: Statistical Modeling: The Two Cultures (with comments and a rejoinder by the author). *Statistical Science* **16**(3), 199 – 231 (2001)
9. Burkart, N., Huber, M.F.: A Survey on the Explainability of Supervised Machine Learning. *Journal of Artificial Intelligence Research* **70**, 245–317 (2021)

10. DeYoung, J., Jain, S., Rajani, N.F., Lehman, E., Xiong, C., Socher, R., Wallace, B.C.: ERASER: A Benchmark to Evaluate Rationalized NLP Models. In: Annual Meeting of the Association for Computational Linguistics (ACL) (2020)
11. ElShawi, R., Sherif, Y., Al-Mallah, M., Sakr, S.: Interpretability in healthcare: A comparative study of local machine learning interpretability techniques. *Computational Intelligence* **37**(4), 1633–1650 (2021)
12. Fisher, A., Rudin, C., Dominici, F.: All Models are Wrong, but Many are Useful: Learning a Variable’s Importance by Studying an Entire Class of Prediction Models Simultaneously. *Journal of Machine Learning Research* **20**(177), 1–81 (2019)
13. Flora, M., Potvin, C., McGovern, A., Handler, S.: Comparing Explanation Methods for Traditional Machine Learning Models Part 1: An Overview of Current Methods and Quantifying Their Disagreement. arXiv preprint arXiv:2211.08943 (2022)
14. Guidotti, R., Ruggieri, S.: Assessing the stability of interpretable models. arXiv preprint arXiv:1810.09352 (2018)
15. Han, T., Srinivas, S., Lakkaraju, H.: Which Explanation Should I Choose? A Function Approximation Perspective to Characterizing Post Hoc Explanations. In: Advances in Neural Information Processing Systems (NeurIPS) (2022)
16. Hancox-Li, L.: Robustness in Machine Learning Explanations: Does It Matter? In: Conference on Fairness, Accountability, and Transparency (FAT*) (2020)
17. Hooker, S., Erhan, D., Kindermans, P.J., Kim, B.: A benchmark for interpretability methods in deep neural networks. In: Advances in Neural Information Processing Systems (NeurIPS) (2019)
18. Koklu, M., Özkan, I.A.: Multiclass classification of dry beans using computer vision and machine learning techniques. *Computers and Electronics in Agriculture* **174**, 105507 (2020)
19. Krishna, S., Han, T., Gu, A., Pombra, J., Jabbari, S., Wu, S., Lakkaraju, H.: The Disagreement Problem in Explainable Machine Learning: A Practitioner’s Perspective. arXiv preprint arXiv:2202.01602 (2022)
20. Leventi-Peetz, A.M., Weber, K.: Rashomon Effect and Consistency in Explainable Artificial Intelligence (XAI). In: Future Technologies Conference (FTC) (2022)
21. Liu, F., Avci, B.: Incorporating Priors with Feature Attribution on Text Classification. In: Annual Meeting of the Association for Computational Linguistics (ACL) (2019)
22. Lundberg, S., Lee, S.I.: A Unified Approach to Interpreting Model Predictions. In: Advances in Neural Information Processing Systems (NeurIPS) (2017)
23. Marx, C.T., Calmon, F.P., Ustun, B.: Predictive multiplicity in classification. In: International Conference on Machine Learning (ICML) (2020)
24. Molnar, C.: *Interpretable Machine Learning*. 2nd edn. (2022)
25. Mücke, S., Pfahler, L.: Check Mate: A Sanity Check for Trustworthy AI. In: *Lernen. Wissen. Daten. Analysen. (LWDA)* (2022)
26. Neely, M., Schouten, S.F., Bleeker, M.J., Lucic, A.: Order in the Court: Explainable AI Methods Prone to Disagreement. arXiv preprint arXiv:2105.03287 (2021)
27. Ribeiro, M.T., Singh, S., Guestrin, C.: ”Why Should I Trust You?”: Explaining the predictions of any classifier. In: International Conference on Knowledge Discovery and Data Mining (KDD) (2016)
28. Roscher, R., Bohn, B., Duarte, M.F., Garcke, J.: Explainable Machine Learning for Scientific Insights and Discoveries. *IEEE Access* **8**, 42200–42216 (2020)
29. Schramowski, P., Stammer, W., Teso, S., Brugger, A., Herbert, F., Shao, X., Luigs, H.G., Mahlein, A.K., Kersting, K.: Making deep neural networks right for the right scientific reasons by interacting with their explanations. *Nature Machine Intelligence* **2**(8), 476–486 (2020)

30. Semenova, L., Rudin, C., Parr, R.: On the Existence of Simpler Machine Learning Models. In: Conference on Fairness, Accountability, and Transparency (FAccT) (2022)
31. Sigillito, V.G., Wing, S.P., Hutton, L.V., Baker, K.B.: Classification of radar returns from the ionosphere using neural networks. Johns Hopkins APL Technical Digest **10**(3), 262–266 (1989)
32. Simonyan, K., Vedaldi, A., Zisserman, A.: Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. In: International Conference on Learning Representations (ICLR) (2014)
33. Smilkov, D., Thorat, N., Kim, B., Viégas, F., Wattenberg, M.: Smoothgrad: removing noise by adding noise. arXiv preprint arXiv:1706.03825 (2017)
34. Sundararajan, M., Taly, A., Yan, Q.: Axiomatic attribution for deep networks. In: International Conference on Machine Learning (ICML) (2017)
35. Watson, M., Hasan, B.A.S., Al Moubayed, N.: Agree to Disagree: When Deep Learning Models With Identical Architectures Produce Distinct Explanations. In: Winter Conference on Applications of Computer Vision (WACV) (2022)
36. Wolberg, W., Street, N., Mangasarian, O.: Breast Cancer Wisconsin (Diagnostic). UCI Machine Learning Repository (1995)
37. Xin, R., Zhong, C., Chen, Z., Takagi, T., Seltzer, M., Rudin, C.: Exploring the Whole Rashomon Set of Sparse Decision Trees. In: Advances in Neural Information Processing Systems (NeurIPS) (2022)
38. Yeh, C., Hsieh, C., Suggala, A.S., Inouye, D.I., Ravikumar, P.: On the (In)Fidelity and Sensitivity of Explanations. In: Advances in Neural Information Processing Systems (NeurIPS) (2019)
39. Zednik, C., Boelsen, H.: Scientific Exploration and Explainable Artificial Intelligence. *Minds and Machines* **32**(1), 219–239 (2022)
40. Zhang, X., Zhao, J.J., LeCun, Y.: Character-level Convolutional Networks for Text Classification. In: Advances in Neural Information Processing Systems (NeurIPS) (2015)