

---

# LOCAL AND GLOBAL EXPLAINABILITY METRICS FOR MACHINE LEARNING PREDICTIONS

---

A PREPRINT

 **Cristian Munoz**  
Holistic AI  
London, UK  
cristian.munoz@holisticai.com

 **Kleyton da Costa**  
Holistic AI  
London, UK  
kleyton.costa@holisticai.com

 **Bernardo Modenesi**  
University of Michigan  
Ann Arbor, USA  
bmodene@umich.edu

 **Adriano Koshiyama**  
Holistic AI  
London, UK  
adriano.koshiyama@holisticai.com

February 24, 2023

## ABSTRACT

Rapid advancements in artificial intelligence (AI) technology have brought about a plethora of new challenges in terms of governance and regulation. AI systems are being integrated into various industries and sectors, creating a demand from decision-makers to possess a comprehensive and nuanced understanding of the capabilities and limitations of these systems. One critical aspect of this demand is the ability to explain the results of machine learning models, which is crucial to promoting transparency and trust in AI systems, as well as fundamental in helping machine learning models to be trained ethically. In this paper, we present novel quantitative metrics frameworks for interpreting the predictions of classifier and regressor models. The proposed metrics are model agnostic and are defined in order to be able to quantify: (i) the interpretability factors based on global and local feature importance distributions; (ii) the variability of feature impact on the model output; and (iii) the complexity of feature interactions within model decisions. We employ publicly available datasets to apply our proposed metrics to various machine learning models focused on predicting customers' credit risk (classification task) and real estate price valuation (regression task). The results expose how these metrics can provide a more comprehensive understanding of model predictions and facilitate better communication between decision-makers and stakeholders, thereby increasing the overall transparency and accountability of AI systems.

**Keywords** Artificial Intelligence · Explainability · Metrics · Risk Management in AI Systems

## 1 Introduction

In recent years, several methods have been proposed in order to bring more clarity to results generated by artificial intelligence (AI) models. The increasing creation of these methods is due to the extensive application of AI in different economic sectors, private and governmental, achieving an unprecedented degree of prediction accuracy, but providing very little understanding regarding how these predictions have been made. The low level of intelligibility of AI models has deemed these models a common nickname: “black boxes”.

A machine learning (ML) model is a program that maps inputs to outputs, with the goal of predicting outputs by learning prediction patterns from data. A “black box”, in the ML context, describes models that cannot be understood by simply analyzing their parameters. In general terms, an explainable or interpretable algorithm is one for which the rules for its prediction decisions can be questioned and explained in a way that is intelligible to humans. In particular, *interpretability* regards the ability to extract causal knowledge about the world from a model and *explainability* regards

being able to explain precisely how a complex model arrived at certain predictions, i.e. its mechanics. Understanding AI models' behavior is important to explain predictions to support the decision-making process, debug unexpected behavior for models (contributing to improving the model accuracy), refine modeling and data mining processes, verify that model behavior is reasonable and fair, present the model's predictions to stakeholders, etc.

In the case of deep learning (DL) models - like artificial neural networks - we can observe high accuracy with a low capacity for explaining the mechanics of either individual or group predictions. This characteristic made the "black-box" ML models a less attractive alternative for policy-makers and market players who need transparency, robustness, and compliance with rules used to support decision-making. Furthermore, ML models are essentially data-driven, and, in this case, the problem of harmful biases in datasets needs to be solved before the models' application.

In this research, we aim to contribute to the new and effervescent literature of explainable AI (XAI) methods. We propose and implement a set of new metrics of explainability in AI models context, each exhibiting distinct features, adding to the growing menu of tools for deciphering "black box" models and empowering AI stakeholders. For this, the paper is organized as follows: Section 2 describes the literature and common methodological approach from which we build our methods; Section 3 presents our proposed methods for explainable AI in local and global levels; Section 4 brings a set of applications for our new methods; and, we conclude the paper in Section 5.

## 2 Literature and Technical Background

The literature in explainable methods for AI models is more recent than the AI literature itself, and it has been gaining prominence as AI models proliferate in virtually all areas of society (Elshawi et al. [2019], ElShawi et al. [2020], Buckmann et al. [2022], Thimoteo et al. [2022]). Despite its infancy, the explainable AI (XAI) literature already possesses distinct taxonomies to classify its methods. In this section we discuss some of the XAI methods and forms to categorize them based on their properties. The discussion is far from being exhaustive, given how effervescent this field is, but the goal of this section is to help the reader understand our proposed XAI metrics and how they related to the rest of literature.

### 2.1 Taxonomies for XAI methods

Explainable AI methods can be categorized with respect to their properties. For instance, these methods can be grouped according to (i) their capacity to elucidate how AI models make predictions locally vs globally; as well as according to (ii) whether they are tailored to work for specific AI models vs they could work in any AI model (i.e. model-agnostic). There is not a single unified way of grouping XAI tools, but these are popular groupings that help elucidate how our tools fit into the literature<sup>1</sup>.

#### 2.1.1 Local vs Global

XAI tools can be categorized as Local or Global, with each of these categories aiming to explain how the AI model makes predictions at different levels. The local set of methods focuses on understanding AI algorithm's behavior at a low hierarchical level. Based often on a single observation (or on small subset of observations), local XAI methods provide metrics relative to how each feature contributes to the final prediction outputted by the AI model. The typical user of local metrics is an individual being targeted by an algorithm – e.g. lending company might need to provide reasons for why certain customer had a mortgage request denied by its AI model (Bracke et al. [2019]). This type of metric is very present in the recent literature of individual and counterfactual fairness, which builds and is built upon local XAI methods (Ge et al. [2022], Albini et al. [2022]).

The global set of XAI, on the other hand, focuses on understanding the AI algorithm's behaviour at a high hierarchical level, i.e. aiding its users to measure how features contribute to predictions over the whole dataset (or collection of datasets). Developers of AI algorithms, regulators, AI managers, science discovery, etc. are usually the main users of these tools, since they tend to be more interested in the general insights and knowledge discovery that the AI model produces, rather focusing at specific individual observations.

#### 2.1.2 Model-specific vs Model-agnostic

Model-specific XAI methods are often derived straightforwardly from the specific AI model that is fully transparent and explainable by design. That is, an additional explainability technique is not required to be overlaid on the AI model in

---

<sup>1</sup>There are other types of groupings, such as whether or not the method uses a model gradient vs a feature perturbation method, but these are beyond the scope of this project.

order to be able to fully explain its prediction process. As example, the result of linear models can be straightforwardly obtained and explained based on estimated coefficients.

Alternatively, the model-agnostic tools can be applied to the outputs of any AI algorithm, including very complex an opaque models, in order to provide interpretation on decision drivers for those models. According to Molnar [2021], this set of tools is *post-hoc*, i.e. they are applied after the model has been trained and they do not require access to its estimates or code, rather it only requires the ability to test the model. In this way, model-agnostic XAI methods are flexible to be used by auditors, researchers and other stakeholders, without having to disclose the model itself.

## 2.2 Relevant approaches in XAI methods

### 2.2.1 Feature Importance

The feature importance is a widely used XAI method and it measures the increase in prediction error after permuting the selected features. Intuitively, if by altering the values of certain feature results in a considerable change in the AI model error, this feature is considered to be important. Alternatively, features are deemed unimportant if the AI model error remains virtually unchanged after altering its values. The Algorithm 1 describes the process to achieve the feature importance of each feature  $j$  in a feature matrix  $X$ .

---

**Algorithm 1:** Permutation Feature Importance algorithm [Molnar, 2021], [Fisher et al., 2019]

---

Initialization: Trained model  $f$ , feature matrix  $X$ , target vector  $y$ , error measure  $\mathcal{L}(y, f)$ ;  
 Estimate the original model error  $\epsilon = \mathcal{L}(y, f(X))$ ;  
 For each feature:  $j = 0, 1, \dots, F$ ;  
     0. generate feature matrix  $X_{pv}$  by permuting values of feature  $j$  in the data  $X$ ;  
     1. estimate error  $\epsilon_{pv} = \mathcal{L}(y, f(X_{pv}))$ ;  
     2. calculate permutated values feature importance  $\mathcal{F}_j = \frac{\epsilon_{pv}}{\epsilon}$  or  $\mathcal{F}_j = \epsilon_{pv} - \epsilon$ ;  
 Sort features by descending  $\mathcal{F}_j$ ;

---

Permutation feature importance is a powerful tool because its interpretation is intuitive and it can be applied to any model – i.e. it is a easy to understand global XAI method. In addition, it does not require retraining the model, nor knowing its estimates and nor its *modus operandi*. On the other hand, there is no consensus in the literature about a training or test set should be used in order to compute the features importance.

### 2.2.2 Partial Dependence curve

A machine learning model  $f$  is often a function of a multitude of features  $x$ , which makes it infeasible to plot the estimated model in a high-dimensional space. Instead, the Partial Dependence curve (PD)<sup>2</sup>, can be used to assess how the predicted outcome of a model  $f$  behaves as a function of values of a particular *selected* feature  $x_s$ , after averaging  $f$  over the values of *all* other features  $x_A$ . Intuitively, the PD can be interpreted as the expected/average model response as a function of the input feature of interest [Pedregosa et al., 2011]. This is helps with model explainability since it makes possible to assess whether the relationship between the outcome and a feature is linear, and/or monotonic, or more complex. For example, if applied to a linear regression, the partial dependence plot always show a linear relationship.

According to Molnar [2021], the PD function for a predictive model  $f$  can be theoretically defined as,

$$f^{x_s}(x_s) := E_{x_A} [f(x_s, x_A)] = \int f(x_s, x_A) d\mathbb{P}(x_A) \quad (1)$$

where  $x_A$  are the remaining features and  $d\mathbb{P}$  their measure. In practice, after the model  $f$  is trained, we have an estimate of the model, denoted by  $\hat{f}(x_s, x_A)$ , and its PD curve with respect to  $x_s$  can be estimated calculating averages in the training data using Bayesian Monte Carlo methods to obtain:

$$\hat{f}^{x_s}(x_s) := \frac{1}{N} \sum_{i=1}^N \hat{f}(x_s, x_A^{(i)}) \quad (2)$$

The estimation assumes that features in  $A$  are not correlated. The violation of this assumption indicates that the averages calculated for the partial dependence may include data points that are implausible.

---

<sup>2</sup>The PD is also referred to as Partial Dependence Plot (PDP).

### 2.2.3 Surrogate model approximation

The relationship between features and the outcome of an AI model is often too complex to be easily summarized in “black-box” models. On the other hand, there are models where explainability is straightforward as discussed in 2.1.2, such as a linear regression or a relatively small decision tree. It turns out that one could use the latter simpler models – also referred to as surrogate models – to approximate a complex “black-box” model, often locally (i.e. for a subset of observations), and reap the explainability properties of the approximating model. In other words, for a subset of observations, it is possible that a simple and explainable model can approximate reasonably well complex predictions made by a complex model, and, in turn, offer a straightforward connection between features and output in the prediction process. This surrogate model approximation can be implemented by fitting a simple explainable model to a dataset corresponding to the same features used to train the complex AI model, while using the outcome generated by the AI model, as opposed to using the observed outcome.

## 3 Proposed XAI methods

In the previous section we laid out taxonomies and tools that we use to build our proposed XAI methods from. All of our 6 new metrics are model-agnostic, applicable to all types of AI models and they do not require knowledge of the model code nor estimates, only the ability to run predictions with these already working models. Furthermore, our proposed methods expand upon already existing XAI approaches, namely (i) Feature Importance, (ii) Partial Dependence curve and (iii) Surrogacy Models, as presented below.

### 3.1 Feature Importance-based metrics

Following the notation presented in 2.2.1, consider a set of features importance,  $\mathcal{F} = \{\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_F\}$ . We start by passing this vector through a normalizing function  $P(\cdot)$ , which normalizes the feature importance vector by the sum of their absolute values. A normalized value for feature importance  $j$  is defined as  $p_j := \frac{|\mathcal{F}_j|}{\sum_{i=1}^F |\mathcal{F}_i|}$ . This implies that the vector of normalized features importance  $P(\mathcal{F})$  constitutes a probability measure:

$$P(\mathcal{F}) := \{p_1, p_2, \dots, p_F\}, \quad \text{where } p_j \geq 0 \quad \text{and} \quad \sum_{i \in F} p_i = 1. \quad (3)$$

In order to evaluate the explainability of an AI model, an analysis of the distribution of the normalized feature weights can be conducted through various metrics, as proposed below.

#### 3.1.1 Metric 1 - Global Feature Importance Spread (GFIS)

The metric GFIS is based on the concept of entropy. More precisely on the entropy of the normalized features measure, which represents the concentration of information within a set of features. Lower entropy values indicate that the majority of the explanation is concentrated in a small number of features, suggesting that the model could be potentially explained more simply. Conversely, higher entropy values suggest that the importance of the features is more evenly distributed, resulting in a decreased level of explainability. The highest degree of complexity in explainability occurs when all feature weights are equal, is in a uniform distribution, which can serve a benchmark entropy for comparison. The uniform benchmark distribution would be defined as  $U = \{\bar{p}, \bar{p}, \dots, \bar{p}\}$ , where  $\bar{p} = 1/F$ .

On the assumption that a lower entropy indicates a more explainable model, a feature importance distribution  $P(\mathcal{F})$  can be compared to a baseline uniform distribution  $U$  to measure its level of explainability. However, simply comparing these two high dimensional objects is often not easy to humans, so we propose a form to summarize the differences into a single number, i.e. the spread of the two distributions can be measured using: (i) the Entropy Ratio; or (ii) the Kullback-Leibler Divergence; or (iii) the Gini coefficient. These metrics allow for a direct comparison between the two distributions, with the goal of determining the degree of explainability with a single measure:

Table 1: Submetrics of Metric 1

Submetric	Measure	Equation
1.1	Entropy Ratio	$S_{ER}(P(\mathcal{F})) = \frac{\sum_{j=1}^F p_j \log p_j}{\sum_{j=1}^F \bar{p} \log \bar{p}}$
1.2	Kullback-Leibler Divergence	$S_{KL}(P(\mathcal{F})) = D_{KL}(P(\mathcal{F})  U) = \sum_{j=1}^F p_j \log \left( \frac{p_j}{\bar{p}} \right)$
1.3	Gini-Coefficient	$S_G(P(\mathcal{F})) = \frac{\sum_{j=1}^F \sum_{j'=1}^F  p_j - p_{j'} }{2F^2 (\sum p_i / F)}$

The best (most easily explained AI) and worst (most complex AI) possible scores for  $S_{ER}$  are 0 and 1, respectively. The Kullback-Leibler Divergence  $S_{KL}$  is a statistical distance metric that measures the dissimilarity between a given feature importance distribution and a uniform distribution. As a distance metric, the range of  $S_{KL}$  is between 0 and  $\infty$ , with the lowest value indicating a uniform distribution and the worst possible explainability score. Lastly, the  $S_G$  coefficient ranges from 0 to 1, being 0 the most complex and 1 the least complex AI model, in contrary to how  $S_{ER}$  is measured.

### 3.1.2 Metric 2 - Local Feature Importance Spread Stability (LFISS)

Local Feature Importance refers to the assignment of feature normalized importance to different regions of the input data space. For a given dataset  $\mathcal{D}$  with  $N$  samples, it is possible to compute a vector of feature importance for each individual observation  $d$ . With observation-individualized measures, we create a feature importance matrix  $\mathcal{P} \in \mathbb{R}^{N \times F}$  where each row represents the normalized feature importance for each instance in the data space. Similar to the concept of global feature importance, the distribution associated with an instance  $d$  is represented by  $P^d(\mathcal{F}) = \mathcal{P}[d, :]$ , i.e. the  $d$ -th row of our matrix. For the entire dataset  $\mathcal{D}$ , the matrix  $\mathcal{P}$  is defined as:

$$\mathcal{P} = [P^1(\mathcal{F})^T \quad P^2(\mathcal{F})^T \quad \dots \quad P^N(\mathcal{F})^T]^T, \quad \text{where} \quad P^d(\mathcal{F}) = [p_1^d \quad p_2^d \quad \dots \quad p_F^d].$$

In comparison to the previous metric presented in 3.1.1, we now face a collection of vectors, as opposed to only one vector. However, we still aim to collapse all these vectors into a single straightforward and intuitive measure that reflects the complexity regarding how features contribute to the final outcome, relative to a benchmark model in which all features contribute equally – i.e. to a uniform feature importance distribution.

In order to accomplish the task in hand, we proceed in three steps. First, we follow the steps taken in 3.1.1 for each vector, i.e. we compute  $S^d := S(P^d)$  for all instance  $d \in \mathcal{D}$ , employing one of the  $S(\cdot)$  functions described in Table 1. We stack these results in a vector  $\mathcal{S} := \{S^1, S^2, \dots, S^N\}$ . In the second step, we apply the normalizing function  $P(\cdot)$  described in equation 3 in our vector, obtaining  $P(\mathcal{S})$ . Finally, we apply once more the  $S(\cdot)$  to measure the spread in  $P(\mathcal{S})$ . The resulting LFISS metric would be given by the following:

$$S_{LFISS} := S(P(\mathcal{S})), \quad \text{where} \quad \mathcal{S} := \{S(P^1), S(P^2), \dots, S(P^N)\}. \quad (4)$$

The interpretation of  $S_{LFISS}$  depends on the function  $S(\cdot)$  chosen from Table 1. For instance, if the KL measure is used, a lower value of  $S_{LFISS}$  suggests that features explain the output in a stable form: either mostly complex across the dataset observations or mostly simple across the observations. In order to know which is the case i.e. mostly complex vs mostly simple, the GFIS can be used. On the other hand, high values of  $S_{LFISS}$  using the KL measure suggests that few features impact the final output strongly for a small subset of the dataset. This can serve to identify for which observations this is the case.

### 3.1.3 Metric 3 - Predictions Groups Contrast (PGC)

The PGC metric compares the top-K ranking of features importance drawn from the entire dataset with the top-K ranking induced from specific subgroups of predictions. It can be applied to both categorical and regression problems, being useful for quantifying how the feature importance priority changes between subgroups. For categorical problems,  $G$  subgroups can represent output categories, and, for regression problems,  $G$  subgroups can be represented by different portions of the output support (e.g. quartiles). Given the top-k features importance indexed by  $I$ , and partial feature importance  $I^g$  computed using only observations from group  $g \in \mathcal{D}$ , two approaches can be used to analyze these two rankings are similar:

$$I_K := \text{argsort}(-P(\mathcal{F}))[1 : K] \quad \text{and} \quad I_K^g = \text{argsort}(-P_g(\mathcal{F}))[1 : K] \quad (5)$$

- Position Overlap Score

$$M_K(\mathcal{F}, g) = \frac{1}{K} \sum_{k \in K} m(I_k, I_k^g), \quad \text{where} \quad m(a, b) = \begin{cases} 1, & \text{if } a = b, \\ 0, & \text{otherwise.} \end{cases}$$

- Position Range Overlap Score

$$M_K(\mathcal{F}, g) = \frac{1}{K} \|I_K^g \cap I_K\|$$

Then, the predictions group contrast PGC is defined as:

$$C(\mathcal{F}) = \frac{1}{GF} \sum_g \sum_{K=1}^F M_K(\mathcal{F}, g)$$

where  $0 < C(\mathcal{F}) < 1$ . A higher contrast value suggests that the model follows a very similar behavior for each group. A lower value suggests that the model has a radical change in the behavior of each group.

### 3.1.4 Metric 4 - $\alpha$ -Feature Importance ( $\alpha_{FI}$ )

The  $\alpha$  Feature Importance metric quantifies the minimum proportion of features required to represent  $\alpha$  of the total importance. In other words, this metric is focused in obtaining the minimum number of features necessary to obtain no less than  $\alpha \times 100\%$  of the total explanation based on the normalized feature importance vector  $P(\mathcal{P})$  presented in equation 3. Define  $p_{(j)}$  as the  $j$ -th largest number of the vector  $P$ , i.e. the  $j$ -th order statistic. Formally, our proposed metric can be computed as follows,

$$\alpha_{FI}(\mathcal{F}) = \frac{\min_J \left\{ F - J \quad , \text{subject to} \quad \sum_{j=F}^J p_{(j)} > \alpha, J \leq F \right\}}{F} \quad (6)$$

where  $0 < \alpha_{FI} < 1$ . If  $\alpha$  is a relatively high number, such as  $4/5$ , then a low  $\alpha_{FI}$  value indicates that a considerable portion of the model’s behavior can be explained with a limited number of representative parameters. Conversely, a higher value indicates that a greater number of features are needed to describe the model’s behavior. In our illustrations, we refer to this metric as the  $FF$  or  $4/5$  metric, pinning the value of  $\alpha$ .

More broadly, an XAI user can calculate the minimal number of features for a given  $\alpha$ , namely  $M^*(\alpha) := F - J^*(\alpha)$ , for various values of  $0 \leq \alpha \leq 1$ . The resulting increasingly monotonic curve  $M^*(\alpha)$  is also informative. It can be used to identify the most important features to explain outputs up to a certain  $\alpha$  degree of explainability. When this curve plateaus, one can identify a set of features that are not so relevant/dispensable to explain the AI model outcomes.

## 3.2 Partial Dependence curve-based metric

### 3.2.1 Metric 5 - Partial Dependence Complexity (PDC)

The Partial Dependence Complexity metric uses the concept of the Partial Dependence curve presented in 2.2.2 to evaluate how simple this curve can be represented. The partial dependence curve is used to show model predictions are affected on average by each feature. Curves represented by linear functions are easy to grasp and interpret, but some features may have a highly non-linear and complex effects on the model outcomes. Our proposed metric aims to measure how much the curve changes across the variable domain. In order to calculate the PDC metric, we suggest the following steps. First, we compute the numerical 2nd derivative of the PD presented in 2. We denote the numerical 2nd derivative by  $\hat{f}^{x_s''}(\cdot)$ , which mathematically represents how fast the 1st derivative of the function changes. If the function is a simple line, the 2nd derivative is equal to zero. On the other hand, if the function is highly non-linear, the 1st derivative changes often, i.e. the 2nd derivative departs often from zero in magnitude. Given this intuition, our proposed metric consist the normalized norm of the vector containing the absolute value of the numerical 2nd derivative of the PD curve, denoted by  $\vec{v}_{x_s}$ , computed over  $q$  points of the  $x_s$  domain:

$$PDC_{x_s} := \left\| \frac{\vec{v}_{x_s}}{\|\vec{v}_{x_s}\|} \right\|, \quad \text{where} \quad \vec{v}_{x_s} := \left\{ |\hat{f}^{x_s''}(x_{s1})|, |\hat{f}^{x_s''}(x_{s2})|, \dots, |\hat{f}^{x_s''}(x_{sq})| \right\}. \quad (7)$$

The closer  $PDC_{x_s}$  is to zero, the more linear the outcome depends on feature  $x_s$ . On the other hand, the closer to one, less linear is the dependence of model outcomes to the feature  $x_s$ . It is also possible to compute the size of the vector  $\vec{v}_{x_s}$  without normalization and compare these values between different model architectures, or between features.

## 3.3 Surrogacy model-based metric

### 3.3.1 Metric 6 - Surrogacy Efficacy Score (SESc)

The Surrogacy Efficacy Score is a technique for gaining a better understanding of the inner workings of complex "black box" models, as discussed in 2.2.3. For example, by using a Tree-based model, this method provides a more interpretable representation of the model’s behavior by partitioning the input data based on the values of certain fields and creating simple rules to approximate the model’s predictions. The Decision Tree model is trained to closely mimic the original model by minimizing the loss between the model’s predictions and surrogate model predictions. Given a model  $f$ , a surrogate model  $DT$ , the input data  $x$ , the true output data  $y$ , the model  $f$  and the surrogate loss function are defined as:

$$\begin{aligned} \text{Loss}_f &= L(y, y'), \quad \text{where} \quad y' = f(x) \\ \text{Loss}_{DT} &= L(y', y''), \quad \text{where} \quad y'' = DT(x) \end{aligned}$$

After estimated our approximation surrogate model, we define our last metric as the following scores:

Problem Type	Equation
Classification	$E = \text{Accuracy}(y', y'')$
Regression	$E = \text{RMSE\_score}(y', y'')$

An added benefit of the Surrogacy Efficacy Score is the ability to derive a list of global feature importances. These importances can then be utilized to compute the metrics described in section 3.1, providing valuable insights into the relative importance of different features in the model’s predictions.

**3.4 Summary of Proposed Metrics**

Table 2: List of proposed metrics and their abbreviations to be used in the results section

Metric	Abbrev.	FI Type	Symbol	Output Domain	Reference
Fourth Fifts	PFI[80/100]	Global	FF	[0,1]	0
Feature Importance Spread - Divergence	PFID	Global	$S_D$	$[0, \infty)$	$\infty$
Feature Importance Spread - Ratio	PFIR	Global	$S_R$	[0,1]	0
Feature Importance Contrast - Overlap	PFICO	Global	$C_O$	[0,1]	1
Feature Importance Contrast - Range Overlap	PFICRO	Global	$C_{RO}$	[0,1]	1
Partial Dependence Interpretability Level	PDEL	Global	$L$	[0,1]	1
Feature Importance Spread Stability - Dataset	FISSD	Local	$T_D$	$[0, \infty)$	0
Feature Importance Spread Stability - Features	FISSF	Local	$T_F$	$[0, \infty)$	0
Surrogacy Efficacy Score	SFI	-	$E$	RMSE or Accuracy domain	-

The metrics presented in the previous section measure different aspects of explainability. This subsection aims to guide users with basic interpretations that can be obtained from these metrics:

- The Feature Importance Spread measures the concentration of feature importance. High concentrations facilitate the interpretability of a group of predictions as it requires prioritizing fewer features.
- The Feature Importance Contrast quantifies the disparity between features used to explain a group of predictions compared to the overall average attributed to the whole model. If a group has a high disparity, it is important to analyze.
- The Partial Dependence Interpretability Level quantifies the average complexity of curves that describe the dependence between the prediction and the input feature. A high level indicates high complexity in the dependence of the variables.
- The Feature Importance Stability measures the variation of importance in the feature space. High stability indicates that the importance is similar throughout the input domain. Low stability indicates high contrast in importance, and global indicators may not be the best scores.
- The Surrogacy Efficacy Score quantifies the veracity of a hypothesis: the prediction of the model can be explained by simple rules.

**4 Experiments**

In this section, we evaluate the benefits of using global and local feature importance metrics proposed in this paper. We employ publicly available datasets to apply our proposed metrics to various machine learning models focused on predicting customers’ credit risk (classification task) and real estate price valuation (regression task). For classification tasks, we consider the models logistic regression (LR), decision tree (DT), K-neighbors classifier (KNC), support vector classifier (SVC), AdaBoost, and multilayer perception (MLP). And for regression tasks, we consider linear regression (Linear), decision tree (DT), K-neighbors regressor (KNR), AdaBoost, and multilayer perceptron (MLP).

**4.1 Classification Tasks**

The German credit dataset is a collection of 1000 records that contain information on 67 different variables related to a customer. These variables provide insights into various aspects of the customer such as their financial status, age, gender, credit history, and other relevant information. The purpose of this dataset is to be used as a training resource for machine learning models. The target feature - credit default - represents the outcome the model is trying to predict.

By analyzing the relationships between the 67 features and the credit default, the model can learn patterns and make informed predictions on future cases.

Table 3: Reference table for the domain, premise, and interpretation based on metrics results

FI Type	Metric Name	Domain	Premise	Lower	Higher
PFI	PIF[80 100]	0-100	Output prediction represented by	Few Features	Many Features
	FISR	0-100	Features importance are	Concentrated	Spread
PFI Top-K	FISR	0-100	Features importance are	Concentrated	Spread
	FICRO[C0]	0-100	Features used to explain Group predictions	Different	Same
	FICRO[C1]	0-100	Features used to explain Group predictions	Different	Same
	PDEL	0-100	Interpret Partial Dependence Variables	Difcult	Easy
SFI	Accuracy	0-100	Represent ouptput prediction by simples rules is	Difcult	Easy
LFI	FISSDR	0-100	Feature Importance spread in the data space	Stable	Unstable
	FISSDM	0-INF	Average Feature Importance is	Spread	Concentrated
	FISSFR	0-100	Feature Importance in the data space	Stable	Unstable
LFI Top-K	FISSDR	0-100	Feature Importance spread in the data space	Stable	Unstable
	FISSDM	0-INF	Average Feature Importance is	Spread	Concentrated
	FISSFR	0-100	Feature Importance in the data space	Stable	Unstable

Table 3 provides a summary of the interpretation of each metric. The third column displays the *domain* of the variable, meaning the range of values that each metric can take. All metrics except for FISSDM (which can take values from zero to infinity) have a range of values from 0% to 100%. The following columns outline the *premise* for each metric and the interpretation that should be made based on the results (whether they are *lower* or *higher*).

Table 4: Explainability Metrics for models trained with German Credit Dataset

FI Type	Metric Name	LR	DT	KNC	SVC	AdaBoost	MLP
PFI	PIF[80 100]	27%	24%	2%	3%	29%	26%
	FISR	82%	75%	23%	32%	78%	71%
PFI Top-K	FISR	96%	88%	-	98%	93%	79%
	FICRO[C0]	74%	91%	100%	75%	86%	100%
	FICRO[C1]	81%	93%	100%	100%	77%	100%
	PDEL	100%	0%	100%	0%	0%	100%
SFI	Accuracy	80%	72%	79%	72%	78%	99%
LFI	DFISR	100%	100%	100%	100%	100%	99%
	FISSDM	0.41	0.36	0.37	0.38	0.39	0.48
	FISSDR	94%	97%	97%	98%	96%	94%
LFI Top-K	DFISR	100%	100%	99%	99%	100%	99%
	FISSDM	0.29	0.29	0.30	0.32	0.33	0.42
	FISSDR	96%	99%	99%	99%	98%	94%

Table 4 shows that LR and AdaBoost have a very similar PFI value, however, the FISR metric indicates that LR spread more importance between features than others. Even for class C0 feature importance spread more than C1 – Appendix X). The interpretability of LR output prediction is more complicated because more variables needed to be disentangled to explain the prediction. In contrast, KNC and SVC have the highest concentration (low FISR value) and few features will be needed to explain the prediction. If we check the metric over the top-80 features can be found that LR, SVC, and AdaBoost spread the importance among the top group. MLP and DT concentrate even more on the importance.

FICO and FICRO metrics quantify the features importance positions shared between group predictions. For instance, KNC and MLP have the same feature importance for classes C0 and C1. So, they use the same features priority to explain their output prediction. Big gaps between the FICRO class metric indicate more dissimilarities in their features priorities, for example, SVC has FICRO gap between C0 (75%) and C1 (100%). SFI accuracy shows that MLP output prediction is simple to explain with simple’s rules and DT and SVC have a more complicated behavior (not explainable with simple rules). It is important to understand that we are not talking about the efficacy of the model output. The surrogate model is trying to fit the model prediction behavior with a DT with 3 levels. FISSDR indicate the stability

level of the feature importance group. All the models have a good stability. Only the MLP model has a little change in their space. The average feature importance is a little unstable in all the models. LR and MLP show the highest importance spread. FISSDR shows that the feature importance value of LR and MLP are more unstable than other models and for MLP it continues happened in the Top-K model, which is a more critical dependency. DT, KNC and SVC show a heigh stability in the Top-k features importance.

## 4.2 Regression Tasks

The Real Estate dataset is a collection of 414 records that provide information on 6 different features related to the pricing of houses. These features include various aspects of a house such as its location, age, and other relevant information that can impact its value. The dataset is used to build ML regression models that aim to estimate the price of a house per unit area.

Table 5: Explainability Metrics for models trained with Real Estate Dataset

FI Type	Metric Name	Linear	DT	KNR	AdaBoost	MLP
PFI	PIF[80 100]	14%	43%	14%	29%	0%
	FISR	37%	86%	36%	70%	14%
PFI Top-K	FISR	-	100%	-	100%	100%
	FICS[CQ0-Q1]	0%	11%	100%	25%	100%
	FICS[CQ1-Q2]	100%	50%	100%	50%	79%
	FICS[CQ2-Q3]	100%	50%	100%	100%	79%
	FICS[CQ3-Q4]	100%	67%	100%	75%	86%
	PDEL	100%	100%	100%	100%	-
SFI	RMSE	3,51	9,65	6,59	4,17	1,23
LFI	DFISR	18%	7%	15%	6%	16%
	DFISM	0.48	0.27	0.26	0.28	0.42
	AFISR	92%	96%	98%	89%	92%
LFI Top-K	DFISR	95%	93%	94%	93%	95%
	DFISM	0.36	0.20	0.23	0.14	0.32
	AFISR	93%	98%	97%	97%	90%

We can see different degrees of important spread (FISR), DT and AdaBoost have the higher spread and MLP and KNR have the higher feature importance concentration. DT, Adaboost and MLP have more complicated output prediction interpretability because their group predictions have different feature priorities. The feature priorities for Linear and KNR are more stable. For Linear model, only the first group selected different features priorities. Partial dependency curves look easy to explain (PDEL). Surrogate model could match easily MLP output prediction and Linear model output prediction. The DFISR are very unstable in the data space so would be recommended use local analysis approaches to explain the model. For the Top-K output prediction a more stable behavior in the feature importance spread (DFISR) and average feature importance (AFISR).

## 5 Conclusion

The necessity of transparency in AI systems and the trade-off between accuracy and explainability in AI models increase the recent discussion about the interpretation of predictions generate with machine learning models. In this paper, we describe some models present in the literature and propose a new set of metrics concerning global and local feature importance. The interpretation structure based on premises improves the interpretation of results because in certain cases, viewing a single number (a metric) is more accurate than making a decision based on a set of graphical observations.

We observe that most strategies focused on explainability leave room for ambiguity in terms of rating or quantifying the model’s understanding and comparability. Furthermore, many strategies are only applicable to a specific architecture. A quantification of interpretability enables a practical understanding of what is happening inside complex AI models. To mitigate this gap, we propose six new explainability metrics based on three strategies - feature importance, partial dependence, and surrogacy. The results found show that the use of explainability metrics in AI systems can improve and empower AI models stakeholders.

We create a reference table to interpret quickly the most important explainability metrics. Between our proposed metrics, feature concentration, features priority shared between predictions groups, complexity represented by simple rules, and local feature importance stability revealed to be the most important in understanding our models.

As for future research, we plan on exploring the statistical properties of the metrics proposed, potentially creating inference tools for model selection and model auditing based on these tools.

## References

- Radwa Elshawi, Mouaz H. Al-Mallah, and Sherif Sakr. On the interpretability of machine learning-based model for predicting hypertension. *BMC Medical Informatics and Decision Making*, 19(1):146, Jul 2019. ISSN 1472-6947. doi:10.1186/s12911-019-0874-0. URL <https://doi.org/10.1186/s12911-019-0874-0>.
- Radwa ElShawi, Youssef Sherif, Mouaz Al-Mallah, and Sherif Sakr. Interpretability in healthcare: A comparative study of local machine learning interpretability techniques. *Computational Intelligence*, n/a(n/a), 2020. doi:<https://doi.org/10.1111/coin.12410>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/coin.12410>.
- Marcus Buckmann, Andreas Joseph, and Helena Robertson. An interpretable machine learning workflow with an application to economic forecasting. Technical report, Bank of England, 2022. URL <https://www.bankofengland.co.uk/working-paper/2022/an-interpretable-machine-learning-workflow-with-an-application-to-economic-forecasting>.
- Lucas M Thimoteo, Marley M Vellasco, Jorge Amaral, Karla Figueiredo, Cátia Lie Yokoyama, and Erito Marques. Explainable artificial intelligence for covid-19 diagnosis through blood test variables. *Journal of Control, Automation and Electrical Systems*, 33(2):625–644, 2022. URL <https://link.springer.com/article/10.1007/s40313-021-00858-y>.
- Philippe Bracke, Anupam Datta, Carsten Jung, and Shayak Sen. Machine learning explainability in finance: an application to default risk analysis. 2019. URL [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3435104](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3435104).
- Yingqiang Ge, Juntao Tan, Yan Zhu, Yinglong Xia, Jiebo Luo, Shuchang Liu, Zuohui Fu, Shijie Geng, Zelong Li, and Yongfeng Zhang. Explainable fairness in recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '22*, page 681–691, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450387323. doi:10.1145/3477495.3531973. URL <https://doi.org/10.1145/3477495.3531973>.
- Emanuele Albini, Jason Long, Danial Dervovic, and Daniele Magazzeni. Counterfactual shapley additive explanations. In *2022 ACM Conference on Fairness, Accountability, and Transparency, FAccT 2022*, page 1054–1070, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393522. doi:10.1145/3531146.3533168. URL <https://doi.org/10.1145/3531146.3533168>.
- Christoph Molnar. *Interpretable Machine Learning*. 2021. URL <https://christophm.github.io/interpretable-ml-book/>.
- Aaron Fisher, Cynthia Rudin, and Francesca Dominici. All models are wrong, but many are useful: Learning a variable’s importance by studying an entire class of prediction models simultaneously, 2019.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.