

Quantifying Interpretability of Arbitrary Machine Learning Models Through Functional Decomposition

Christoph Molnar, Giuseppe Casalicchio, and Bernd Bischl

Department of Statistics, LMU Munich,
Ludwigstr. 33, 80539 Munich, Germany
`christoph.molnar@stat.uni-muenchen.de`

Abstract. To obtain interpretable machine learning models, either interpretable models are constructed from the outset – e.g. shallow decision trees, rule lists, or sparse generalized linear models – or post-hoc interpretation methods – e.g. partial dependence or ALE plots – are employed. Both approaches have disadvantages. While the former can restrict the hypothesis space too conservatively, leading to potentially suboptimal solutions, the latter can produce too verbose or misleading results if the resulting model is too complex, especially w.r.t. feature interactions. We propose to make the compromise between predictive power and interpretability explicit by quantifying the complexity / interpretability of machine learning models. Based on functional decomposition, we propose measures of number of features used, interaction strength and main effect complexity. We show that post-hoc interpretation of models that minimize the three measures becomes more reliable and compact. Furthermore, we demonstrate the application of such measures in a multi-objective optimization approach which considers predictive power and interpretability at the same time.

Keywords: Interpretable Machine Learning · Explainable AI · Accumulated Local Effects · Multi-Objective Optimization

1 Introduction

Machine learning models are optimized for predictive performance, but it is often required to understand models in order to e.g. debug them, gain trust in the predictions and satisfy regulatory requirements. Therefore, performance often has to be traded off for interpretability. In areas such as life sciences and social sciences, it is common to restrict model selection to interpretable models such as (generalized) linear regression models and decision trees [23]. This often relies on an intuitive notion of interpretability, leading to an avoidance of "black boxes" such as tree ensembles and neural networks [5]. A restriction to structurally simpler models has the drawback that better performing models are often excluded a priori from model selection. An alternative is to allow any model and apply post-hoc interpretation methods to explain model behavior

and predictions. Interpretation methods quantify effects that features have on predictions, compute feature importances or explain individual predictions, see [24] for an overview. While model-agnostic post-hoc interpretation methods can – in general – be used regardless of model complexity, their reliability and compactness deteriorates when models use a high number of features, have strong feature interactions and complex feature main effects.

Model-agnostic interpretability measures are needed to make the compromise between interpretability and predictive performance explicit when selecting models [5, 31]. Instead of fixing the trade-off by preselecting an interpretable model class, model-agnostic measures would allow informed model selection with the desired balance between interpretability and predictive performance [13]. Interpretability is not well defined [23] and depends on user preferences and domain [14, 20, 30, 31]. This supports the conclusion in [4] that we cannot summarize interpretability with a single metric.

Contributions. We propose three model-agnostic measures of machine learning model interpretability. The measures can be used to compare trained models or to explicitly optimize interpretability during hyperparameter tuning and model selection. First we review related work on interpretability measures and the background of functional decomposition, on which our proposed measures are based. For the **number of features used** by the model, we propose an estimation heuristic. Based on the decomposition of the prediction function, we suggest measures for **interaction strength** and for **average complexity of the feature main effects**. We argue that minimizing these three measures improves the reliability and compactness of post-hoc interpretation methods. Finally, we illustrate the use of our proposed measures in multi-objective optimization and discuss implications of interpretability measures for the field of interpretable machine learning.

2 Related Work and Background

In this section we introduce the notation, review related work and describe the functional decomposition on which we base the proposed complexity measures.

Notation: We consider machine learning prediction functions $f : x \mapsto y$, where $x \in \mathbb{R}^p$ is a p -dimensional feature vector and $y \in \mathbb{R}$ is the prediction (e.g. regression output or a classification score). For the decomposition of this function, we write $f_S : x_S \mapsto y$, $S \subseteq \{1, \dots, p\}$, $x_S \in \mathbb{R}^{|S|}$ to denote a function that maps a vector with a subset of features to a marginal prediction. If subset S contains a single feature, we write f_j . We refer to the training data for the machine learning model with the tuples $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$ and refer to the value of the j -th feature from the i -th instance as $x_j^{(i)}$. We write X_j to refer to the j -th feature as a random variable.

2.1 Interpretability Measures

This section provides a non-exhaustive overview of approaches for measuring and optimizing interpretability. Many measures of interpretability are model-specific, i.e. only models of the same class can be compared (e.g. decision trees). Model size is often used as a measure for interpretability (e.g. number of decision rules, tree depth, ...) [4, 20, 31, 35]. Akaike's Information Criterion (AIC) [1] and the Bayesian Information Criterion (BIC) [33] are more widely applicable measures for the trade-off between goodness of fit and degrees of freedom. AIC and BIC fix a certain compromise between interpretability and performance, and consider only one dimension of interpretability, the degrees of freedom. In [36] the authors propose an interpretability evaluation model that considers the (model-specific) structural complexity of machine learning models. In additive models (e.g. linear regression), the number of features is often used as a measure of interpretability [32]. In [26] the authors propose model-agnostic measures of model stability, based on the semantic similarity of predictions when the model is re-trained on different subsamples of the training data. Similar to our approach, the measures are based on the predictions, not on the structure of the model.

In [27] the authors propose explanation fidelity and explanation stability metrics of local explanation models [29]. They propose to incorporate the metrics as a regularizer into the loss function of a neural network to simultaneously optimize for predictive performance and higher quality of local explanations. Their local explainability metrics complement ours, since we consider global model properties.

Further approaches measure interpretability as the usability of a (interpretable) model to support a human in a task, usually measured in a survey as response time, correctness of the response and task difficulty [10, 20, 36]. In [15] an interpretability measure based on runtime operation count is proposed and evaluated in user studies.

2.2 Functional Decomposition

Any high-dimensional prediction function can be decomposed into a sum of components with increasing dimensionality: an intercept, first-order feature effects, second-order effects and so on up to the p -th order effect:

$$f(x) = \underbrace{f_0}_{\text{Intercept}} + \underbrace{\sum_{j=1}^p f_j(x_j)}_{\text{1st order effects}} + \underbrace{\sum_{j < k}^p f_{jk}(x_j, x_k)}_{\text{2nd order effects}} + \dots + \underbrace{f_{1,\dots,p}(x_1, \dots, x_p)}_{\text{p-th order effect}} \quad (1)$$

This decomposition is only unique with additional constraints regarding the components. Stone [34] suggested orthogonality constraints and approximating the prediction function with weighted integrals. Hooker [18] defined centering, orthogonality and variance decomposition as desirable properties, resulting in

unique and hierarchically orthogonal components under the correlation inner product.

Accumulated Local Effects (ALE) were proposed in [3] as a tool for visualizing feature effects (e.g. Figure 1) and as an alternative unique decomposition of the prediction function with components $f_S = f_{S,ALE}$. The ALE decomposition is unique under an orthogonality-like property further described in [3].

The ALE main effect $f_{j,ALE}$ of a feature $x_j, j \in \{1, \dots, p\}$ for a prediction function f is defined as

$$f_{j,ALE}(x_j) = \int_{z_{0,j}}^{x_j} \mathbb{E} \left[\frac{\delta f(X_1, \dots, X_p)}{\delta X_j} \middle| X_j = z_j \right] dz_j - c_j \quad (2)$$

Here, $z_{0,j}$ is a lower bound of X_j (usually the minimum observed value of x_j) and the expectation \mathbb{E} is computed conditional on the value for x_j and over the marginal distribution of all other features. The constant c_j is chosen so that the mean of $f_{j,ALE}(x_j)$ with respect to the marginal distribution of X_j is zero. The ALE main effects are defined as the gradients of f with respect to the features, but are estimated with finite differences, i.e. access to the gradients of the model is not required. For the estimation we refer to [3]. We base our proposed measures on the ALE decomposition, because ALE are computationally cheap (worst case $O(n)$ for a feature main effect), the effects can be computed sequentially instead of simultaneously as in [18] and, most importantly, ALEs do not require knowledge of the joint data distribution. Additionally, ALE have software implementations [2, 25].

3 Functional Complexity

In this section we motivate complexity measures based on functional decomposition. Based on Equation 1, we decompose the prediction function into a constant (estimated as $f_0 = \frac{1}{n} \sum_{i=1}^n f(x^{(i)})$), plus the main effects (estimate with ALE), and a remainder term containing interactions (difference between full model and f_0 + main effects).

$$f(x) = f_0 + \underbrace{\sum_{j=1}^p \overbrace{f_{j,ALE}(x_j)}^{\text{MEC: How complex?}} + \overbrace{IA(x)}^{\text{IAS: Interaction strength?}}}_{\text{NF: How many features were used?}} \quad (3)$$

This arrangement of components emphasizes a decomposition of the prediction function into a main effect model and an interaction remainder. The main effect model itself can be used as a prediction function and we can analyze how well it approximates f , which is the idea behind the interaction measure IAS. The average main effect complexity (MEC) captures how many parameters are needed to describe the one-dimensional main effects on average. The number of features used (NF) describes how many features were used in the full prediction function.

3.1 Number of Features (NF)

We propose an approach based on feature permutation to determine how many features are used by a model. We regard features as "used" by the model when changing a feature changes the prediction, which may differ from the numbers of features available during training.

If available, a model-specific method for extracting the number of features used by the model is preferable, e.g. counting the number of non-zero weights in a sparse linear regression model. A model-agnostic heuristic is useful when the prediction function is accessible but not the internal structure of the model (e.g. prediction via API call), or when combining preprocessing steps and models complicates programmatic extraction (e.g. training a decision tree on sparse principal components).

The proposed procedure is formally described in Algorithm 1. To estimate whether the j -th feature was used, we sample instances from data \mathcal{D} , replace their j -th feature values with random values from the distribution of X_j (e.g. by sampling x_j from other instances from \mathcal{D}), and observe whether the predictions change. If the prediction of any sample changes, the feature was used for the prediction. The rate of false positives is zero, i.e. the probability that the heuristic

Algorithm 1: Number of Features Used (NF)

Input: Number of samples M , data \mathcal{D}

```

1 NF = 0
2 for  $j \in 1, \dots, p$  do
3   Draw  $M$  instances  $\{x^{(m)}\}_{m=1}^M$  from dataset  $\mathcal{D}$ 
4   Create  $\{x^{(m)*}\}_{m=1}^M$  as a copy of  $\{x^{(m)}\}_{m=1}^M$ 
5   for  $m \in 1, \dots, M$  do
6     Sample  $x_j^{(new)}$  from  $\{x_j^{(i)}\}_{i=1}^n$  with the constraint that  $x_j^{(new)} \neq x_j^{(m)}$ 
7     Set  $x_j^{(m)*} = x_j^{(new)}$ 
8   if  $f(x^{(m)*}) \neq f(x^{(m)})$  for any  $m \in \{1, \dots, M\}$  then  $NF = NF + 1$ .
9 return NF

```

counts a feature as used, but the model did not use the feature is zero. The probability of a false negative, i.e. the heuristic overlooks a feature, depends on the number of samples M , the model function f and the data distribution. Let P_{dep}^j be the probability that the prediction of a random instance depends on the value of x_j . For an instance that depends on x_j for its prediction, let P_{change}^j be the probability that a sample from X_j changes the prediction for an instance i . Then the probability of overlooking feature j is: $P_{fn}^j = (1 - P_{dep}^j + P_{dep}^j(1 - P_{change}^j))^M$. With the simplifying assumption that $P_{fn}^j = P_{fn} \forall j \in 1, \dots, p$, the probability that we miss at least one feature is $1 - (1 - P_{fn})^p$. For a linear model without interactions and only numerical features, the false negative rate

is 0: $P_{dep}^j = 1$ and $P_{change}^j = 1$, so that $P_{fn}^j = (1 - 1 + 0)^M = 0$. Let us assume a non-linear model where only one percent of instances rely on feature x_j ($P_{dep}^j = 0.01$) and these instances have a probability of 0.02 that the feature permutation changes the prediction ($P_{change}^j = 0.02$). If we set $M = 100$, then $P_{fn}^j = (0.99 + 0.01 \cdot 0.01)^{100} \approx 0.37$. If we increase M to 500, the probability that NF counts too few features drops to ≈ 0.007 .

We tested the NF heuristic with the Boston Housing data. We trained decision trees (CART) with maximum depths $\in \{1, 2, 10\}$ leading to 1, 2 and 4 features used and LASSO with penalty $\lambda \in \{10, 5, 2, 1, 0.1, 0.001\}$ leading to 0, 2, 3, 4, 11 and 13 features used. For each model we estimated NF with sample sizes $M \in \{10, 50, 500\}$ and repeated each estimation 100 times. For the elastic net models, NF was always equal to the true number of features. For the CART models the mean absolute differences between NF and the true number of features were 0.280 ($M = 10$), 0.020 ($M = 50$) and 0.000 ($M = 500$).

3.2 Interaction Strength (IAS)

Interactions between features mean that the prediction cannot be expressed as a sum of independent feature effects, but the effect of a feature depends on values of other features [24]. We propose to measure interaction strength as the scaled approximation error between the ALE main effect model and the prediction function f . Based on the ALE decomposition, the ALE main effect model is defined as the sum of first order ALE effects:

$$f_{ALE1st}(x) = f_0 + f_{1,ALE}(x_1) + \dots + f_{p,ALE}(x_p)$$

We define interaction strength as the approximation error measured with loss L :

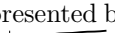

$$IAS = \frac{\mathbb{E}(L(f, f_{ALE1st}))}{\mathbb{E}(L(f, f_0))} \geq 0 \quad (4)$$

Here, f_0 is the mean of the predictions and can be interpreted as the functional decomposition where all feature effects are set to zero. IAS with the $L2$ loss equals 1 minus the R-squared measure, where the true targets y_i are replaced with $f(x^{(i)})$:

$$IAS = \frac{\sum_{i=1}^n (f(x^{(i)}) - f_{ALE1st}(x^{(i)}))^2}{\sum_{i=1}^n (f(x^{(i)}) - f_0)^2} = 1 - R^2$$

If $IAS = 0$, then $L(f, f_{ALE1st}) = 0$, which means that the first order ALE model perfectly approximates f and the model has no interactions. IAS can be larger than 0 for additive models for which we would expect $IAS = 0$, as observed in e.g. Table 1 (e.g. $IAS = 0.01$). This small deviation can occur when the true ALEs (see Equation 2) are not perfectly approximated by finite differences.

3.3 Main Effect Complexity (MEC)

To determine the average shape complexity of ALE main effects $f_{j,ALE}$, we propose the main effect complexity (MEC) measure. For a single ALE main effect, we define MEC_j as the number of parameters needed to approximate the curve with linear segments. For the entire model, MEC is the average MEC_j over all main effects, weighted with their variance. Figure 1 shows an ALE plot (= main effect) and its approximation with two linear segments. In the remainder, main effect ALE curves are represented by sparklines, e.g. Figure 1 by . Vertical bars, if present, e.g. , show borders between the segments of an approximation.

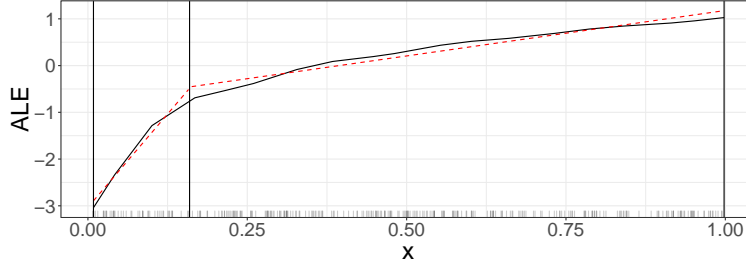



Fig. 1. ALE curve (solid line) approximated by two linear segments (dotted line).

Through the approximation with linear segments, the degrees of freedom required for describing a main effect curve become measurable. We measure the degrees of freedom as the number of non-zero coefficients for intercepts and slopes of the linear segments. The approximation allows some error, e.g. an almost linear main effect like  may have $MEC_j = 1$, even if dozens of parameters would be needed to describe it perfectly. The approximation quality is measured with R-squared, i.e. the proportion of variance of $f_{j,ALE}$ that is explained by the approximation with linear segments. An approximation has to reach an R-squared of at least $1 - \epsilon$, where ϵ is the user defined maximum approximation error. We also introduced parameter max_{seg} , the maximum number of segments. In the case that an approximation cannot reach an R-squared above $1 - \epsilon$ with a given max_{seg} , MEC_j is computed with the maximum number of segments. The selected maximum approximation error ϵ should be small, but not too small. We found ϵ between 0.01 and 0.1 visually meaningful. We apply a post-processing step that greedily sets slopes of the linear segments to zero, as long as R-squared $\in \{1 - \epsilon, 1\}$. This post-processing step potentially decreases the MEC_j , especially for models with constant segments like decision trees or rule-based models. MEC_j is averaged over all features to obtain the main effect complexity for the model. Each MEC_j is weighted with the variance of the corresponding ALE main effect to give more weight to features that contribute more to the prediction. For

example, the three main effect curves --- , ~ and ^ differ from each other by a scaling factor in the y-axis, but each requires five segments to be approximated with $R^2 \geq 0.95$: +++ , ~ and ^ . When weighted with the variance, --- gets a weight of 0.02, ~ a weight of 2.03 and ^ a weight of 18.28. Algorithm 2 describes the MEC computation in detail.

Algorithm 2: Main Effect Complexity (MEC).

Input: Prediction function f , approximation error ϵ , maximum number of segments max_{seg} , data \mathcal{D}

- 1 Define $R^2(g_j, f_{j,ALE}) := \sum_{i=1}^n (g_j(x_j^{(i)}) - f_{j,ALE}(x_j^{(i)}))^2 / \sum_{i=1}^n (f_{j,ALE}(x_j^{(i)}))^2$
- 2 **for** $j \in \{1, \dots, p\}$ **do**
- 3 Estimate $f_{j,ALE}$
- 4 // Approximate ALE with linear model
- 5 Fit $g_j(x_j) = \beta_0 + \beta_1 x_j$ predicting $f_{j,ALE}(x_j^{(i)})$ from $x_j^{(i)}$, $i \in 1, \dots, n$
- 6 Set $K = 1$ // Increase number of segments until approximation is good enough
- 7 **while** $K < max_{seg}$ **AND** $R^2(g_j, f_{j,ALE}) < (1 - \epsilon)$ **do**
- 8 // Optimize intervals Z_k via generalized simulated annealing, estimate β 's per segment with ordinary least squares
- 9 // For categorical feature, set slopes $\beta_{1,k}$ to zero
- 10 $g_j(x_j) = \sum_{k=1}^{K+1} \mathbb{I}_{x_j \in Z_k} \cdot (\beta_{0,k} + \beta_{1,k} x_j)$
- 11 Set $K = K + 1$
- 12 // Post-processing of slopes
- 13 **for** $k \in 1, \dots, K$ **do**
- 14 // n_k is number of instances in interval Z_k
- 15 Set $\beta_{0,k} = \frac{1}{n_k} \sum_{i: x_j^{(i)} \in Z_k} f_{j,ALE}(x_j^{(i)})$ and $\beta_{1,k} = 0$ in g_j
- 16 **if** $R^2(g_j, f_{j,ALE}) > (1 - \epsilon)$ **then** Use new $\beta_{0,k}$, $\beta_{1,k}$ in g_j
- 17 **else** Keep old $\beta_{0,k}$, $\beta_{1,k}$ in g_j
- 18 // Sum of non-zero coefficients minus first intercept
- 19 $MEC_j = K + \sum_{k=1}^K \mathbb{I}_{\beta_{1,k} > 0} - 1$
- 20 $V_j = \frac{1}{n} \sum_{i=1}^n (f_{j,ALE}(x_j^{(i)}))^2$
- 21 **return** $MEC = \frac{1}{\sum_{j=1}^p V_j} \sum_{j=1}^p V_j \cdot MEC_j$

4 Improving Post-hoc Interpretation

Minimizing the number of features (NF), the interaction strength (IAS) and the main effect complexity (AMEC) improves reliability and compactness of post-hoc interpretation methods such as partial dependence plots, ALE plots, feature importance, interaction effects and local surrogate models.

The less features, the less verbose the interpretations. Our NF measure improves the readability of post-hoc analysis results. The computational complexity and output size of most interpretation methods scales with $O(NF)$, like feature effect plots [3, 16] or feature importance [8, 12]. As shown in Table 2, a model with fewer features has a more compact representation and if additionally $IAS = 0$, the ALE main effects fully characterize the prediction function. Interpretation methods that analyze 2-way feature interactions scale with $O(NF^2)$. A complete functional decomposition [3, 18] would require to estimate $\sum_{k=1}^{NF} \binom{NF}{k}$ components which has a computational complexity of $O(2^{NF})$.

The less interaction, the more reliable feature effects. Feature effect plots, such as partial dependence plots and ALE plots visualize the marginal relationship between a feature and the prediction. The estimated effects are averages across instances. The effects can vary greatly for individual instances and even have opposite directions when the model includes feature interactions.

In the following simulation, we trained three models with different capabilities of modeling interactions between features: a linear regression model, a support vector machine (radial basis kernel, $C=0.05$) and gradient boosted trees. We simulated 500 data points with 4 features and a continuous target based on Friedman et. al (1991)[17]. The features are uniformly distributed in the following intervals: $0 \leq x_1 \leq 100$, $40\pi \leq x_2 \leq 560\pi$, $0 \leq x_3 \leq 1$, $1 \leq x_4 \leq 11$. The target was simulated as: $y = (x_1^2 + (x_2 \cdot x_3 - (1/(x_2 \cdot x_4)))^2)^{0.5} + e$, where $e \sim N(0, 125)$. Figure 2 shows an increasing interaction strength depending on the model used. This means that we prefer models with less interaction when using feature effect plots.

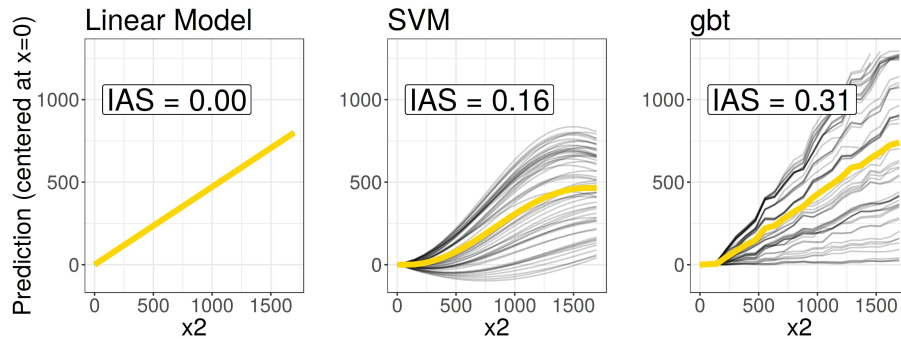


Fig. 2. The higher the interaction strength in a model (IAS increases from left to right), the less representative the Partial Dependence Plot (light thick line) becomes for individual instances represented by their Individual Conditional Expectation curves (dark thin lines).

The less complex the main effects, the better summarizable. In linear models, a feature effect can be expressed by a single number, the regression coefficient. If effects are non-linear the method of choice is visualization [3, 16]. Summarizing the effects with a single number (e.g. average marginal effects [22]) can be misleading, e.g. if the effect has a U-shape, the average effect might be zero. As a by-product of MEC, there is a third option: Instead of reporting a single number, the coefficients of the segmented linear model can be reported. Minimizing MEC means preferring models with main effects that can be described with fewer numbers, offering a more compact model description.

5 Optimization of Performance and Interpretability

As one of the main applications of the proposed interpretability measures, we demonstrate model selection for performance and interpretability in a multi-objective optimization approach.

Predicting Wine Quality. We used the wine quality dataset [9] which contains physical-chemical properties such as alcohol and residual sugar of 4870 white wines. The goal was to predict wine quality on a scale of 0 to 10, assessed by the median of three blind ratings.

Motivation. As [14] emphasizes, it is difficult to know the desired compromise between interpretability and performance before modeling the data and suggests multi-objective optimization. This stands in contrast to a priori selecting an interpretable model class (e.g. decision rules) and optimizing within this class or exclusively optimizing performance and applying post-hoc interpretations. We suggest searching over a wide spectrum of model classes and hyper parameter settings, presenting the set of Pareto optimal models, and allowing the practitioner to choose a suitable compromise between interpretability and performance. The three measures of interpretability provide a detailed characterization of machine learning models, which enables making informed decisions (e.g. how much does performance suffer if we use a model without interactions?).

Optimization Setup. We used the mlrMBO model-based optimization framework [19] to find the best model based on four objectives: number of features used by the model (NF), main effect complexity (MEC), interaction strength (IAS) and the cross-validated mean absolute error (MAE). We optimized over the space of following model classes (and hyperparameters): **CART** (maximum tree-depth and pruning cp), support **vector machine** (cost C and sigma), **elastic net** regression (regularization alpha and penalization lambda), **gradient boosted trees** (maximum depth, number of iterations), **gradient boosted generalized additive model** (number of iterations) and **random forest** (mtry).

Model-based Optimization Setup. We used the ParEGO algorithm [21] for multi-objective optimization. Within the fitness function, the MAE was estimated using 5-fold cross-validation and the other measures (NF, MEC, IAS) were estimated using all data instances. We set the number of iterations of ParEGO to 350. For all other parameters of the model-based, multi-objective optimization, we relied on the sensitive defaults provided by [7].

Table 1. Pareto front of models minimizing mean absolute error (MAE), number of features (NF), main effect complexity (MEC) and interaction strength (IAS).

Model (Hyperparameters)	MAE	MEC	IAS	NF
1 gbt (max_depth:10,nrounds:780)	0.40	3.50	0.71	11.00
2 gbt (max_depth: 8,nrounds:266)	0.41	4.10	0.64	11.00
3 rf (mtry: 5)	0.43	2.40	0.50	11.00
4 rf (mtry: 2)	0.44	2.40	0.48	11.00
5 rf (mtry: 1)	0.45	2.80	0.47	11.00
6 gbt (max_depth: 3,nrounds:617)	0.48	7.30	0.41	11.00
7 gbt (max_depth: 2,nrounds:931)	0.51	8.10	0.26	11.00
8 gbt (max_depth: 2,nrounds:100)	0.54	3.40	0.10	11.00
9 gbt (max_depth: 1,nrounds:949)	0.55	4.40	0.02	11.00
10 gamb (mstop:265)	0.57	1.70	0.00	11.00
11 svm (C:738.6223,sigma:2e-04)	0.57	1.10	0.05	11.00
12 svm (C:126.3303,sigma:2e-04)	0.58	1.00	0.01	11.00
13 gbt (max_depth: 1,nrounds: 43)	0.58	2.40	0.01	10.00
14 CART (maxdepth: 7,cp:0.0038)	0.58	2.30	0.27	10.00
15 CART (maxdepth:12,cp:0.0057)	0.59	2.00	0.21	5.00
16 elastic net (alpha:0.4723,lambda:0.0526)	0.59	1.00	0.00	8.00
17 elastic net (alpha:0.6471,lambda:0.0856)	0.60	1.00	0.00	6.00
18 CART (maxdepth:20,cp:0.0073)	0.60	2.00	0.20	4.00
19 elastic net (alpha:0.8768,lambda:0.0908)	0.61	1.00	0.00	2.00
20 elastic net (alpha:0.8681,lambda:0.2227)	0.63	1.00	0.00	1.00
21 median	0.67	0.00	0.00	0.00

Results. Table 1 shows the set of Pareto-optimal models along with their MAE, NF, IAS and MEC. If two models from the same class with different hyperparameter values had exactly the same MAE, NF, IAS and MEC, we kept only one and dropped the others. We also dropped constant models (e.g. elastic net regression with strong penalization), with exception of the median model. For a more informative visualization, we propose to visualize the main effects together with the measures in Table 2. The four selected models show different trade-offs between the four measures.

Performance-Interpretability Trade-off. The complexity measures allow to study the trade-off between interpretability and performance across different model classes and hyperparameter settings. We mapped each measure to

Table 2. A selection of four models from the Pareto optimal set, along with their ALE main effect curves. From left to right, the columns show models with 1) lowest MAE, 2) lowest MAE when $MEC = 1$, 3) lowest MAE when $IAS \leq 0.1$, and 4) lowest MAE with $NF \leq 7$. Corresponding hyperparameters can be found in Table 1.

	gbt [row 1]	svm [row 12]	gbt [row 8]	CART [row 15]
MAE	0.4	0.58	0.54	0.59
MEC	3.5	1	3.4	2
IAS	0.71	0.01	0.1	0.21
NF	11	11	11	5
fixed.acidity				
volatile.acidity				
citric.acid				
residual.sugar				
chlorides				
free.sulfur.dioxide				
total.sulfur.dioxide				
density				
pH				
sulphates				
alcohol				

the interval $[0, 1]$ by scaling each measure M with meaningful upper and lower bounds:

$$M_{scaled}(M) = \frac{M - M_{inf}}{M_{sup} - M_{inf}}$$

For MAE, we set M_{inf} to the lowest observed MAE of all models and $M_{sup} = \frac{1}{n} \sum_{i=1}^n |y^{(i)} - median(y)|$. For NF, we set $M_{inf} = 0$ and $M_{sup} = p$. For MEC, we set $M_{inf} = 0$ and M_{sup} to the highest observed MEC of all models. For IAS, we set $M_{inf} = 0$ and M_{sup} to the highest observed IAS of all models. To combine the three measures in a single dimension, we mapped interpretability ad-hoc as: $Interpretability = 3 - (NF_{scaled} + IAS_{scaled} + MEC_{scaled})$. This weights all three (scaled) measures equally. The maximum interpretability is 3 for the constant model that always predicts the median wine quality. The theoretical minimum interpretability is 0 for a model that uses all features and has the highest interaction strength and highest main effect complexity measure among all Pareto optimal models. Figure 3 maps each model and hyperparameter configuration from the Pareto set to the performance / interpretability space.

6 Discussion

We proposed three model-agnostic measures for machine learning model complexity based on functional decomposition: number of features used, interaction

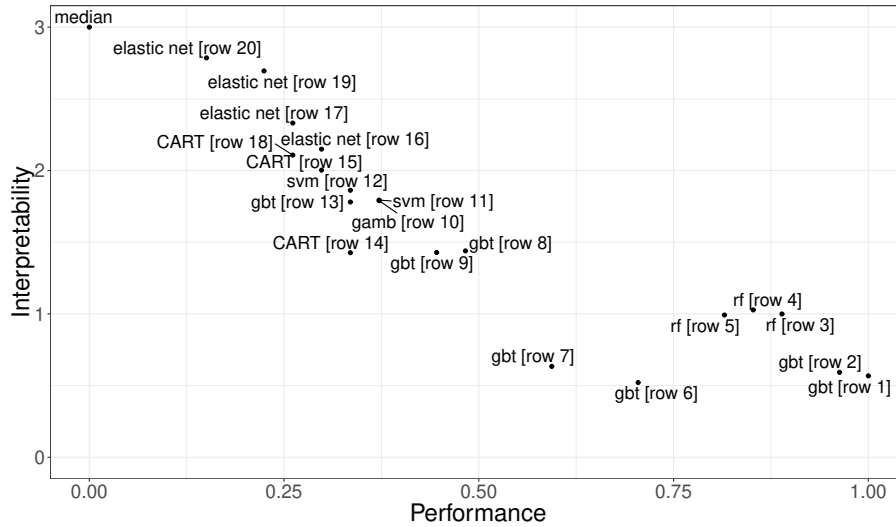


Fig. 3. Performance vs. interpretability tradeoff for predicting wine quality. Corresponding hyperparameters and measures are shown in Table 1.

strength and main effect complexity. Due to their model-agnostic nature, the measures allow model comparison across different model classes. We argued that minimizing these measures for a machine learning model improves its post-hoc interpretation. We demonstrated that the measures can be optimized directly with multi-objective optimization to make the trade-off between interpretability and performance explicit. Our proposed measures can be applied to a wide range of models since they are model-agnostic and work for regression and binary classification (based on classification scores / probabilities). We formulated the measures with both continuous and categorical features in mind, but leave an in-depth investigation of categorical features open for future work. The measures can be used for model selection, for model benchmarks and as objectives in automated machine learning frameworks.

Limitations. The proposed decomposition of the prediction function and definition of the complexity measures will not be appropriate in every situation. For example, all higher order effects are combined into a single interaction strength measure that does not distinguish between two-way interactions and higher order interactions. Two models can have the same IAS, but one has a single two-way interaction, the other many different higher order interactions. However, the framework of ALE decomposition allows to estimate higher order effects and to construct different interaction measures. The main effect complexity measure only considers linear segments but not e.g. seasonal components or other structures. Furthermore, the complexity measures quantify machine learning models

from a functional point of view and ignore the structure of the model (e.g. whether it can be represented by a tree).

The bigger picture. Interpretability is a high-dimensional concept (sparsity, additivity, fidelity, human simulability, ...) and we need several approaches to make interpretability measurable. In this context we see our work complementary to other approaches [10, 15, 27], which together form a basis for a more rigorous definition of interpretability as demanded by [11, 23]. Availability of different interpretability measures also fits in with the notion that interpretability depends on the audience and the context [30]. Different situations require differently weighted interpretability measures. In some situations we might prefer sparseness and a lack of interactions, in others it might be important that we can represent the model as a decision list. A multi-dimensional view of interpretability solves the lack of definition of interpretability and supports researcher to make quantified, verifiable and clearer statements about interpretability.

Implementation. The code for this paper is available at https://github.com/compstat-lmu/paper_2019_iml_measures. For the examples and experiments we relied on the `mlr` package [6] in R [28].

Acknowledgements. This work is funded by the Bavarian State Ministry of Science and the Arts in the framework of the Centre Digitisation.Bavaria (ZD.B) and supported by the German Federal Ministry of Education and Research (BMBF) under Grant No. 01IS18036A. The authors of this work take full responsibilities for its content.

References

- [1] Akaike, H.: Information theory and an extension of the maximum likelihood principle. In: Selected papers of hirotugu akaike, pp. 199–213. Springer (1998)
- [2] Apley, D.: ALEPlot: Accumulated Local Effects (ALE) Plots and Partial Dependence(PD) Plots (2017), <https://CRAN.R-project.org/package=ALEPlot>, r package version 1.0
- [3] Apley, D.W.: Visualizing the effects of predictor variables in black box supervised learning models. arXiv preprint arXiv:1612.08468 (2016)
- [4] Askira-Gelman, I.: Knowledge discovery: comprehensibility of the results. In: Proceedings of the thirty-first Hawaii international conference on system sciences. vol. 5, pp. 247–255. IEEE (1998)
- [5] Bibal, A., Frénay, B.: Interpretability of machine learning models and representations: an introduction. In: Proceedings on ESANN. pp. 77–82 (2016)
- [6] Bischl, B., Lang, M., Kothhoff, L., Schiffner, J., Richter, J., Studerus, E., Casalicchio, G., Jones, Z.M.: `mlr`: Machine learning in R. *J. Mach. Learn. Res.* **17**(170), 1–5 (2016)

- [7] Bischl, B., Richter, J., Bossek, J., Horn, D., Thomas, J., Lang, M.: mlrMBO: A Modular Framework for Model-Based Optimization of Expensive Black-Box Functions (2017), <http://arxiv.org/abs/1703.03373>
- [8] Casalicchio, G., Molnar, C., Bischl, B.: Visualizing the feature importance for black box models. In: Joint European Conference on Machine Learning and Knowledge Discovery in Databases. pp. 655–670. Springer (2018)
- [9] Cortez, P., Cerdeira, A., Almeida, F., Matos, T., Reis, J.: Modeling wine preferences by data mining from physicochemical properties. *Decision Support Systems* **47**(4), 547–553 (2009)
- [10] Dhurandhar, A., Iyengar, V., Luss, R., Shanmugam, K.: Tip: Typifying the interpretability of procedures. arXiv preprint arXiv:1706.02952 (2017)
- [11] Doshi-Velez, F., Kim, B.: Towards a rigorous science of interpretable machine learning. arXiv preprint arXiv:1702.08608 (2017)
- [12] Fisher, A., Rudin, C., Dominici, F.: All models are wrong but many are useful: Variable importance for black-box, proprietary, or misspecified prediction models, using model class reliance. arXiv preprint arXiv:1801.01489 (2018)
- [13] Freitas, A.A.: A critical review of multi-objective optimization in data mining: a position paper. *ACM SIGKDD Explorations Newsletter* **6**(2), 77–86 (2004)
- [14] Freitas, A.A.: Comprehensible classification models: a position paper. *ACM SIGKDD explorations newsletter* **15**(1), 1–10 (2014)
- [15] Friedler, S.A., Roy, C.D., Scheidegger, C., Slack, D.: Assessing the local interpretability of machine learning models. arXiv preprint arXiv:1902.03501 (2019)
- [16] Friedman, J.H.: Greedy function approximation: a gradient boosting machine. *Annals of statistics* pp. 1189–1232 (2001)
- [17] Friedman, J.H., et al.: Multivariate adaptive regression splines. *The annals of statistics* **19**(1), 1–67 (1991)
- [18] Hooker, G.: Generalized functional anova diagnostics for high-dimensional functions of dependent variables. *J. Comput. Graph. Stat.* **16**(3) (2007)
- [19] Horn, D., Bischl, B.: Multi-objective parameter configuration of machine learning algorithms using model-based optimization. In: 2016 IEEE Symposium Series on Computational Intelligence (SSCI). pp. 1–8. Ieee (2016)
- [20] Huysmans, J., Dejaeger, K., Mues, C., Vanthienen, J., Baesens, B.: An empirical evaluation of the comprehensibility of decision table, tree and rule based predictive models. *Decision Support Systems* **51**(1), 141–154 (2011)
- [21] Knowles, J.: Parego: a hybrid algorithm with on-line landscape approximation for expensive multiobjective optimization problems. *IEEE Transactions on Evolutionary Computation* **10**(1), 50–66 (2006)
- [22] Leeper, T.J.: Interpreting regression results using average marginal effects with r’s margins. Available at the comprehensive R Archive Network (CRAN) (2017)
- [23] Lipton, Z.C.: The mythos of model interpretability. ICML WHI ’16 (2016)
- [24] Molnar, C.: Interpretable Machine Learning (2019), <https://christophm.github.io/interpretable-ml-book/>

- [25] Molnar, C., Bischl, B., Casalicchio, G.: iml: An r package for interpretable machine learning. *JOSS* **3**(26), 786 (2018). <https://doi.org/10.21105/joss.00786>, <http://joss.theoj.org/papers/10.21105/joss.00786>
- [26] Philipp, M., Rusch, T., Hornik, K., Strobl, C.: Measuring the stability of results from supervised statistical learning. *Journal of Computational and Graphical Statistics* **27**(4), 685–700 (2018)
- [27] Plumb, G., Al-Shedivat, M., Xing, E., Talwalkar, A.: Regularizing black-box models for improved interpretability. arXiv preprint arXiv:1902.06787 (2019)
- [28] R Core Team: R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria (2018), <https://www.R-project.org/>
- [29] Ribeiro, M.T., Singh, S., Guestrin, C.: Why should i trust you?: Explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. pp. 1135–1144. ACM (2016)
- [30] Rudin, C.: Please stop explaining black box models for high stakes decisions. arXiv preprint arXiv:1811.10154 (2018)
- [31] Rüping, S., et al.: Learning interpretable models (2006)
- [32] Schielzeth, H.: Simple means to improve the interpretability of regression coefficients. *Methods in Ecology and Evolution* **1**(2), 103–113 (2010)
- [33] Schwarz, G., et al.: Estimating the dimension of a model. *The annals of statistics* **6**(2), 461–464 (1978)
- [34] Stone, C.J.: The use of polynomial splines and their tensor products in multivariate function estimation. *The Annals of Statistics* pp. 118–171 (1994)
- [35] Yang, H., Rudin, C., Seltzer, M.: Scalable bayesian rule lists. In: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. pp. 3921–3930. JMLR. org (2017)
- [36] Zhou, Q., Liao, F., Mou, C., Wang, P.: Measuring interpretability for different types of machine learning models. In: *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. pp. 295–308 (2018)