
Inherent Inconsistencies of Feature Importance

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The rapid advancement and widespread adoption of machine learning-driven
2 technologies have underscored the practical and ethical need for creating in-
3 terpretable artificial intelligence systems. Feature importance, a method that
4 assigns scores to the contribution of individual features on prediction outcomes,
5 seeks to bridge this gap as a tool for enhancing human comprehension of these
6 systems. Feature importance serves as an explanation of predictions in diverse
7 contexts, whether by providing a global interpretation of a phenomenon across
8 the entire dataset or by offering a localized explanation for the outcome of a
9 specific data point. Furthermore, feature importance is being used both for
10 explaining models and for identifying plausible causal relations in the data,
11 independently from the model. However, it is worth noting that these various
12 contexts have traditionally been explored in isolation, with limited theoretical
13 foundations.

14 This paper presents an axiomatic framework designed to establish coherent
15 relationships among the different contexts of feature importance scores. Notably,
16 our work unveils a surprising conclusion: when we combine the proposed
17 properties with those previously outlined in the literature, we demonstrate the
18 existence of an inconsistency. This inconsistency highlights that certain essential
19 properties of feature importance scores cannot coexist harmoniously within a
20 single framework.

21 1 Introduction

22 *Feature Importance* scores gauge the contribution of each feature to an outcome of a model.
23 Most model-agnostic feature importance scores use a two-step process: in the first step, value is
24 assigned to subsets of the features. In the second step, the score of individual features is derived
25 from the values of subsets. This two-step process allows for a discussion about the expected
26 behavior of the value function and the feature importance score. Many feature importance scores
27 have been proposed in the literature: the bivariate-association [1] evaluates a feature's importance
28 based on its conditional attributes, independent of other features, ablation-studies [2, 3, 4]
29 quantify a feature's significance by assessing its contribution when removed from the entire
30 feature set, SHAP [5] computes feature importance as the mean of its contributions across
31 various subsets of features, and MCI [6] determines importance as the maximal contribution
32 among all possible feature subsets (see Table 1). SHAP and MCI use an axiomatic approach in
33 which the expected behaviors are defined as properties, and the functions are derived to satisfy
34 these properties.

35 Feature importance scores can be categorized by two main attributes: the scope, i.e. *local*
36 vs *global*, and the objective, i.e. *data* vs *model*. Methods focusing on local interpretations
37 seek to explain individual predictions (e.g., the role of each feature in a patient's diagnosis [7]).

Table 1: Examples of common feature importance scores. φ denotes the importance function, which takes ν , the value function, and a feature $f \in \mathcal{F}$ as inputs, and assigns an importance score.

Name	Feature importance score: $\varphi(\nu, f)$
Bivariate	$\nu(\{f\})$
Ablation	$\nu(\mathcal{F}) - \nu(\mathcal{F} \setminus \{f\})$
Shapley	$\sum_{S \subseteq \mathcal{F} \setminus \{f\}} \frac{ S !(\mathcal{F} - S -1)!}{ \mathcal{F} !} \cdot (\nu(S \cup \{f\}) - \nu(S))$
MCI	$\max_{S \subseteq \mathcal{F} \setminus \{f\}} (\nu(S \cup \{f\}) - \nu(S))$

38 Conversely, methods focusing on global interpretation try to understand how each feature affects
39 a phenomenon (e.g., the role of each gene in a particular disease [8, 9]). Along the second axis,
40 the data and the model are distinguished by the type of conclusion required. The objective of
41 explaining the data is to infer conclusions about the world that are encoded in the data, as the
42 scientist does in his research [10, 11, 12]. The objective of explaining the model, however, is to
43 use an explanation to monitor and debug a model, to ensure it is working as intended (e.g., as
44 the engineer does for security purposes [13, 14]).

45 Table 2 maps feature importance research according to the local vs. global and data vs. model
46 settings. Most feature importance scores thus far have focused on explaining models, although
47 the data scenario has also been gaining increased attention in recent years. However, the quadrant
48 of the data-local setting is still unexplored in the field of explainable AI. Perhaps this is due to
49 the challenge of providing an accurate explanation as to why a specific outcome (rather than
50 an average result) came into being (rather than being calculated by a model). For example,
51 which characteristic of John Doe is responsible for the fact that he did, or did not, suffer a
52 stroke? These types of questions pertain to individual causal effects that are notoriously difficult
53 to estimate [15, 16].

54 Several studies have examined the relations between the different settings: Lundberg et al. [17]
55 presented a global score that is computed by combining local scores, hence indicating that at
56 least the local and the global settings are not independent. Covert et al. [1] proposed a method
57 of assigning global importance to features, which draws a connection with the local feature
58 importance score of SHAP [5]. Chen et al. [18] defined distinctions between the data and the
59 model and argued that the nature of an explanation depends on what one seeks to explain – the
60 data or the model. Nevertheless, most studies focus solely on one setting. The studies that do
61 consider multiple settings, often do not present an explicit set of expectations for the relations
62 between importance scores under the different settings.

63 In this work, we establish the expected behavior of feature importance scores across diverse
64 contexts. Our objective is to formalize a set of properties that capture the anticipated consistency
65 between local and global interpretations, as well as the alignment between data-driven and
66 model-based assessments. An intriguing extension of this framework is the introduction of the
67 data-local scenario, which, in theory, can be achieved by integrating our properties with axiomatic
68 methods, wherein expected behaviors are rigorously defined, and functions are derived accordingly.
69 In the global-data scenario, we employ a set of properties introduced by Catav et al. [6], implying
70 the MCI importance score. Leveraging our proposed local-global relation, one can derive the
71 expected data-local importance score. Conversely, in the local-model scenario, Lundberg and
72 Lee [5] present the SHAP importance score as the only function that satisfies their proposed
73 set of properties. Here, utilizing our proposed data-model relation, one can similarly obtain the
74 expected data-local importance score. However, to our surprise, in both cases, even a modest set
75 of requirements leads to contradictions.

76 The paper is structured as follows: Section 2 consists of a formulation of the framework that
77 generalizes the two-step process of feature importance to all settings. In Section 3 we focus
78 on the local-global consistency: we present its properties and then demonstrate that they are
79 incompatible with a previous result that defined the data-global setting. At the end of the
80 section, we provide a brief discussion of the nature of the contradiction. Section 4 follows a
81 similar structure as the latter, except addressing the case of the data-model consistency, which

Table 2: Examples of feature importance scores and their categorization according to the global/local and data/model settings.

	Global	Local
Model	Additive-Importance-Measures [1] Bivariate-Association [1] Ablation-Studies [2, 3, 4] FIRM [19] Tree-Shap [17]	SHAP [5] Lime [20] Gradient-Based-Localization [21] Relevance-Propagation [22] TreeExplainer [17]
Data	True-To-Data [18] MCI [6] UMFI [23]	

82 contradicts a previous result that defined the model-local setting. Due to space constraints, we
 83 provide supporting proofs for our claims in the Appendix, along with an extended version of
 84 Theorem 1 that allows a clear demarcation between local and global importance scores.

85 This work makes two key contributions: (1) We introduce a unified axiomatic framework that
 86 encompasses feature importance analysis in diverse settings, including global vs. local and
 87 model vs. data contexts. (2) We rigorously demonstrate inconsistencies within these settings,
 88 shedding light on disparities between global and local interpretations and between model-based
 89 and data-centric evaluations. These findings enhance our understanding of the nuances and
 90 challenges in the theory of feature importance analysis within machine learning interpretability.

91 2 Framework

92 We begin by introducing some notation: the setting consists of an input space \mathcal{X} , an output
 93 space \mathcal{Y} , so that given a pair $(x, y) \sim (\mathcal{X} \times \mathcal{Y})$, the learning task is to predict y by observing x .
 94 Without loss of generality, $\mathcal{X} \subseteq \mathbb{R}^{|\mathcal{F}|}$, where \mathcal{F} is the set of available features. The explanation
 95 task is aimed to assign a score to each feature, based on its contribution to prediction. It consists
 96 of a two-step process: a **value function** is a function $\nu : \{x, \mathcal{X}\} \times 2^{\mathcal{F}} \rightarrow \mathbb{R}$ that assigns a scalar
 97 to each subset of features $S \subseteq \mathcal{F}$, where $\nu(x, S)$ denotes the local value of a subset of features S
 98 for a given pair (x, y) and $\nu(\mathcal{X}, S)$ denotes the global value of this subset. A **feature importance**
 99 **function** is a function $\varphi : 2^{\mathcal{F}} \times \mathcal{F} \rightarrow \mathbb{R}$ that receives the output of a value function and assigns a
 100 feature importance score to each feature. For simplicity, we denote the importance of the feature
 101 f for both the global and the local importance functions as $\varphi(\nu(z), f)$ for $z \in \{x, \mathcal{X}\}$, where the
 102 actual input for the value function differs between them. An elaborated version of this notation
 103 appears in Section C.1 of the Appendix. For the data-model discussion, we add several notations:
 104 the data itself, \mathcal{D} , which is a probability measure over $\mathcal{X} \times \mathcal{Y}$; \mathcal{M} , which is a predictor over \mathcal{D} ;
 105 and $\nu^{\mathcal{D}}, \nu^{\mathcal{M}}$ which are indicators for the current mode of evaluation in the value function.

106 We say that ν is a valid value function if it satisfies: $\nu(\mathcal{X}, \emptyset) = 0$, and in the global setting,
 107 we further require that monotonicity, i.e. for any subset $S \subseteq T \Rightarrow \nu(\mathcal{X}, S) \leq \nu(\mathcal{X}, T)$. This
 108 reflects the intuition that adding features to a model can not decrease the amount of information
 109 regarding the target variable, and thus can not decrease the prediction ability of a model. These
 110 two properties imply that $\forall S \subseteq \mathcal{F}, \nu(\mathcal{X}, S) \geq 0$ in the global setting. We do not assume these
 111 conditions generally hold for the local setting, implying that $\nu(x, S)$ may be negative for some
 112 $x \in \mathcal{X}$. Finally, \mathbb{E} denotes the expected value function with respect to \mathcal{D} .

113 3 The Local-Global relation

114 It is natural to anticipate that a global phenomenon is an aggregate of local phenomena. This
 115 anticipated consistency can be illustrated intuitively: we find it confusing if a model that predicts
 116 loan repayment by lenders considers the age of the lenders to be a crucial factor in the global
 117 sense, yet at the same time declares that age is not a factor in predicting repayment for any
 118 specific lender. To avoid such scenarios, we require a small set of properties that ensure a
 119 meaningful relation between local and global settings in the framework of feature importance.

120 3.1 Expected properties

121 In this section, we formulate two consistency properties that we require to hold between the local
122 and the global settings. We use these properties to prove the first inconsistency theorem.

Property 1 (Value Consistency). ν is Value Consistent if

$$\forall \mathcal{S} \subseteq \mathcal{F}, \nu(\mathcal{X}, \mathcal{S}) = \mathbb{E}[\nu(x, \mathcal{S})]$$

123 In property 1, to establish the relation between the local and the global value functions, the
124 global value of each subset is constrained to be the expectancy taken over the inputs of the local
125 value on this subset.

Property 2 (Importance Consistency). A tuple $\{\nu, \varphi\}$ is Importance Consistent if

$$\forall f \in \mathcal{F}, \varphi(\nu(\mathcal{X}), f) = \mathbb{E}[\varphi(\nu(x), f)]$$

126 In property 2, to establish the local-global relations of the importance score, a consistency
127 requirement analogous to the one above is made for the feature importance function: the global
128 importance of a feature is the expected value of the local feature importance of this feature.

129 The two properties above define the expected relations between feature importance in local and
130 global settings. We say that a tuple $\{\nu, \varphi\}$ is local-global consistent to denote that the *Value*
131 *Consistency* and *Importance Consistency* properties hold.

132 3.2 The local-global inconsistency

We use the MCI function [6] to demonstrate the discrepancy between local and global settings.
This function relies on a pre-defined set of properties which the importance score is expected
to maintain. Apparently, the only function that satisfied these properties is the MCI function,
defined as follows:

$$MCI(\nu, f) = \max_{\mathcal{S} \subseteq \mathcal{F} \setminus \{f\}} (\nu(\mathcal{S} \cup \{f\}) - \nu(\mathcal{S}))$$

133 Remarkably, the MCI score is the only function that uniquely satisfies the MCI properties, detailed
134 in Section A.1. Our analysis leads us to demonstrate the following inconsistency:

135 **Theorem 1.** *properties 1,2, and MCI properties do not hold simultaneously.*

Proof sketch. Let $\{\nu, \varphi\}$ be local-global consistent tuple such that ν is non-decreasing. Assume
that MCI properties hold, i.e. φ is the MCI function. From the local-global consistency, we get
that $\forall f \in \mathcal{F}$:

$$MCI(\mathbb{E}[\nu(x)], f) = MCI(\nu(\mathcal{X}), f) = \mathbb{E}[MCI(\nu(x), f)]$$

136 This leads to a contradiction since MCI uses the max operator and therefore is a non-linear
137 function of the value function. A proof by counterexample is attached in Section B.1. \square

138 The proof sketch presented here is a simplified version, in which the local importance function
139 and the global importance function are identical. A more detailed version of the proof, which
140 does not assume that, can be found in Section C.

141 3.3 Discussion of the local-global relation

142 While the global-data setting is defined by Marginal Contribution Importance (MCI), the local-data
143 setting (as the fourth quadrant in Table 2 demonstrates) is much harder to interpret and define.
144 To tackle this issue, the approach adopted in this study was to use MCI's definition of global-data
145 and define the local-global expected relation. However, this led to an inconsistency theorem. The
146 source of inconsistency lies in the different considerations of ambiguous information: MCI ensures
147 that meaningful information is not missed by attributing the maximum contribution to each
148 feature, regardless of the contribution of other correlated features. This differs from methods
149 such as SHAP [5], where contributions are split between correlated features.

150 4 The Data-Model relation

151 When explaining data, the focus is on understanding the underlying process generating them; while
152 when explaining the model, the focus is on understanding how the model is making predictions

153 based on the data. However, these settings are intertwined – models are often used as proxies by
 154 which nature can be explored. In cases where the model predictions are identical to the data, we
 155 expect conclusions reached from analyzing the model to hold with regard to the data. Therefore,
 156 we expect that the data and model will agree on each feature’s importance. Nonetheless, this
 157 expected property implies a degenerate case where $\nu^{\mathcal{D}} \equiv 0$, which implies that for any importance
 158 function, the importance score of all features becomes zero, rendering them insignificant.

159 4.1 Expected properties

160 In this section, we formulate another consistency property, that expresses the expected relations
 161 between the data and the model settings. Then, we show that fulfillment of this property, along
 162 with known previous results, is only possible in a degenerate case.

Property 3 (Data-Model Consistency). *Let \mathcal{M} be a model that predicts over \mathcal{D} . A tuple $\{\mathcal{D}, \mathcal{M}, \nu, \varphi\}$ is Data-Model Consistent if $\forall x, y \sim \mathcal{D}, \mathcal{M}(x) = y$ and $\forall z \in \{x, \mathcal{X}\}$ it holds that*

$$\forall f \in \mathcal{F}, \quad \varphi(\nu^{\mathcal{D}}(z), f) = \varphi(\nu^{\mathcal{M}}(z), f)$$

163 The *Data-Model Consistency* property 3 states that if a model predicts the target perfectly, then
 164 the data and model importance scores of each feature are identical.

165 4.2 The data-model inconsistency

We use the SHAP function [5] to demonstrate the discrepancy between model and data settings. This function relies on a pre-defined set of properties which the importance score is expected to maintain. Apparently, the only function that satisfied these properties is the SHAP function, defined as follows:

$$\text{SHAP}(\nu, f) = \sum_{\mathcal{S} \subseteq \mathcal{F} \setminus \{f\}} \frac{|\mathcal{S}|!(|\mathcal{F}| - |\mathcal{S}| - 1)!}{|\mathcal{F}|!} \cdot (\nu(\mathcal{S} \cup \{f\}) - \nu(\mathcal{S}))$$

166 Notably, the SHAP score implies additional properties, detailed in Section A.2. This introduction
 167 leads us to the following inconsistency:

168 **Theorem 2.** *If a tuple $\{\mathcal{D}, \mathcal{M}, \nu, \varphi\}$ satisfies Data-Model Consistency (property 3) and SHAP*
 169 *properties, then $\nu^{\mathcal{D}} \equiv 0$.*

170 Our proof is based on a difference between models and the real world. Specifically, when the
 171 data contain correlated features, e.g. height measured in centimeters and inches, a model may
 172 learn based on only one of the features, resulting in different feature importance scores for each
 173 feature in the model. However, in the real world, both features are equally important. A detailed
 174 proof of the theorem is attached in Section B.2.

175 4.3 Discussion of the data-model relation

176 The need to link the data and model settings is not only theoretical. It is motivated by the need
 177 to use models to understand how the world works. Feature importance is often used, even if not
 178 stated explicitly, as a proxy for causal analysis. Unfortunately, the known limitations of trying to
 179 establish causal relations from observational data apply to feature importance too. The example
 180 we used to prove the inconsistency of the data-model often appears in real-world problems. Two
 181 features can be highly similar because a common, unobserved variable, caused them, or one of
 182 them caused the other. For example, when continuously measuring a variable of interest but
 183 only recording its mean and maximum values as observed variables. This problem of lacking
 184 information to disentangle the effect of two variables is known as unidentifiability.

185 To illustrate this in our context, consider two penalized regression models that are trained on
 186 two identical features. The first model employs an L1 regularization (lasso regression), and the
 187 second model employs an L2 regularization (ridge regression). The predictions of the two models
 188 are identical. However, assigning feature importance may lead to different results between the
 189 models - lasso regression will result in assigning all the importance to one of the features, whereas
 190 ridge regression will result in assigning equal importance to both features.

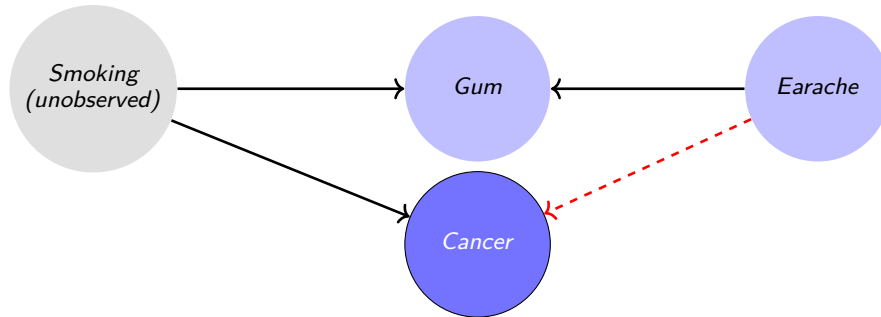


Figure 1: An example of a directed acyclic graph with a collider variable *Gum*.

191 Another situation that may lead to unexpected outcomes from feature importance scores is when
 192 a collider (also known as an inverted fork) exists in the data [15, 16]. For example consider
 193 the situation illustrated in Figure 1: Smoking cigarettes (*Smoking*) causes cancer (*Cancer*),
 194 but also increases chewing gum consumption (*Gum*). Assume also that doctors recommend
 195 people with earaches (*Earache*) to chew gum. Now, imagine scenarios in which a researcher is
 196 developing a model to predict *Cancer* using different subsets of the features *Gum* and *Earache*,
 197 but lacks information on *Smoking*. In the first scenario, the researcher uses only the *Earache*
 198 feature. Since earache and cancer are independent, any value-based feature importance score
 199 will assign zero importance to *Earache*. In the second scenario, where only the *Gum* feature is
 200 present, the researcher will conclude that *Gum* is an important feature since it is correlated with
 201 *Cancer*. In the third scenario, where a model that contains both *Earache* and *Gum* is considered,
 202 the researcher will infer that *Earache* has non-zero importance. This results from conditioning
 203 on *Gum*, creating an association between *Earache* and *Cancer* due to the presence of a collider.
 204 Intuitively, a person who chews gum and does not have an earache is more likely to be a smoker
 205 (notice that the smoking feature is unobserved), and hence at high risk of cancer. Therefore,
 206 the feature importance score might mislead a naïve researcher into thinking that earaches are
 207 predictive of cancer and that gum chewing is a cure for the disease.

208 The situations described here have been studied in the causality literature and there is no recipe
 209 for overcoming them that does not involve additional information about the world [15, 16].

210 5 Conclusion

211 In this work, we investigated the possibility to create a unified framework of feature importance
 212 scores, by defining their expected properties. Surprisingly, we found that it is impossible to define
 213 feature importance scores that are consistent between different settings. Specifically, the expected
 214 consistency between local and global scores contradicts properties of the data-global setting.
 215 Furthermore, there is no guarantee that feature importance scores of a model that perfectly
 216 predicts the data will reflect the feature importance of the data themselves.

217 Our inconsistency result is reminiscent of Kleinberg [24], which proves a similar result for clustering.
 218 Analogously, we do not argue that we have defined the only possible set of relevant properties
 219 for the various settings. We did, however, attempt to define a set of properties that we believe
 220 are essential. Yet, even these requirements led to inconsistencies. Future research can tackle
 221 which further assumptions can be made about feature importance scores, or other explainability
 222 methods, that are meaningful and yet can still be consistent.

223 In the meantime, our results show that feature importance scores should be used cautiously,
 224 aligning with recent research that has attempted to measure the quality and usefulness of
 225 explainability tools for different applications [25, 26, 27]. As such, our work tries to promote
 226 substantive discussions and accurate definitions of explainability, as previously advocated, for
 227 example, by Lipton [28] and Kumar et al. [29]. Hence, we hope that our work will contribute to
 228 stimulating additional research that will result in a solid theoretical foundation for explainable AI.

229 **References**

- 230 [1] Ian Covert, Scott Lundberg, and Su-In Lee. Understanding global feature contributions with
231 additive importance measures. *Advances in Neural Information Processing Systems*, 33,
232 2020.
- 233 [2] VA Casagrande and IT Diamond. Ablation study of the superior colliculus in the tree shrew
234 (tupaia glis). *Journal of Comparative Neurology*, 156(2):207–237, 1974.
- 235 [3] Eric Bengtson and Dan Roth. Understanding the value of features for coreference resolution.
236 In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*,
237 pages 294–303, Honolulu, Hawaii, October 2008. Association for Computational Linguistics.
238 URL <https://www.aclweb.org/anthology/D08-1031>.
- 239 [4] Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will
240 Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining
241 improvements in deep reinforcement learning. In *Thirty-Second AAAI Conference on Artificial
242 Intelligence*, 2018.
- 243 [5] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In
244 *Advances in neural information processing systems*, pages 4765–4774, 2017.
- 245 [6] Amnon Catav, Boyang Fu, Yazeed Zoabi, Ahuva Libi Weiss Meilik, Noam Shomron, Jason
246 Ernst, Sriram Sankararaman, and Ran Gilad-Bachrach. Marginal contribution feature
247 importance—an axiomatic approach for explaining data. In *International Conference on
248 Machine Learning*, pages 1324–1335. PMLR, 2021.
- 249 [7] Todd Kulesza, Margaret Burnett, Weng-Keen Wong, and Simone Stumpf. Principles of
250 explanatory debugging to personalize interactive machine learning. In *Proceedings of the
251 20th international conference on intelligent user interfaces*, pages 126–137, 2015.
- 252 [8] Sarah Tan, Rich Caruana, Giles Hooker, Paul Koch, and Albert Gordo. Learning global
253 additive explanations for neural nets using model distillation. *arXiv preprint arXiv:1801.08640*,
254 2018.
- 255 [9] Juan M Fontana, Muhammad Farooq, and Edward Sazonov. Estimation of feature importance
256 for food intake detection based on random forests classification. In *2013 35th Annual
257 International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*,
258 pages 6756–6759. IEEE, 2013.
- 259 [10] Padideh Danaee, Reza Ghaeini, and David A Hendrix. A deep learning approach for cancer
260 detection and relevant gene identification. In *Pacific symposium on biocomputing*, pages
261 219–229. World Scientific, 2017.
- 262 [11] Brett A McKinney, David M Reif, Marylyn D Ritchie, and Jason H Moore. Machine learning
263 for detecting gene-gene interactions. *Applied bioinformatics*, 5(2):77–88, 2006.
- 264 [12] Kadri Haljas, Azmeraw T Amare, Behrooz Z Alizadeh, Yi-Hsiang Hsu, Thomas Mosley,
265 Anne Newman, Joanne Murabito, Henning Tiemeier, Toshiko Tanaka, Cornelia Van Duijn,
266 et al. Bivariate genome-wide association study of depressive symptoms with type 2 diabetes
267 and quantitative glycemc traits. *Psychosomatic medicine*, 80(3):242, 2018.
- 268 [13] Bojan Karlaš, David Dao, Matteo Interlandi, Bo Li, Sebastian Schelter, Wentao Wu, and
269 Ce Zhang. Data debugging with shapley importance over end-to-end machine learning
270 pipelines. *arXiv preprint arXiv:2204.11131*, 2022.
- 271 [14] Wenbo Guo, Dongliang Mu, Jun Xu, Purui Su, Gang Wang, and Xinyu Xing. Lemna:
272 Explaining deep learning based security applications. In *Proceedings of the 2018 ACM
273 SIGSAC Conference on Computer and Communications Security*, pages 364–379, 2018.
- 274 [15] Judea Pearl. *Causality*. Cambridge university press, 2009.
- 275 [16] Miguel A Hernán and James M Robins. *Causal inference*, 2010.

- 276 [17] Scott M Lundberg, Gabriel Erion, Hugh Chen, Alex DeGrave, Jordan M Prutkin, Bala Nair,
277 Ronit Katz, Jonathan Himmelfarb, Nisha Bansal, and Su-In Lee. From local explanations
278 to global understanding with explainable ai for trees. *Nature machine intelligence*, 2(1):
279 2522–5839, 2020.
- 280 [18] Hugh Chen, Joseph D Janizek, Scott Lundberg, and Su-In Lee. True to the model or true
281 to the data? *arXiv preprint arXiv:2006.16234*, 2020.
- 282 [19] Alexander Zien, Nicole Krämer, Sören Sonnenburg, and Gunnar Rätsch. The feature
283 importance ranking measure. In *Joint European Conference on Machine Learning and*
284 *Knowledge Discovery in Databases*, pages 694–709. Springer, 2009.
- 285 [20] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?"
286 explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD*
287 *international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.
- 288 [21] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi
289 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-
290 based localization. In *Proceedings of the IEEE international conference on computer vision*,
291 pages 618–626, 2017.
- 292 [22] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert
293 Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions
294 by layer-wise relevance propagation. *PLoS one*, 10(7):e0130140, 2015.
- 295 [23] Joseph Janssen and Vincent Guan. Ultra marginal feature importance. *arXiv preprint*
296 *arXiv:2204.09938*, 2022.
- 297 [24] Jon Kleinberg. An impossibility theorem for clustering. *Advances in neural information*
298 *processing systems*, 15, 2002.
- 299 [25] Cynthia Rudin and Joanna Radin. Why are we using black box models in ai when we don't
300 need to? a lesson from an explainable ai competition. 2019.
- 301 [26] Chun-Hao Chang, Sarah Tan, Ben Lengerich, Anna Goldenberg, and Rich Caruana. How
302 interpretable and trustworthy are gams? In *Proceedings of the 27th ACM SIGKDD*
303 *Conference on Knowledge Discovery & Data Mining*, pages 95–105, 2021.
- 304 [27] Giorgio Visani, Enrico Bagli, Federico Chesani, Alessandro Poluzzi, and Davide Capuzzo.
305 Statistical stability indices for lime: Obtaining reliable explanations for machine learning
306 models. *Journal of the Operational Research Society*, 73(1):91–101, 2022.
- 307 [28] Zachary C Lipton. The mythos of model interpretability. *Queue*, 16(3):31–57, 2018.
- 308 [29] I Elizabeth Kumar, Suresh Venkatasubramanian, Carlos Scheidegger, and Sorelle Friedler.
309 Problems with shapley-value-based explanations as feature importance measures. In *Interna-*
310 *tional Conference on Machine Learning*, pages 5491–5500. PMLR, 2020.
- 311 [30] Eric J Friedman. Paths and consistency in additive cost sharing. *International Journal of*
312 *Game Theory*, 32(4):501–518, 2004.
- 313 [31] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks.
314 In *International conference on machine learning*, pages 3319–3328. PMLR, 2017.

315 A Additional Properties

316 In this section, we present the additional properties mentioned in the local-global section and the
317 data-model section.

318 A.1 MCI properties

319 Catav et al. [6] introduced the following three properties to define the expected behavior of
320 feature importance scores in the global-data setting:

321 **Property 4** (Marginal Contribution). *A tuple $\{\nu, \varphi\}$ satisfies the Marginal Contribution property*
322 *when the importance of a feature is equal to or higher than the increase in the value function*
323 *when adding it to all the other features, i.e. $\varphi(\nu(\mathcal{X}, f)) \geq \nu(\mathcal{X}, \mathcal{F}) - \nu(\mathcal{X}, \mathcal{F} \setminus f)$.*

324 The *Marginal Contribution* property states that the importance of a feature is at least its
325 contribution to the value function when adding the latter to the set of all other features.

Property 5 (Elimination). *A tuple $\{\nu, \varphi\}$ satisfies the Elimination property when eliminating*
features from \mathcal{F} can only decrease the importance of each feature. i.e., if $T \subseteq \mathcal{F}$ and $\bar{\nu}$ is the
value function which is obtained by eliminating T from \mathcal{F} then

$$\forall f \in \mathcal{F} \setminus T, \quad \varphi(\nu(\mathcal{X}, f)) \geq \varphi(\bar{\nu}(\mathcal{X}, f))$$

326 The *Elimination* property states that the importance of a feature does not become smaller when
327 other features are removed from the calculation. The process of elimination is defined as follows:

Definition. (Elimination operation) *Let \mathcal{F} be a set of features and ν be a value function.*
Eliminating the set $T \subset \mathcal{F}$ creates a new set of features $\mathcal{F}' = \mathcal{F} \setminus T$ and a new value function
 $\nu' : 2^{\mathcal{F}'} \rightarrow \mathbb{R}$ such that

$$\forall S \subseteq \mathcal{F}', \quad \nu'(S) = \nu(S)$$

Property 6 (Minimalism). *A tuple $\{\nu, \varphi\}$ satisfies the Minimalism property when for every*
function $\bar{\varphi} : \mathbb{R}^{2^{\mathcal{F}}} \rightarrow \mathbb{R}^{\mathcal{F}}$ for which properties 4 and 5 hold, then

$$\forall f \in \mathcal{F}, \quad \varphi(\nu(\mathcal{X}, f)) \leq \bar{\varphi}(\nu(\mathcal{X}, f))$$

328 The *Minimalism* property states that among all the functions that satisfy properties 4 and 5, the
329 feature importance scoring function should be minimal.

330 Using these properties, Catav et al. [6] prove the following:

331 **Theorem 3.** *The MCI feature importance score (see Table 1) is the only score for which the*
332 *Marginal Contribution property, the Elimination property, and the Minimalism property (Properties*
333 *4,5,6) hold simultaneously.*

334 A.2 SHAP properties

335 The following properties stem naturally from the SHAP function, which is the only function that
336 satisfies the SHAP properties proposed in Lundberg and Lee [5]:

337 **Property 7** (Triviality). *A tuple $\{\nu, \varphi\}$ satisfies the Triviality Property if the following conditions*
338 *hold:*

- 339 1. *For all $S \subseteq \mathcal{F}$, if $\nu(x, S) \neq 0$, then there exists a feature $f \in S$ such that $\varphi(\nu(x), f) \neq 0$.*
- 340 2. *If $\varphi(\nu(x), f) \neq 0$, then there exists a subset $S \subseteq \mathcal{F}$ such that $\nu(x, S \cup \{f\}) \neq \nu(x, S)$.*

341 *The Triviality Property establishes a non-trivial relationship between the value and the importance*
342 *functions. It requires that if a subset of features has any value, it will be reflected in the*
343 *importance of at least one feature from this subset. Conversely, it demands that if any feature is*
344 *important (i.e., has non-zero importance), it must be included in some valuable subset. Notably,*
345 *if a feature f satisfies $\nu(x, S) = \nu(x, S \cup \{f\})$ for any subset of features, then f has zero*
346 *importance.*

347 **Property 8** (Dummy Feature). Let \mathcal{M} be a model that predicts over \mathcal{D} . A tuple $\{\nu, \varphi\}$ satisfies
 348 the Dummy Feature Property if, for all $f \in \mathcal{F}$ and for all $x, x' \in \mathcal{X}$ such that x differs from x'
 349 only by the f 'th feature $\mathcal{M}(x) = \mathcal{M}(x')$, then

$$\varphi(\nu^{\mathcal{M}}, f) = 0$$

350 The Dummy Feature Property implies that if changing the value of a feature has no effect on
 351 a model's output, then the importance of that feature is zero. This property also had been
 352 recognized in previous works such as Friedman [30] and Sundararajan et al. [31].

353 B Inconsistencies Proofs

354 In this section, we present proofs for the inconsistency theorems.

355 B.1 Theorem 1

Let $\{\nu, \varphi\}$ be local-global consistent tuple such that ν is non-decreasing. Assume that MCI
 properties hold, i.e. φ is the MCI function. From the local-global consistency, we get that
 $\forall f \in \mathcal{F}$:

$$MCI(\mathbb{E}[\nu(x)], f) = MCI(\nu(\mathcal{X}), f) = \mathbb{E}[MCI(\nu(x), f)]$$

356 This leads to a contradiction since MCI uses the max operator, and therefore is a non-linear
 357 function of the value function.

358 Now, we aim to demonstrate, by way of a counter-example, that the MCI function is not linear.
 359 This will lead to a contradiction between the properties of the data-global, as defined in [6]
 360 setting and $\{\nu, \varphi\}$ Consistency properties. Formally, we contradict the following equality:

361 For any ν which is a valid value function,

$$\alpha \cdot MCI(\nu(x_0)) + (1 - \alpha) \cdot MCI(\nu(x_1)) = MCI(\alpha \cdot \nu(x_0) + (1 - \alpha) \cdot \nu(x_1)) \quad (1)$$

362 **Counter-example.** Let \mathcal{X} be a dataset consisting of two samples: x_0 and x_1 , over the feature
 363 space $\mathcal{F} = \{f_0, f_1\}$. We define the value function ν as follows:

$$\nu = \begin{pmatrix} & x_0 & x_1 & \mathcal{X} \\ \{\emptyset\} : & 0 & 0 & 0 \\ \{f_0\} : & 0 & 1 & 1.5 \\ \{f_1\} : & 1 & 1 & 1 \\ \{f_0, f_1\} : & 2 & 1 & 1.5 \end{pmatrix}$$

364 Now, let $\alpha = \frac{1}{2}$. We will evaluate the left-hand side of equation (1) and the right-hand side
 365 separately.

366 **Left-hand side evaluation:**

$$\begin{aligned} & \alpha \cdot MCI(\nu(x_0)) + (1 - \alpha) \cdot MCI(\nu(x_1)) \\ &= \frac{1}{2} \cdot MCI \left(\begin{pmatrix} 0 \\ 0 \\ 1 \\ 2 \end{pmatrix} \right) + \frac{1}{2} \cdot MCI \left(\begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \end{pmatrix} \right) \\ &= \frac{1}{2} \cdot \begin{pmatrix} 1 \\ 1.5 \end{pmatrix} + \frac{1}{2} \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1.25 \end{pmatrix} \end{aligned}$$

367 **Right-hand side evaluation:**

$$\begin{aligned} & MCI(\alpha \cdot \nu(x_0) + (1 - \alpha) \cdot \nu(x_1)) \\ &= MCI(\nu(\mathcal{X})) \\ &= MCI \left(\begin{pmatrix} 0 \\ 0.5 \\ 1 \\ 1.5 \end{pmatrix} \right) = \begin{pmatrix} 0.5 \\ 1 \end{pmatrix} \end{aligned}$$

368 Hence, we have found a counter-example for equation (1), which contradicts the claimed linearity
 369 of the MCI function. This concludes the proof. \square

370

371 B.2 Theorem 2

372 We establish that any tuple $\{\nu, \varphi\}$ satisfying the Triviality, Dummy-Feature, and Data-Model
 373 Consistency properties (Properties 7, 8, 3) inevitably encounters a scenario where $\nu^{\mathcal{D}} \equiv 0$. This
 374 scenario implies that for any importance function that considers the value, the importance score of
 375 all features becomes zero, rendering them insignificant.

376 **Proof.** Let \mathcal{D} be a probability measure over \mathcal{X} such that its feature space contains two
 377 duplicate features of a random variable, which solely dictate the target. Formally, $\rho \in [0, 1]$,
 378 $\{f_0, f_1\} \subseteq \mathcal{F}$, and $\forall x \in \mathcal{X}$, $f_0(x) = f_1(x) = \rho$. The target is defined as $\mathcal{D}(x) = h(\rho)$, where
 379 h is some function of ρ . Let $\mathcal{M}_0, \mathcal{M}_1$ be two models s.t each model focuses on one feature and
 380 neglects the other:

$$\forall i \in \{0, 1\}, \mathcal{M}_i(x) = h(f_i(x))$$

381 Let the tuple $\{\nu, \varphi\}$ satisfy Triviality, Dummy-Feature and Data-Model Consistency (properties
 382 7, 8, 3). By definition, $\mathcal{M}_0, \mathcal{M}_1$ predict the data perfectly, and therefore by the *Data-Model*
 383 *Consistency*, it holds that

$$\forall i \in \{0, 1\} \text{ and } \forall f \in \mathcal{F}, \varphi(\nu^{\mathcal{D}}, f) = \varphi(\nu^{\mathcal{M}_i}, f)$$

384 The *Dummy Feature* Axiom implies that

$$\forall i \in \{0, 1\}, \varphi(\nu^{\mathcal{M}_i}, f_{1-i}) = 0$$

385 Combining the last two implies that

$$\forall f \in \mathcal{F}, \varphi(\nu^{\mathcal{D}}, f) = 0$$

386 By the *Triviality* Axiom, the only value function that satisfies the above is $\nu^{\mathcal{D}} \equiv 0$. \square

387

388 C Distinguish between the local and global importance functions

389 In this section, we will reformulate our properties to distinguish between the local and global
 390 importance functions.

391 C.1 Framework

392 Before proceeding, we will provide a more detailed definition of the value and importance functions
 393 to ensure precision in describing these functions:

394 A value function is represented as $\nu : (\mathcal{D} \times \{x, \mathcal{X}\} \times 2^{\mathcal{F}}) \rightarrow \mathbb{R}$. It assigns a scalar to each
 395 subset of features $\mathcal{S} \subseteq \mathcal{F}$. Here, $\nu(\mathcal{D}, x, \mathcal{S})$ signifies the value of a feature subset \mathcal{S} for the local
 396 instance x , drawn from a probability measure \mathcal{D} . Additionally, $\nu(\mathcal{D}, \mathcal{X}, \mathcal{S})$ represents the value of
 397 the same feature subset over the entire sample space.

398 On the other hand, a feature importance function is denoted as $\varphi : (\{\text{local, global}\} \times 2^{\mathcal{F}}) \rightarrow \mathbb{R}^{\mathcal{F}}$. It
 399 takes an indicator specifying whether it operates in the local or global context and the output
 400 of the value function. This function assigns feature importance scores to individual features.
 401 Specifically, $\varphi(\text{local}, \nu(\mathcal{D}, x), f)$ indicates the importance of feature f for the instance x , while
 402 $\varphi(\text{global}, \nu(\mathcal{D}, \mathcal{X}), f)$ signifies the importance of the same feature across the entire sample space.

403 Note that the monotonicity property of ν holds in the global setting, but not necessarily in the
 404 local setting. For example, consider the case where the prediction target is whether a person
 405 has cancer and one of the features is whether the person carries a lighter in their pocket. This
 406 feature may be globally important, since it may correlate with smoking. However, it is possible
 407 that some people carry lighters but do not smoke, in which case this feature might lead to an
 408 erroneous prediction and hence has a negative local contribution. Globally admissible denotes
 409 a case where an instance x has only a non-negative contribution. Formally, $\nu(\mathcal{D}, x)$ is globally
 410 admissible if $\nu(\mathcal{D}, x)$ is monotonic non-decreasing and $\nu(\mathcal{D}, x, \emptyset) = 0$.

Figure 2: **consistency diagram**: In local-global consistency the global value is the expectation of the local values, while the global importance is the expected value of local importances.

$$\begin{array}{ccc}
 \nu(x) & \longrightarrow & \varphi(x, \nu(x)) \\
 \downarrow \mathbb{E} & & \downarrow \mathbb{E} \\
 \nu(\mathcal{X}) & \longrightarrow & \varphi(\mathcal{X}, \nu(\mathcal{X}))
 \end{array}$$

411 C.2 Expected properties

412 **Property 9** (Value Consistency). ν is Value Consistent if for every \mathcal{D}

- 413 1. $\forall S \subseteq \mathcal{F}, \quad \nu(\mathcal{D}, \mathcal{X}, S) = \mathbb{E}[\nu(\mathcal{D}, x, S)]$
 414 2. $\exists x^*$ and $\exists \mathcal{D}^*$ such that \mathcal{D}^* is a Dirac measure and $\nu(\mathcal{D}^*, x^*) = \nu(\mathcal{D}, x)$

415 To establish the relation between the local and the global value functions, two complementary
 416 conditions are required: First, the global value of each subset is constrained to be the expectancy
 417 taken over the inputs of the local value on this subset. Second, the local value is constrained to
 418 be able to realize the global value.

419 **Property 10** (Importance Consistency). A tuple $\{\nu, \varphi\}$ is Importance Consistent for every \mathcal{D}

- 420 1. $\forall f \in \mathcal{F}, \quad \varphi(\text{global}, \nu(\mathcal{D}, \mathcal{X}), f) = \mathbb{E}[\varphi(\text{local}, \nu(\mathcal{D}, x), f)]$
 421 2. ν is Value Consistent.

422 To establish the local-global relations of the importance score, a consistency requirement analogous
 423 to the one above is made for the feature importance function: The global importance of a feature
 424 is the expected value of the local feature importance of this feature.

425 Consistency implies a commutative diagram, which is presented in Figure 2.

426 C.3 Detailed proof for Theorem 1

427 The proof uses the following two lemmas. Combining these lemmas implies that φ is a linear
 428 function, and the rest of the proof is identical to the abbreviated version that appears above.

Lemma 1. Let \mathcal{D} be a probability measure over \mathcal{X} . If $\{\nu, \varphi\}$ is importance consistent and $x \in \mathcal{X}$ such that $\nu(\mathcal{D}, x)$ is globally admissible, then

$$\varphi(\text{global}, \nu(\mathcal{D}, x)) = \varphi(\text{local}, \nu(\mathcal{D}, x))$$

429 **Proof.** [of Lemma 1] Let $\{\nu, \varphi\}$ be important consistent tuple. Let \mathcal{D} be a probability measure
 430 over \mathcal{X} and let $x \in \mathcal{X}$ be a globally admissible instance. Denote \mathcal{D}' as the corresponding
 431 probability measure, i.e $\nu(\mathcal{D}, x) = \nu(\mathcal{D}', \mathcal{X}')$. By the Value Consistency property there exist a
 432 Dirac measure \mathcal{D}^* such that $\nu(\mathcal{D}', \mathcal{X}') = \nu(\mathcal{D}^*, x^*)$. Hence,

$$\varphi(\text{global}, \nu(\mathcal{D}, x)) = \varphi(\text{global}, \nu(\mathcal{D}', \mathcal{X}')) \quad (2)$$

$$= \varphi(\text{global}, \nu(\mathcal{D}^*, \mathcal{X}^*)) \quad (3)$$

$$= \varphi(\text{local}, \nu(\mathcal{D}^*, x^*)) \quad (4)$$

$$= \varphi(\text{local}, \nu(\mathcal{D}', x')) \quad (5)$$

$$= \varphi(\text{local}, \nu(\mathcal{D}, x)) \quad (6)$$

433 where (3) is from the global admissibility of $\nu(\mathcal{D}, x)$, (4) follows from the consistency and from
 434 the fact that \mathcal{D}^* is Dirac, and the following equations follow from the definition of \mathcal{D}' . \square

435

Lemma 2. Let \mathcal{D} be a probability measure over \mathcal{X} and let ν be a monotonic non-decreasing value function (i.e. $\forall x \in \mathcal{X}, \nu(\mathcal{D}, x)$ is globally admissible). If $\{\nu, \varphi\}$ is local-global consistent then

$$\varphi(\text{global}, \mathbb{E}[\nu(\mathcal{D}, x)]) = \mathbb{E}[\varphi(\text{global}, \nu(\mathcal{D}, x))]$$

436 **Proof.** [of Lemma 2] Let $\{\nu, \varphi\}$ be a local-global consistent tuple and let \mathcal{D} be such that $\nu(\mathcal{D}, x)$
 437 is globally admissible for every x in the support of \mathcal{D} . Therefore,

$$\mathbb{E}[\varphi_{(\text{global}, \nu(\mathcal{D}, x))}] \tag{7}$$

$$= \mathbb{E}[\varphi_{(\text{local}, \nu(\mathcal{D}, x))}] \tag{8}$$

$$= \varphi_{(\text{global}, \nu(\mathcal{D}, \mathcal{X}))} \tag{9}$$

$$= \varphi_{(\text{global}, \mathbb{E}[\nu(\mathcal{D}, x)])} \tag{10}$$

438 where (8) is valid by Lemma 1, and (9) and (10) are valid by the importance consistency 2. \square
 439