

Evaluating and Aggregating Feature-based Model Explanations

Umang Bhatt^{1,2*}, Adrian Weller^{1,3} and José M. F. Moura²

¹University of Cambridge

²Carnegie Mellon University

³The Alan Turing Institute

{usb20, aw665}@cam.ac.uk, moura@ece.cmu.edu

Abstract

A feature-based model explanation denotes how much each input feature contributes to a model’s output for a given data point. As the number of proposed explanation functions grows, we lack quantitative evaluation criteria to help practitioners know when to use which explanation function. This paper proposes quantitative evaluation criteria for feature-based explanations: low sensitivity, high faithfulness, and low complexity. We devise a framework for aggregating explanation functions. We develop a procedure for learning an aggregate explanation function with lower complexity and then derive a new aggregate Shapley value explanation function that minimizes sensitivity.

1 Introduction

There has been great interest in understanding black-box machine learning models via post-hoc explanations. Much of this work has focused on feature-level importance scores for how much a given input feature contributes to a model’s output. These techniques are popular amongst machine learning scientists who want to sanity check a model before deploying it in the real world [Bhatt *et al.*, 2020]. Many feature-based explanation functions are gradient-based techniques that analyze the gradient flow through a model to determine salient input features [Shrikumar *et al.*, 2017; Sundararajan *et al.*, 2017]. Other explanation functions perturb input values to a reference output and measure the change in the model’s output [Štrumbelj and Kononenko, 2014; Lundberg and Lee, 2017].

With many candidate explanation functions, machine learning practitioners find it difficult to pick which explanation function best captures how a model reaches a specific output for a given input. Though there has been work in qualitatively evaluating feature-based explanation functions on human subjects [Lage *et al.*, 2019], there has been little exploration into formalizing quantitative techniques for evaluating model explanations. Recent work has created auxiliary tasks to test if attribution is assigned to relevant inputs [Yang

and Kim, 2019] and has developed tools to verify if the features important to an explanation function are relevant to the model itself [Camburu *et al.*, 2019].

Borrowing from the humanities, we motivate three criteria for assessing a feature-based explanation: sensitivity, faithfulness, and complexity. Philosophy of science research has advocated for explanations that vary proportionally with changes in the system being explained [Lipton, 2003]; as such, explanation functions should be insensitive to perturbations in the model inputs, especially if the model output does not change. Capturing relevancy faithfully is helpful in an explanation [Ruben, 2015]. Since humans cannot process a lot of information at once, some have argued for minimal model explanations that contain only relevant and representative features [Batterman and Rice, 2014]; therefore, an explanation should not be complex (i.e., use few features).

In this paper, we first define these three distinct criteria: low sensitivity, high faithfulness, and low complexity. With many explanation function choices, we then propose methods for learning an aggregate explanation function that combines explanation functions. If we want to find the simplest explanation from a set of explanations, then we can aggregate explanations to minimize the complexity of the resulting explanation. If we want to learn a smoother explanation function that varies slowly as inputs are perturbed, we can leverage an aggregation scheme that learns a less sensitive explanation function. To the best of our knowledge, we are the first to rigorously explore aggregation of various explanations, while placing explanation evaluation on an objective footing. To that end, we highlight the contributions of this paper:

- We describe three desirable criteria for feature-based explanation functions: low sensitivity, high faithfulness, and low complexity.
- We develop an aggregation framework for combining explanation functions.
- We create two techniques that reduce explanation complexity by aggregating explanation functions.
- We derive an approximation for Shapley-value explanations by aggregating explanations from a point’s nearest neighbors, minimizing explanation sensitivity and resembling how humans reason in medical settings.

*Contact Author

2 Preliminaries

Restricting to supervised classification settings, let f be a black box predictor that maps an input $\mathbf{x} \in \mathbb{R}^d$ to an output $f(\mathbf{x}) \in \mathcal{Y}$. An explanation function g from a family of explanation functions, \mathcal{G} , takes in a predictor f and a point of interest \mathbf{x} and returns importance scores $g(f, \mathbf{x}) = \phi_{\mathbf{x}} \in \mathbb{R}^d$ for all features, where $g(f, \mathbf{x})_i = \phi_{\mathbf{x}, i}$ (simplified to ϕ_i in context) is the importance of (or attribution for) feature x_i of \mathbf{x} . By g_j , we refer to a particular explanation function, usually from a set of explanation functions $\mathcal{G}_m = \{g_1, g_2, \dots, g_m\}$.

We denote $D : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}_{\geq 0}$ to be a distance metric over explanations, while $\rho : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}_{\geq 0}$ denotes a distance metric over the inputs. An evaluation criterion μ takes in a predictor f , explanation function g , and input \mathbf{x} , and outputs a scalar: $\mu(f, g; \mathbf{x})$. $\mathcal{D} = \{(\mathbf{x}^i, y^i)\}_{i=1}^n$ refers to a dataset of input-output pairs, and \mathcal{D}_x denotes all \mathbf{x}^i in \mathcal{D} .

3 Evaluating Explanations

With the number of techniques to develop feature level explanations growing in the explainability literature, picking which explanation function g to use can be difficult. In order to study the aggregation of explanation functions, we define three desiderata of an explanation function g .

3.1 Desideratum: Low Sensitivity

We want to ensure that, if inputs are near each other and their model outputs are similar, then their explanations should be close to each other. Assuming f is differentiable, we desire an explanation function g to have low sensitivity in the region around a point of interest \mathbf{x} , implying local smoothness of g . While [Melis and Jaakkola, 2018] codified the property, [Ghorbani *et al.*, 2019] empirically tested explanation function sensitivity. We follow the convention of the former and define *max sensitivity* and *average sensitivity* in the neighborhood of a point of interest \mathbf{x} .

Let $\mathcal{N}_r = \{z \in \mathcal{D}_x \mid \rho(\mathbf{x}, z) \leq r, f(\mathbf{x}) = f(z)\}$ be a neighborhood of datapoints within a radius r of \mathbf{x} .

Definition 1 (Max Sensitivity). Given a predictor f , an explanation function g , distance metrics D and ρ , a radius r , and a point \mathbf{x} , we define the max sensitivity of g at \mathbf{x} as:

$$\mu_M(f, g, r; \mathbf{x}) = \max_{z \in \mathcal{N}_r} D(g(f, \mathbf{x}), g(f, z))$$

Definition 2 (Average Sensitivity). Given a predictor f , an explanation function g , distance metrics D and ρ , a radius r , a distribution $\mathbb{P}_{\mathbf{x}}(\cdot)$ over the inputs centered at point \mathbf{x} , we define the average sensitivity of g at \mathbf{x} as:

$$\mu_A(f, g, r; \mathbf{x}) = \int_{z \in \mathcal{N}_r} D(g(f, \mathbf{x}), g(f, z)) \mathbb{P}_{\mathbf{x}}(z) dz$$

3.2 Desideratum: High Faithfulness

Faithfulness has been defined in [Yeh *et al.*, 2019]. The feature importance scores from g should correspond to the important features of \mathbf{x} for f ; as such, when we set particular features \mathbf{x}_s to a baseline value $\bar{\mathbf{x}}_s$, the change in predictor's output should be proportional to the sum of attribution scores

of features in \mathbf{x}_s . We measure this as the correlation between the sum of the attributions of \mathbf{x}_s and the difference in output when setting those features to a reference baseline. For a subset of indices $S \subseteq \{1, 2, \dots, d\}$, $\mathbf{x}_s = \{x_i, i \in S\}$ denotes a sub-vector of input features that partitions the input, $\mathbf{x} = \mathbf{x}_s \cup \mathbf{x}_c$. $\mathbf{x}_{[\mathbf{x}_s = \bar{\mathbf{x}}_s]}$ denotes an input where \mathbf{x}_s is set to a reference baseline while \mathbf{x}_c remains unchanged: $\mathbf{x}_{[\mathbf{x}_s = \bar{\mathbf{x}}_s]} = \bar{\mathbf{x}}_s \cup \mathbf{x}_c$. When $|S| = d$, $\mathbf{x}_{[\mathbf{x}_s = \bar{\mathbf{x}}_s]} = \bar{\mathbf{x}}$.

Remark (Reference Baselines). Recent work has discussed how to pick a proper reference baseline $\bar{\mathbf{x}}$. [Sundararajan *et al.*, 2017] suggests using a baseline where $f(\bar{\mathbf{x}}) \approx 0$, while others have proposed taking the baseline to be the mean of the training data. [Chang *et al.*, 2019] notes that the baseline can be learned using generative modeling.

Definition 3 (Faithfulness). Given a predictor f , an explanation function g , a point \mathbf{x} , and a subset size $|S|$, we define the faithfulness of g to f at \mathbf{x} as:

$$\mu_F(f, g; \mathbf{x}) = \text{corr}_{S \in \binom{[d]}{|S|}} \left(\sum_{i \in S} g(f, \mathbf{x})_i, f(\mathbf{x}) - f(\mathbf{x}_{[\mathbf{x}_s = \bar{\mathbf{x}}_s]}) \right)$$

For our experiments, we fix $|S|$ then randomly sample subsets \mathbf{x}_s of the fixed size from \mathbf{x} to estimate correlation. Since we do not see all $\binom{[d]}{|S|}$ subsets in our calculation of faithfulness, we may not get an accurate estimate of the criterion. Though hard to codify and even harder to aggregate, faithfulness is desirable, as it demonstrates that an explanation captures which features the predictor uses to generate an output for a given input. Learning global feature importances that highlight, in expectation, which features a predictor relies on is a challenging problem left to future work.

3.3 Desideratum: Low Complexity

A complex explanation is one that uses all d features in its explanation of which features of \mathbf{x} are important to f . Though this explanation may be faithful to the model (as defined above), it may be too difficult for the user to understand (especially if d is large). We define a fractional contribution distribution, where $|\cdot|$ denotes absolute value:

$$\mathbb{P}_g(i) = \frac{|g(f, \mathbf{x})_i|}{\sum_{j \in [d]} |g(f, \mathbf{x})_j|}; \mathbb{P}_g = \{\mathbb{P}_g(1), \dots, \mathbb{P}_g(d)\}$$

Note that \mathbb{P}_g is a valid probability distribution. Let $\mathbb{P}_g(i)$ denote the fractional contribution of feature x_i to the total magnitude of the attribution. If every feature had equal attribution, the explanation would be complex (even if it is faithful). The simplest explanation would be concentrated on one feature. We define complexity as the entropy of \mathbb{P}_g .

Definition 4 (Complexity). Given a predictor f , explanation function g , and a point \mathbf{x} , the complexity of g at \mathbf{x} is:

$$\mu_C(f, g; \mathbf{x}) = \mathbb{E}_i[-\ln(\mathbb{P}_g)] = -\sum_{i=1}^d \mathbb{P}_g(i) \ln(\mathbb{P}_g(i))$$

4 Aggregating Explanations

Given a trained predictor f , a set of explanation functions $\mathcal{G}_m = \{g_1, \dots, g_m\}$, a criterion to optimize μ , and a set of

inputs \mathcal{D}_x , we want to find an aggregate explanation function \mathbf{g}_{agg} that satisfies μ at least as well as any $\mathbf{g}_i \in \mathcal{G}_m$. Let $h(\cdot)$ represent some function that combines m explanations into a consensus $\mathbf{g}_{\text{agg}} = h(\mathcal{G}_m)$. We now explore different candidates for $h(\cdot)$.

4.1 Convex Combination

Suppose we have two different explanation functions \mathbf{g}_1 and \mathbf{g}_2 and have chosen a criterion μ to evaluate a \mathbf{g} . Consider an aggregate explanation, $\mathbf{g}_{\text{agg}} = h(\mathbf{g}_1, \mathbf{g}_2)$. A potential $h(\cdot)$ is a convex combination where $\mathbf{g}_{\text{agg}} = h(\mathbf{g}_1, \mathbf{g}_2) = w\mathbf{g}_1 + (1-w)\mathbf{g}_2 = \mathbf{w}^\top \mathcal{G}_m$.

Proposition 1. *If D is the ℓ_2 distance and $\mu = \mu_A$ (average sensitivity), the following holds:*

$$\mu_A(\mathbf{g}_{\text{agg}}) \leq w\mu_A(\mathbf{g}_1) + (1-w)\mu_A(\mathbf{g}_2)$$

Proof. Assuming $\mathbb{P}_x(z)$ is uniform, we can apply the triangle inequality and the convexity of D to arrive at the above. \square

A convex combination of explanation functions thus yields an aggregate explanation function that is at most as sensitive as any of the explanation functions taken alone. In order to learn w given \mathbf{g}_1 and \mathbf{g}_2 , we set up an objective as follows.

$$w^* = \arg \min_w \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_x} [\mu_A(\mathbf{g}_{\text{agg}}(\mathbf{f}, \mathbf{x}))] \quad (1)$$

Assuming a uniform distribution around all $\mathbf{x} \in \mathcal{D}_x$, we can rewrite this as:

$$w^* = \arg \min_w \int_{\mathbf{x} \sim \mathcal{D}_x} \int_{\mathbf{z} \in \mathcal{N}_r} D(\mathbf{g}_{\text{agg}}(\mathbf{x}), \mathbf{g}_{\text{agg}}(\mathbf{z})) \mathbb{P}_x(\mathbf{z}) d\mathbf{z} d\mathbf{x}$$

By Cauchy-Schwartz, we get the following:

$$w^* \leq \arg \min_w \int_{\mathbf{x} \sim \mathcal{D}_x} \int_{\mathbf{z} \in \mathcal{N}_r} D(a, b) d\mathbf{z} d\mathbf{x}$$

where $a = w\mathbf{g}_1(\mathbf{f}, \mathbf{x}) + (1-w)\mathbf{g}_2(\mathbf{f}, \mathbf{x})$ and $b = w\mathbf{g}_1(\mathbf{f}, \mathbf{z}) + (1-w)\mathbf{g}_2(\mathbf{f}, \mathbf{z})$. This implies that w^* will be minimal when one element of w^* is 0 and the other is 1. Therefore, a convex combination of two explanation functions, found by solving Equation (1), will be at most as sensitive as the least sensitive explanation function.

4.2 Centroid Aggregation

Another sensible candidate for $h(\cdot)$ to combine m explanation functions is based on centroids with respect to some distance function $D : \mathcal{G} \times \mathcal{G} \mapsto \mathbb{R}$, so that:

$$\mathbf{g}_{\text{agg}} \in \arg \min_{\mathbf{g} \in \mathcal{G}} \mathbb{E}_{\mathbf{g}_i \in \mathcal{G}_m} [D(\mathbf{g}, \mathbf{g}_i)^p] = \arg \min_{\mathbf{g} \in \mathcal{G}} \sum_{\mathbf{g}_i \in \mathcal{G}_m} D(\mathbf{g}, \mathbf{g}_i)^p$$

where p is a positive constant. The simplest examples of distances are the ℓ_2 and ℓ_1 distances with real-valued attributions where $\mathcal{G} \subseteq \mathbb{R}^d$.

Proposition 2. *When D is the ℓ_2 distance and $p = 2$, the aggregate explanation is the feature-wise sample mean.*

$$\mathbf{g}_{\text{agg}}(\mathbf{f}, \mathbf{x}) = \mathbf{g}_{\text{avg}}(\mathbf{f}, \mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \mathbf{g}_i(\mathbf{f}, \mathbf{x}) \quad (2)$$

Proposition 3. *When D is the ℓ_1 distance and $p = 1$, the aggregate explanation is the feature-wise sample median.*

$$\mathbf{g}_{\text{agg}}(\mathbf{f}, \mathbf{x}) = \text{med}\{\mathcal{G}_m\}$$

Propositions 2 and 3 follow from standard results in statistics that the mean minimizes the sum of squared differences and the median minimizes the sum of absolute deviations [Berger, 2013].

We could obtain rank-valued attributions by taking any quantitative vector-valued attributions and ranking features according to their values. If D is the Kendall-tau distance with rank-valued attributions where $\mathcal{G} \subseteq \mathcal{S}_d$ (the set of permutations over d features), then the resulting aggregation mechanism via computing the centroid is called the Kemeny-Young rule. For rank-valued attributions, any aggregation mechanism falls under the rank aggregation problem in social choice theory for which many practical ‘‘voting rules’’ exist [Bhatt *et al.*, 2019a].

We analyze the error of a candidate \mathbf{g}_{agg} . Suppose the optimal explanation for \mathbf{x} using \mathbf{f} is $\mathbf{g}^*(\mathbf{f}, \mathbf{x})$ and suppose \mathbf{g}_{agg} is the mean explanation for \mathbf{x} in Equation (2). Let $\epsilon_{i,\mathbf{x}} = \|\mathbf{g}^*(\mathbf{f}, \mathbf{x}) - \mathbf{g}_i(\mathbf{f}, \mathbf{x})\|$ be the error between the optimal explanation and the i^{th} explanation function.

Proposition 4. *The error between the aggregate explanation $\mathbf{g}_{\text{agg}}(\mathbf{f}, \mathbf{x})$ and the optimal explanation $\mathbf{g}^*(\mathbf{f}, \mathbf{x})$ satisfies:*

$$\epsilon_{\text{agg}} \leq \frac{\sum_{i=1}^n \sum_{j=1}^m \epsilon_{j,\mathbf{x}^i}}{mn}$$

Proof. For a fixed \mathbf{x} , we have:

$$\begin{aligned} \epsilon_{\text{agg},\mathbf{x}} &= \|\mathbf{g}^*(\mathbf{f}, \mathbf{x}) - \mathbf{g}_{\text{agg}}(\mathbf{f}, \mathbf{x})\| \\ &= \left\| \frac{m\mathbf{g}^*(\mathbf{f}, \mathbf{x})}{m} - \frac{1}{m} \sum_{i=1}^m \mathbf{g}_i(\mathbf{f}, \mathbf{x}) \right\| \\ &\leq \frac{1}{m} \sum_{i=1}^m \|\mathbf{g}^*(\mathbf{f}, \mathbf{x}) - \mathbf{g}_i(\mathbf{f}, \mathbf{x})\| = \frac{\sum_{i=1}^m \epsilon_{i,\mathbf{x}}}{m} \end{aligned}$$

Averaging across \mathcal{D}_x , we obtain the result. \square

Hence, by aggregating, we do better than when using one explanation function alone. Many gradient-based explanation functions fit to noise [Hooker *et al.*, 2019]. One way to reduce noise would be to aggregate by ensembling or averaging. As proven in Proposition 4, the typical error of the aggregate is less than the expected error of each function alone.

5 Lowering Complexity Via Aggregation

In this section, we describe iterative algorithms for aggregating explanation functions to obtain $\mathbf{g}_{\text{agg}}(\mathbf{f}, \mathbf{x})$ with lower complexity whilst combining m candidate explanation functions $\mathcal{G}_m = \{\mathbf{g}_1, \dots, \mathbf{g}_m\}$. We desire a $\mathbf{g}_{\text{agg}}(\mathbf{f}, \mathbf{x})$ that contains information from all candidate explanations $\mathbf{g}_i(\mathbf{f}, \mathbf{x})$ yet has entropy less than or equal to that of each explanation $\mathbf{g}_i(\mathbf{f}, \mathbf{x})$. As discussed, a reasonable candidate for an aggregate explanation function is the sample mean given by Equation (2). We may want $\mathbf{g}_{\text{agg}}(\mathbf{f}, \mathbf{x})$ to approach the sample mean, $\mathbf{g}_{\text{avg}}(\mathbf{f}, \mathbf{x})$; however, the sample mean may have greater complexity than that of each $\mathbf{g}_i(\mathbf{f}, \mathbf{x})$.

For example, let $g_1(\mathbf{f}, \mathbf{x}) = [-1, 0]^T$ and $g_2(\mathbf{f}, \mathbf{x}) = [0, 1]^T$. The sample mean is $g_{\text{avg}}(\mathbf{f}, \mathbf{x}) = [-0.5, 0.5]^T$. Both g_1 and g_2 have the minimum possible complexity of 0, while g_{avg} has the maximum possible complexity, $\log(2)$. Our aggregation technique must ensure that $g_{\text{agg}}(\mathbf{f}, \mathbf{x})$ approaches $g_{\text{avg}}(\mathbf{f}, \mathbf{x})$ while guaranteeing $g_{\text{agg}}(\mathbf{f}, \mathbf{x})$ has complexity less than or equal to that of each $g_i(\mathbf{f}, \mathbf{x})$. We now present two approaches for learning a lower complexity explanation, visually represented in Figure 1.

5.1 Gradient-Descent Style Method

Our first approach is similar to gradient descent. Starting from each $g_i(\mathbf{f}, \mathbf{x})$, we iteratively move towards $g_{\text{avg}}(\mathbf{f}, \mathbf{x})$ in each of the d directions (i.e., changing the k th feature by a small amount) if the complexity decreases with that move. We stop moving when the complexity no longer decreases or $g_{\text{avg}}(\mathbf{f}, \mathbf{x})$ is reached. Simultaneously, we start from $g_{\text{avg}}(\mathbf{f}, \mathbf{x})$ and iteratively move towards each $g_i(\mathbf{f}, \mathbf{x})$ in each of the d directions if the complexity decreases. We stop moving when the complexity no longer decreases or any of the $g_i(\mathbf{f}, \mathbf{x})$ are reached. The final $g_{\text{agg}}(\mathbf{f}, \mathbf{x})$ is the location that has the smallest complexity from these $2d$ different walks. Since we only move if the complexity decreases and start from each $g_i(\mathbf{f}, \mathbf{x})$, the entropy of $g_{\text{agg}}(\mathbf{f}, \mathbf{x})$ is guaranteed to be less than or equal to the entropy of all $g_i(\mathbf{f}, \mathbf{x})$.

5.2 Region Shrinking Method

In our second approach, we consider the closed region, \mathbf{R} , which is the convex hull of all the explanation functions, $g_i(\mathbf{f}, \mathbf{x})$. Notice region \mathbf{R} initially contains g_{avg} . We consider an iterative approach to find the global minimum in the region \mathbf{R} . As before, we consider the convex combination formed by two explanation functions, g_i and g_j . Using convex optimization, we find the value on the line segment between g_i and g_j that has the minimum complexity; essentially, we iteratively shrink the region. For the region shrinking method, the convex combination formed by g_i and g_j is:

$$w(g_i) + (1 - w)(g_j), w \in [0, 1]$$

For every pair of functions in \mathcal{G}_m , we find the functions that produces the minimum complexity in the convex combination of the functions, producing a new set of candidates \mathcal{G}'_m . g_{agg} is the element in set \mathcal{G}'_m with minimal complexity after K iterations. In each iteration, a function is chosen if it has the minimum complexity of all the functions in a convex combination. Thus, the minimum complexity of the set \mathcal{G}'_m decreases or remains constant with each iteration.

6 Lowering Sensitivity Via Aggregation

To construct an aggregate explanation function g that minimizes sensitivity, we would need to ensure that a test point's explanation is a function of the explanations of its nearest neighbors under ρ . This is a natural analog for how humans reason: we use past similar events (training data) and facts about the present (individual features) to make decisions [Bhatt *et al.*, 2019b]. We now contribute a new explanation function g_{AVA} that combines the Shapley value explanations of a test point's nearest neighbors to explain the test point.

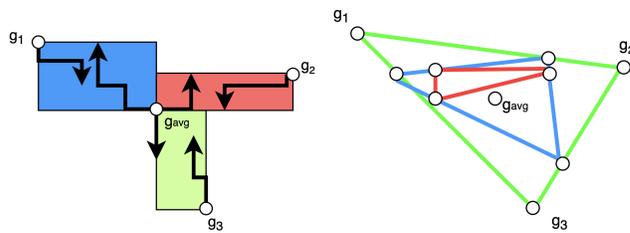


Figure 1: Visual examples of the two complexity lowering aggregation algorithms: gradient-descent style (a) and region shrinking (b) methods using explanation functions g_1, g_2, g_3

6.1 Shapley Value Review

Borrowing from game theory, Shapley values denote the marginal contributions of a player to the payoff of a coalitional game. Let T be the number of players and let $v : 2^T \rightarrow \mathbb{R}$ be the characteristic function, where $v(S)$ denotes the worth (contribution) of the players in $S \subseteq T$. The Shapley value of player i 's contribution (averaging player i 's marginal contributions to all possible subsets S) is:

$$\phi_i(v) = \frac{1}{|T|} \sum_{S \subseteq T \setminus \{i\}} \binom{T-1}{S}^{-1} (v(S \cup \{i\}) - v(S))$$

Let $\Phi \in \mathbb{R}^T$ be a Shapley value contribution vector for all players in the game, where $\phi_i(v)$ is the i th element of Φ .

6.2 Shapley Values as Explanations

In the feature importance literature, we formulate a similar problem to where the game's payoff is the predictor's output $y = \mathbf{f}(\mathbf{x})$, the players are the d features of \mathbf{x} , and the ϕ_i values represent the contribution of x_i to the game $\mathbf{f}(\mathbf{x})$. Let the characteristic function be the importance score of a subset of features x_s , where $\mathbb{E}_Y[\cdot | \mathbf{x}]$ is an expectation over $\mathbb{P}_f(\cdot | \mathbf{x})$:

$$v_x(S) = \mathbb{E}_Y \left[-\log \frac{1}{\mathbb{P}_f(Y | \mathbf{x}_s)} \middle| \mathbf{x} \right]$$

This characteristic function denotes the negative of the expected number of bits required to encode the predictor's output based on the features in a subset S [Chen *et al.*, 2019]. Shapley value contributions can be approximated via Monte Carlo sampling [Štrumbelj and Kononenko, 2014] or via weighted least squares [Lundberg and Lee, 2017].

6.3 Aggregate Valuation of Antecedents

We now explore how to explain a test point in terms of the Shapley value explanations of its neighbors. Termed Aggregate Valuation of Antecedents (AVA), we derive an explanation function that explains a data point in terms of the explanations of its neighbors. We do the following: suppose we want to find an explanation function $g_{\text{AVA}}(\mathbf{f}, \mathbf{x}_{\text{test}})$ for a point of interest \mathbf{x}_{test} . First we find the k nearest neighbors of \mathbf{x}_{test} under ρ denoted by $\mathcal{N}_k(\mathbf{x}_{\text{test}}, \mathcal{D})$.

$$\mathcal{N}_k(\mathbf{x}_{\text{test}}, \mathcal{D}) = \arg \min_{\mathcal{N} \subset \mathcal{D}, |\mathcal{N}|=k} \sum_{\mathbf{z} \in \mathcal{N}} \rho(\mathbf{x}_{\text{test}}, \mathbf{z})$$

We define $\mathbf{g}_{\text{AVA}}(\mathbf{f}, \mathbf{x}_{\text{test}}) = \Phi_{\mathbf{x}_{\text{test}}}$ as the explanation function where:

$$\begin{aligned} \mathbf{g}_{\text{AVA}}(\mathbf{f}, \mathbf{x}_{\text{test}})_i &= \phi_i(v_{\text{AVA}}) = \sum_{\mathbf{z} \in \mathcal{N}_k(\mathbf{x}_{\text{test}})} \frac{\mathbf{g}_{\text{SHAP}}(\mathbf{f}, \mathbf{z})_i}{\rho(\mathbf{x}_{\text{test}}, \mathbf{z})} \\ &= \sum_{\mathbf{z} \in \mathcal{N}_k(\mathbf{x}_{\text{test}})} \frac{\phi_i(v_{\mathbf{z}})}{\rho(\mathbf{x}_{\text{test}}, \mathbf{z})} \end{aligned}$$

In essence, we weight each neighbor’s Shapley value contribution by the inverse distance from the neighbor to the test point. AVA is closely related to bootstrap aggregation from classical statistics, as we take an average of model outputs to improve explanation function stability.

Theorem 5. $\mathbf{g}_{\text{AVA}}(\mathbf{f}, \mathbf{x}_{\text{test}})$ is a Shapley value explanation.

Proof. We want to show that $\mathbf{g}_{\text{AVA}}(\mathbf{f}, \mathbf{x}_{\text{test}}) = \Phi_{\mathbf{x}_{\text{test}}}$ is indeed a vector of Shapley values. Let $\mathbf{g}_{\text{SHAP}}(\mathbf{f}, \mathbf{z}) = \Phi_{\mathbf{z}}$ be the vector of Shapley value contributions for a point $\mathbf{z} \in \mathcal{N}_k$. By [Lundberg and Lee, 2017], we know $\mathbf{g}_{\text{SHAP}}(\mathbf{f}, \mathbf{z})_i = \phi_i(v_{\mathbf{z}})$ is a unique Shapley value for the characteristic function $v_{\mathbf{z}}$. By linearity of Shapley values [Shapley, 1953], we know that:

$$\phi_i(v_{\mathbf{z}_1} + v_{\mathbf{z}_2}) = \phi_i(v_{\mathbf{z}_1}) + \phi_i(v_{\mathbf{z}_2}) \quad (3)$$

This means that the $\Phi_{\mathbf{z}_1} + \Phi_{\mathbf{z}_2}$ will yield a unique Shapley value contribution vector for the characteristic function $v_{\mathbf{z}_1} + v_{\mathbf{z}_2}$. By linearity (or additivity), we know for any scalar α :

$$\alpha \phi_i(v_{\mathbf{z}}) = \phi_i(\alpha v_{\mathbf{z}}) \quad (4)$$

This means $\alpha \Phi_{\mathbf{z}}$ will yield a unique Shapley value contribution vector for the characteristic function $\alpha v_{\mathbf{z}}$. Now define:

$$\Phi_{\mathbf{x}_{\text{test}}} = \sum_{\mathbf{z} \in \mathcal{N}_k(\mathbf{x}_{\text{test}})} \frac{\Phi_{\mathbf{z}}}{\rho(\mathbf{x}_{\text{test}}, \mathbf{z})}$$

We can conclude that $\Phi_{\mathbf{x}_{\text{test}}}$ is a vector of Shapley values. \square

While [Sundararajan *et al.*, 2017] takes a path integral from a fixed reference baseline $\bar{\mathbf{x}}$ and [Lundberg and Lee, 2017] only considers attribution along the straight line path between $\bar{\mathbf{x}}$ and \mathbf{x}_{test} , AVA takes a weighted average of attributions along paths from training points in \mathcal{N}_k to \mathbf{x}_{test} . AVA can similarly be thought of as a convex combination of explanation functions where the explanation functions are the explanations of the neighbors of \mathbf{x}_{test} and the weights are $\rho(\mathbf{x}_{\text{test}}, \mathbf{z})^{-1}$. Though the weights are guaranteed to be non-negative, we normalize the weights to sum to 1 and edit the AVA formulation to be: $\mathbf{g}_{\text{AVA}}(\mathbf{f}, \mathbf{x}_{\text{test}}) = \rho_{\text{tot}} \Phi_{\mathbf{x}_{\text{test}}}$ where $\rho_{\text{tot}} = \sum_{\mathbf{z} \in \mathcal{N}_k(\mathbf{x}_{\text{test}})} \rho(\mathbf{x}_{\text{test}}, \mathbf{z})^{-1}$. Notice this formulation is a specific convex combination as described before; therefore, AVA will result in a lower sensitivity than $\mathbf{g}_{\text{SHAP}}(\mathbf{f}, \mathbf{x})$ alone.

6.4 Medical Connection

Similar to how a model uses input features to reach an output, medical professionals learn how to proactively search for risk predictors in a patient. Medical professionals not only use patient attributes (e.g., vital signs, personal information) to make a diagnosis but also leverage experiences with past patients; for example, if a doctor treated a rare disease over a decade ago, then that experience can be crucial

when attributes alone are uninformative about how to diagnose [Goold and Lipkin Jr, 1999]. This is the analogous to “close” training points affecting a predictor’s output. AVA combines the attributions of past training points (past patients) to explain an unseen test point (current patient). When using the MIMIC dataset [Johnson *et al.*, 2016], AVA models the aforementioned intuition.

7 Experiments

We now report some empirical results. We evaluate models trained on the following datasets: Adult, Iris [Dua and Graff, 2017], MIMIC [Johnson *et al.*, 2016], and MNIST [LeCun *et al.*, 1998]. We use the following explanation functions: SHAP [Lundberg and Lee, 2017], Shapley Sampling (SS) [Štrumbelj and Kononenko, 2014], Gradient Saliency (Grad) [Baehrens *et al.*, 2010], Grad*Input (G*I) [Shrikumar *et al.*, 2017], Integrated Gradients (IG) [Sundararajan *et al.*, 2017], and DeepLift (DL) [Shrikumar *et al.*, 2017].

For all tabular datasets, we train a multilayer perceptron (MLP) with leaky-ReLU activation using the ADAM optimizer. For Iris [Dua and Graff, 2017], we train our model to 96% test accuracy. For Adult [Dua and Graff, 2017], our model has 82% test accuracy. As motivated in Section 6.4, we use MIMIC (Medical Information Mart for Intensive Care III) [Johnson *et al.*, 2016]. We extract seventeen real-valued features deemed critical, per [Purushotham *et al.*, 2018], for sepsis prediction. Our model gets 91% test accuracy on the task. For MNIST [LeCun *et al.*, 1998], our model is a convolutional neural network and has 90% test accuracy.

For experiments with a baseline $\bar{\mathbf{x}}$, zero baseline implies that we set features to 0 and average baseline uses the average feature value in \mathcal{D} . Before doing aggregation, we unit norm all explanations. For the complexity criterion, we take the positive ℓ_1 norm. We set $D = \ell_2$ and $\rho = \ell_\infty$.

7.1 Faithfulness μ_{F}

In Table 2, we report results for faithfulness for various explanation functions. When evaluating, we take the average of multiple runs where, in each run, we see at least 50 datapoints; for each datapoint, we randomly select $|S|$ features and replace them with baseline values. We then calculate the Pearson’s correlation coefficient between the predicted logits of each modified test point and the average explanation attribution for only the subset of features. We notice that, as subset size increases, faithfulness increases until the subset is large enough to contain all informative features. We find that Shapley values, approximated with weighted least squares, are the most faithful explanation function for smaller datasets.

7.2 Max and Avg Sensitivity μ_{M} and μ_{A}

In Table 3, we report the max and average sensitivities for various explanation functions. To evaluate the sensitivity criterion, we sample a set of test points from \mathcal{D} and an additional larger set of training points. We then find the training points that fall within a radius r neighborhood of each test point and find the distance between each nearby training point explanation and the test point explanation to get a mean and max. We average over ten random runs of this procedure. Sensitivity is

INPUT	BEST (DEEPLIFT)	CONVEX	GRADIENT-DESCENT	REGION-SHRINKING
	$\mu_C = 3.688$	$\mu_C = 3.685$	$\mu_C = 3.575$	$\mu_C = 3.208$

Table 1: Qualitative example of aggregation to lower complexity (μ_C): We show that it is possible to lower complexity slightly with both of our approaches; note that achieving lowest complexity on an image would imply that all attribution is placed on a single pixel.

METHOD SUBSET	ADULT 2	IRIS 2	MIMIC 10	MIMIC 20
SHAP	(62, 60)	(67, 68)	(31, 36)	(37, 47)
SS	(46, 27)	(32, 36)	(59, 58)	(38, 45)
GRAD	(30, 53)	(14, 16)	(37, 41)	(28, 63)
G*I	(38, 39)	(27, 30)	(54, 48)	(59, 43)
IG	(47, 33)	(60, 57)	(66, 51)	(68, 51)
DL	(58, 43)	(46, 48)	(84, 54)	(43, 45)

Table 2: Faithfulness μ_F averaged over a test set: (Zero Baseline, Training Average Baseline). Exact quantities can be obtained by dividing table entries by 10^2

METHOD RADIUS	ADULT 2	IRIS 0.2	MIMIC 4
SHAP	(60, 54)	(310, 287)	(6, 5)
SS	(191, 168)	(477, 345)	(83, 81)
GRAD	(60, 50)	(68, 66)	(28, 28)
G*I	(86, 71)	(298, 279)	(77, 50)
IG	(19, 17)	(495, 462)	(19, 15)
DL	(74, 74)	(850, 820)	(135, 111)

Table 3: Sensitivity: (Max μ_M , Avg μ_A). Exact quantities can be obtained by dividing table entries by 10^3

highly dependent on the dimensionality d and on the radius r . We find that as sensitivity decreases as r increases. Empirically, for MIMIC, Shapley values approximated by weighted least squares (SHAP) are the least sensitive.

7.3 MNIST Complexity μ_C

In Table 1, we provide a qualitative example for the gradient descent-style and region-shrinking methods for lowering complexity of explanations from a model trained on MNIST. We show an example with images since it illustrates the notion of lower complexity well; however, other data types (tabular) might be better suited for complexity optimization.

7.4 AVA

Our empirical findings support use of an AVA explanation if low sensitivity is desired. [Ghorbani *et al.*, 2019] note that perturbation-based explanations (like g_{SHAP}) are less sensitive than their gradient-based counterparts. In Table 4, we show that AVA explanations not only have lower sensitivities in all experiments but also have less complex explanations (depending on the radius r and number of features d). After

METHOD	ADULT	IRIS	MIMIC
$\mu_A(\mathbf{f}, g_{\text{SHAP}})$	0.16 ± 0.11	0.22 ± 0.25	0.47 ± 0.12
$\mu_A(\mathbf{f}, g_{\text{AVA}})$	0.07 ± 0.07	0.13 ± 0.18	0.31 ± 0.13
$\mu_M(\mathbf{f}, g_{\text{SHAP}})$	0.68 ± 0.13	1.20 ± 0.36	0.83 ± 0.17
$\mu_M(\mathbf{f}, g_{\text{AVA}})$	0.52 ± 0.11	1.18 ± 0.28	0.72 ± 0.22
$\mu_C(\mathbf{f}, g_{\text{SHAP}})$	1.94 ± 0.26	1.36 ± 0.36	2.33 ± 0.23
$\mu_C(\mathbf{f}, g_{\text{AVA}})$	1.93 ± 0.24	1.24 ± 0.32	2.61 ± 0.29

Table 4: AVA lowers the sensitivity of Shapley value explanations across all datasets. When d is small (fewer features), AVA explanations are slightly less complex.

finding the average distance between pairs of points, we use $r = 1$ for Adult, $r = 0.3$ for Iris, and $r = 10$ for MIMIC.

8 Conclusion

Borrowing from earlier work in social science and the philosophy of science, we codify low sensitivity, high faithfulness, and low complexity as three desirable properties of explanation functions. We define these three properties for feature-based explanation functions, develop an aggregation scheme for learning combinations of various explanation functions, and devise schemes to learn explanations with lower complexity (iterative approaches) and lower sensitivity (AVA). We hope that this work will provide practitioners with a principled way to evaluate feature-based explanations and to learn an explanation which aggregates and optimizes for criteria desired by end users. Though we consider one criterion at a time, future work could further axiomatize our criteria, explore the interaction between different evaluation criteria, and devise a multi-objective optimization approach to finding a desirable explanation; for example, can we develop a procedure for learning a less sensitive and less complex explanation function simultaneously?

Acknowledgements

We thank reviewers for their feedback. We thank Pradeep Ravikumar, John Shi, Brian Davis, Kathleen Ruan, Javier Anoran, James Allingham, and Adithya Raghuraman for their comments and help. UB acknowledges support from DeepMind and the Leverhulme Trust via the Leverhulme Center for the Future of Intelligence (CFI) and from the Partnership on AI. AW acknowledges support from the David MacKay Newton Research Fellowship at Darwin College, The Alan Turing Institute under EPSRC grant EP/N510129/1 & TU/B/000074, and the Leverhulme Trust via the CFI.

References

- [Adebayo *et al.*, 2018] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages 9505–9515, 2018.
- [Ancona *et al.*, 2018] Marco Ancona, Enea Ceolini, Cengiz Oztireli, and Markus Gross. Towards better understanding of gradient-based attribution methods for deep neural networks. In *6th International Conference on Learning Representations (ICLR 2018)*, 2018.
- [Baehrens *et al.*, 2010] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Muller. How to explain individual classification decisions. *JMLR*, 11(Jun):1803–1831, 2010.
- [Batterman and Rice, 2014] Robert W Batterman and Collin C Rice. Minimal model explanations. *Philosophy of Science*, 81(3):349–376, 2014.
- [Berger, 2013] James O Berger. *Statistical decision theory and Bayesian analysis*. Springer Science & Business Media, 2013.
- [Bhatt *et al.*, 2019a] Umang Bhatt, Pradeep Ravikumar, et al. Building human-machine trust via interpretability. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9919–9920, 2019.
- [Bhatt *et al.*, 2019b] Umang Bhatt, Pradeep Ravikumar, and José M. F. Moura. Towards aggregating weighted feature attributions. *arXiv:1901.10040*, 2019.
- [Bhatt *et al.*, 2020] Umang Bhatt, Alice Xiang, Shubham Sharma, Adrian Weller, Ankur Taly, Yunhan Jia, Joydeep Ghosh, Ruchir Puri, José M. F. Moura, and Peter Eckersley. Explainable machine learning in deployment. *ACM Conference on Fairness, Accountability, and Transparency (FAT*)*, 2020.
- [Bylinskii *et al.*, 2018] Zoya Bylinskii, Tilke Judd, Aude Oliva, Antonio Torralba, and Frédo Durand. What do different evaluation metrics tell us about saliency models? *IEEE transactions on pattern analysis and machine intelligence*, 41(3):740–757, 2018.
- [Camburu *et al.*, 2019] Oana-Maria Camburu, Eleonora Giunchiglia, Jakob Foerster, Thomas Lukasiewicz, and Phil Blunsom. Can I trust the explainer? Verifying post-hoc explanatory methods. *arXiv:1910.02065*, 2019.
- [Carter *et al.*, 2019] Brandon Carter, Jonas Mueller, Siddhartha Jain, and David Gifford. What made you do this? understanding black-box decisions with sufficient input subsets. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 567–576, 2019.
- [Chang *et al.*, 2019] Chun-Hao Chang, Elliot Creager, Anna Goldenberg, and David Duvenaud. Explaining image classifiers by counterfactual generation. In *International Conference on Learning Representations*, 2019.
- [Chen *et al.*, 2018] Jianbo Chen, Le Song, Martin Wainwright, and Michael Jordan. Learning to explain: An information-theoretic perspective on model interpretation. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 883–892, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- [Chen *et al.*, 2019] Jianbo Chen, Le Song, Martin J Wainwright, and Michael I Jordan. L-shapley and C-shapley: Efficient model interpretation for structured data. *International Conference on Learning Representations*, 2019.
- [Davis *et al.*, 2020] B. Davis, U. Bhatt, K. Bhardwaj, R. Marculescu, and J. M. F. Moura. On network science and mutual information for explaining deep neural networks. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8399–8403, 2020.
- [Dua and Graff, 2017] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [Ghorbani *et al.*, 2019] Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3681–3688, 2019.
- [Gilpin *et al.*, 2018] Leilani H Gilpin, David Bau, Ben Z Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal. Explaining explanations: An overview of interpretability of machine learning. In *2018 IEEE 5th International Conference on data science and advanced analytics (DSAA)*, pages 80–89. IEEE, 2018.
- [Goold and Lipkin Jr, 1999] Susan Dorr Goold and Mack Lipkin Jr. The doctor–patient relationship: challenges, opportunities, and strategies. *Journal of general internal medicine*, 14(Suppl 1):S26, 1999.
- [Grabska-Barwińska, 2020] Agnieszka Grabska-Barwińska. Measuring and improving the quality of visual explanations. *arXiv preprint arXiv:2003.08774*, 2020.
- [Hazard *et al.*, 2019] Christopher J Hazard, Christopher Fusting, Michael Resnick, Michael Auerbach, Michael Meehan, and Valeri Korobov. Natively interpretable machine learning and artificial intelligence: Preliminary results and future directions. *arXiv preprint arXiv:1901.00246*, 2019.
- [Hind *et al.*, 2019] Michael Hind, Dennis Wei, Murray Campbell, Noel CF Codella, Amit Dhurandhar, Aleksandra Mojsilović, Karthikeyan Natesan Ramamurthy, and Kush R Varshney. Ted: Teaching ai to explain its decisions. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 123–129, 2019.
- [Honegger, 2018] Milo Honegger. Shedding Light on Black Box Algorithms. Master’s thesis, Karlsruhe Institute of Technology, Germany, 2018.
- [Hooker *et al.*, 2019] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability methods in deep neural networks. In *Advances in Neural Information Processing Systems*, pages 9734–9745, 2019.

- [Johnson *et al.*, 2016] Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific Data*, 2016.
- [Kindermans *et al.*, 2019] Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. The (un) reliability of saliency methods. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pages 267–280. Springer, 2019.
- [Lage *et al.*, 2019] Isaac Lage, Emily Chen, Jeffrey He, Menaka Narayanan, Been Kim, Samuel J Gershman, and Finale Doshi-Velez. Human evaluation of models built for interpretability. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 7, pages 59–67, 2019.
- [LeCun *et al.*, 1998] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [Lipton, 2003] Peter Lipton. *Inference to the best explanation*. Routledge, 2003.
- [Lundberg and Lee, 2017] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems 30 (NeurIPS 2017)*, pages 4765–4774, 2017.
- [Melis and Jaakkola, 2018] David Alvarez Melis and Tommi Jaakkola. Towards robust interpretability with self-explaining neural networks. In *Advances in Neural Information Processing Systems (NeurIPS 2018)*, 2018.
- [Montavon *et al.*, 2018] Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*, 73:1–15, 2018.
- [Osman *et al.*, 2020] Ahmed Osman, Leila Arras, and Wojciech Samek. Towards ground truth evaluation of visual explanations. *arXiv preprint arXiv:2003.07258*, 2020.
- [Plumb *et al.*, 2018] Gregory Plumb, Denali Molitor, and Ameet S Talwalkar. Model agnostic supervised local explanations. In *Advances in Neural Information Processing Systems*, pages 2515–2524, 2018.
- [Poursabzi-Sangdeh *et al.*, 2018] Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Vaughan, and Hanna Wallach. Manipulating and measuring model interpretability. *arXiv preprint arXiv:1802.07810*, 2018.
- [Purushotham *et al.*, 2018] Sanjay Purushotham, Chuizheng Meng, Zhengping Che, and Yan Liu. Benchmarking deep learning models on large healthcare datasets. *Journal of Biomedical Informatics*, 83:112–134, 2018.
- [Ribeiro *et al.*, 2016] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 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*. ACM, 2016.
- [Rieger and Hansen, 2020] Laura Rieger and Lars Kai Hansen. Irof: a low resource evaluation metric for explanation methods. *arXiv preprint arXiv:2003.08747*, 2020.
- [Ruben, 2015] David-Hillel Ruben. *Explaining explanation*. Routledge, 2015.
- [Samek *et al.*, 2016] Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE transactions on neural networks and learning systems*, 28(11):2660–2673, 2016.
- [Shapley, 1953] Lloyd S Shapley. A value for n-person games. In *Contributions to the Theory of Games II*, pages 307–317, 1953.
- [Shrikumar *et al.*, 2017] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70 (ICML 2017)*, pages 3145–3153. Journal of Machine Learning Research, 2017.
- [Štrumbelj and Kononenko, 2014] Erik Štrumbelj and Igor Kononenko. Explaining prediction models and individual predictions with feature contributions. *Knowledge and Information Systems*, 41(3):647–665, 2014.
- [Sundararajan *et al.*, 2017] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70 (ICML 2017)*, pages 3319–3328. Journal of Machine Learning Research, 2017.
- [Wang *et al.*, 2020] Zifan Wang, Piotr Mardziel, Anupam Datta, and Matt Fredrikson. Interpreting interpretations: Organizing attribution methods by criteria. *arXiv preprint arXiv:2002.07985*, 2020.
- [Warnecke *et al.*, 2019] Alexander Warnecke, Daniel Arp, Christian Wressnegger, and Konrad Rieck. Evaluating explanation methods for deep learning in security. *arXiv preprint arXiv:1906.02108*, 2019.
- [Yang and Kim, 2019] Mengjiao Yang and Been Kim. BIM: Towards quantitative evaluation of interpretability methods with ground truth. *arXiv:1907.09701*, 2019.
- [Yang *et al.*, 2019] Fan Yang, Mengnan Du, and Xia Hu. Evaluating explanation without ground truth in interpretable machine learning. *arXiv preprint arXiv:1907.06831*, 2019.
- [Yeh *et al.*, 2019] Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Suggala, David I Inouye, and Pradeep K Ravikumar. On the (in) fidelity and sensitivity of explanations. In *Advances in Neural Information Processing Systems*, pages 10965–10976, 2019.
- [Zhang *et al.*, 2019] Hao Zhang, Jiayi Chen, Haotian Xue, and Quanshi Zhang. Towards a unified evaluation of explanation methods without ground truth. *arXiv preprint arXiv:1911.09017*, 2019.

A Additional Evaluation Criteria

In addition to the aforementioned three criteria, there are many other desirable criteria for a g . To assist practitioners, we now collect and list these additional quantitative evaluation criteria for feature-level explanations. It is possible to evaluate all criteria for both perturbation-based explanations [Štrumbelj and Kononenko, 2014; Lundberg and Lee, 2017] and gradient-based explanations [Sundararajan *et al.*, 2017; Shrikumar *et al.*, 2017]. Note we omit evaluation criteria that assume access to ground-truth explanations for training points; for a thorough treatment on this topic, see [Hind *et al.*, 2019; Osman *et al.*, 2020]. We do not delve into human-centered evaluation of explanation functions either; see [Gilpin *et al.*, 2018; Poursabzi-Sangdeh *et al.*, 2018; Yang *et al.*, 2019] for detailed discussions.

Predictability of Explanations

We would want to ensure that explanations from g are predictable. As such, $g(\mathbf{f}, \mathbf{x})$ ought not vary over function calls. [Honegger, 2018] notes that identical inputs should give the identical explanations.

Definition 5 (Identity). Given a predictor \mathbf{f} , an explanation function g , and distance metrics D and ρ , we define the identity criterion for g on \mathcal{D} as:

$$\begin{aligned} \mu_{\text{IDENTITY}}(\mathbf{f}, g) &= \mathbb{E}_{\mathbf{x} \in \mathcal{D}_x} [D(g(\mathbf{f}, \mathbf{x}), g(\mathbf{f}, \mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_x} [\|g(\mathbf{f}, \mathbf{x}) - g(\mathbf{f}, \mathbf{x})\|_0] \end{aligned}$$

Note the above two are equivalent and we take the ℓ_0 norm of the difference between two separate calls to g with the same input \mathbf{x} . The identity criterion favors non-stochastic explanation functions. We would want to ensure that any non-identical inputs should have non-identical explanations.

Definition 6 (Separability). Given a predictor \mathbf{f} , an explanation function g , and distance metrics D and ρ , we define the separability of g on \mathcal{D} as:

$$\begin{aligned} \mu_{\text{SEP}}(\mathbf{f}, g) &= \mathbb{E}_{\mathbf{x}, z \in \mathcal{D}_x, \mathbf{x} \neq z} [D(g(\mathbf{f}, \mathbf{x}), g(\mathbf{f}, z))] \\ &= \mathbb{E}_{\mathbf{x}, z \in \mathcal{D}_x, \mathbf{x} \neq z} [\|g(\mathbf{f}, \mathbf{x}) - g(\mathbf{f}, z)\|_0] \end{aligned}$$

We would also want to know how surprising an explanation $g(\mathbf{f}, \mathbf{x})$ is compared to explanations for training data. [Hazard *et al.*, 2019] defines conviction of an input \mathbf{x} with respect to \mathcal{D}_x for k -Nearest Neighbor algorithms; similarly, we define the conviction of $g(\mathbf{f}, \mathbf{x})$ to explanations of training points, \mathcal{D}_x , using g .

Definition 7 (Conviction). Given a predictor \mathbf{f} , an explanation function g , a probability distribution over explanations $\mathbb{P}_\phi(\cdot)$, and a data point \mathbf{x} , we define the conviction of g at \mathbf{x} for \mathcal{D} as:

$$\mu_{\text{CON}}(\mathbf{f}, g, \mathbb{P}_\phi; \mathbf{x}) = \frac{\mathbb{E}_{z \sim \mathcal{D}_x} [I(g(\mathbf{f}, z))]}{I(g(\mathbf{f}, \mathbf{x}))}$$

where $I(g(\mathbf{f}, \mathbf{x})) = -\ln(\mathbb{P}_\phi(g(\mathbf{f}, \mathbf{x})))$

$\mu_{\text{CON}} = 0$ means that $g(\mathbf{f}, \mathbf{x})$ is surprising. As $\mu_{\text{CON}} \rightarrow \infty$, $g(\mathbf{x})$ contains an expected amount of surprisal and can reasonably occur. We desire a higher μ_{CON} , implying that g behaves predictably. By changing the distribution

to $\mathbb{P}_\phi(\cdot | y = \mathbf{f}(\mathbf{x}))$, the numerator to conditional entropy where $\mathbf{f}(z) = \mathbf{f}(\mathbf{x})$, and self-information to $I(g(\mathbf{f}, \mathbf{x})) = -\ln(\mathbb{P}_\phi(g(\mathbf{f}, \mathbf{x}) | y = \mathbf{f}(\mathbf{x})))$, we define the *conditional conviction* of $g(\mathbf{f}, \mathbf{x})$ to explanations of the same predicted class.

Other techniques have also argued that $g(\mathbf{f}, \mathbf{x})$ should recover the output of the original predictor, $\mathbf{f}(\mathbf{x})$. Deemed compatibility, this criterion attempts to use g as a simple proxy for reproducing the outputs of the complex \mathbf{f} .

Definition 8 (Compatibility). Given a predictor \mathbf{f} and an explanation function g , we define the completeness of g for a dataset \mathcal{D} as:

$$\mu_{\text{COM}}(\mathbf{f}, g) = \frac{1}{N} \sum_{\mathbf{x} \in \mathcal{D}_x} \left| \left(\sum_{i=1}^d g(\mathbf{f}, \mathbf{x})_i \right) - \mathbf{f}(\mathbf{x}) \right|$$

The closer μ_{COM} is to 0, the more compatible the explanation function is; that is, the explanation function recovers the complex model’s outputs well. This criterion is related to the *completeness* axiom of some explanation functions [Sundararajan *et al.*, 2017]. An explanation functions can be built to be compatible with the original model (or complete with respect to \mathbf{f}). This is also related to the notion of *post-hoc accuracy* discussed in [Chen *et al.*, 2018].

Importance of Explanations

Not only do we want to ensure that the g faithfully identifies the most important features, but we also want to understand how well \mathbf{f} performs when \mathbf{x}_s is unobserved (or set to a baseline $\mathbf{x}_s = \bar{\mathbf{x}}_s$). In particular, we craft S to contain the indices for the $|S|$ features with the highest $\text{abs}(g(\mathbf{f}, \mathbf{x})_i)$.

$$S = \arg \max_{S \subset [d], |S|=k} \sum_{i \in S} \text{abs}(g(\mathbf{f}, \mathbf{x})_i)$$

Therefore, \mathbf{x}_s is now a sub-vector of the most important features according to a specific g . As done in [Chang *et al.*, 2019], we define a score s_f for how confidently \mathbf{f} predicts an output y in terms of log-odds.

$$s_f(y|\mathbf{x}) = \log(\hat{\mathbb{P}}_f(y|\mathbf{x})) - \log(1 - \hat{\mathbb{P}}_f(y|\mathbf{x}))$$

Definition 9 (Deletion). Given a predictor \mathbf{f} , an explanation function g , a point of interest \mathbf{x} , a predicted output y , and a subset of important features S , we define the deletion score for \mathbf{f} at \mathbf{x} as:

$$\mu_{\text{DEL}}(\mathbf{f}, g; \mathbf{x}, y) = s_f(y|\mathbf{x}) - s_f(y|\mathbf{x}_{[x_s=\bar{x}_s]})$$

Definition 10 (Addition). Given a predictor \mathbf{f} , an explanation function g , a point of interest \mathbf{x} , a predicted output y , and a subset of important features S , we define the addition score for \mathbf{f} at \mathbf{x} as:

$$\mu_{\text{ADD}}(\mathbf{f}, g; \mathbf{x}, y) = s_f(y|\mathbf{x}_{[x_s=\bar{x}_s]}) - s_f(y|\bar{\mathbf{x}})$$

While the deletion score conveys how the log-odds change when we delete the subset of important features from \mathbf{x} , the addition score tells us how much the log-odds change when we add the subset to the baseline. Instead of re-scoring (via change in log-odds) a modified input like $\mathbf{x}_{[x_s=\bar{x}_s]}$, we can retrain the predictor \mathbf{f} based on a dataset of modified inputs $\mathcal{D}_{\mathbf{x}_{[x_s=\bar{x}_s]}}$. Addition and Deletion are closely related to *explanation selectivity*, described in [Montavon *et al.*, 2018].

Let $\mathbf{f}_{\bar{x}_s}$ denote the predictor trained on the modified inputs with the most important pixels removed. As in [Hooker *et al.*, 2019], we define the ROAR score as the difference in accuracy between the original predictor and the modified predictor. We can also train a predictor where the least important features (those in \mathbf{x}_c) are removed. We denote that predictor to be $\mathbf{f}_{\bar{x}_c}$ and define a KAR score, as proposed in [Hooker *et al.*, 2019].

Definition 11 (ROAR). Given a predictor \mathbf{f} , an explanation function \mathbf{g} , a modified predictor $\mathbf{f}_{\bar{x}_s}$, and a subset of important features S , we define the ROAR score for \mathbf{g} on a dataset \mathcal{D} as:

$$\mu_{\text{ROAR}}(\mathbf{f}, \mathbf{g}, \mathbf{f}_{\bar{x}_s}) = \frac{1}{N} \sum_{\mathbf{x} \in \mathcal{D}_x} \mathbb{1}[\mathbf{f}(\mathbf{x}) = y] - \mathbb{1}[\mathbf{f}_{\bar{x}_s}(\mathbf{x}) = y]$$

Definition 12 (KAR). Given a predictor \mathbf{f} , an explanation function \mathbf{g} , a modified predictor $\mathbf{f}_{\bar{x}_c}$, and a subset of important features S , we define the KAR score for \mathbf{g} on a dataset \mathcal{D} as:

$$\mu_{\text{KAR}}(\mathbf{f}, \mathbf{g}, \mathbf{f}_{\bar{x}_c}) = \frac{1}{N} \sum_{\mathbf{x} \in \mathcal{D}_x} \mathbb{1}[\mathbf{f}(\mathbf{x}) = y] - \mathbb{1}[\mathbf{f}_{\bar{x}_c}(\mathbf{x}) = y]$$

Other Connections

We can also draw parallels between the three criteria proposed in the main paper and existing criteria in the literature.

Low sensitivity is discussed as *stability* in [Melis and Jaakkola, 2018], as *explanation continuity* in [Montavon *et al.*, 2018], as *sensitivity* in [Yeh *et al.*, 2019], and as *reliability* in [Kindermans *et al.*, 2019].

High faithfulness appears as *relevance* in [Samek *et al.*, 2016], as *gold set* in [Ribeiro *et al.*, 2016], as *faithfulness* in [Plumb *et al.*, 2018], as *sensitivity-n* in [Ancona *et al.*, 2018], and as *infidelity* in [Yeh *et al.*, 2019].

Low complexity is loosely related to *information gain* from [Bylinskii *et al.*, 2018] and to *descriptive sparsity* from [Warnecke *et al.*, 2019].

Moreover, very recent literature has also tried to develop various other explanation function criteria: parameter randomization [Adebayo *et al.*, 2018], clustering-based interpretations [Carter *et al.*, 2019], existence of “unexplainable components” [Zhang *et al.*, 2019], variants of perturbation techniques [Grabska-Barwińska, 2020], variants of mutual information measures [Davis *et al.*, 2020], impact of iterative feature removal [Rieger and Hansen, 2020], and necessity and sufficiency of attributions [Wang *et al.*, 2020].

B Proofs

For thoroughness, we elaborate on the proofs from the main paper here.

B.1 Proof of Proposition 1

Proof. Assuming we fix the predictor \mathbf{f} , let $\mathbf{g}(\mathbf{x}) = \mathbf{g}(\mathbf{f}, \mathbf{x})$ and let \int represent $\int_{\rho(\mathbf{x}, \mathbf{z}) \leq R}$ for the rest of this proof.

$$\begin{aligned} \mu_A(\mathbf{g}_{agg}) &= \int D(\mathbf{g}_{agg}(\mathbf{x}), \mathbf{g}_{agg}(\mathbf{z})) \mathbb{P}_{\mathbf{x}}(\mathbf{z}) d\mathbf{z} \\ &= \int \|\mathbf{g}_{agg}(\mathbf{x}) - \mathbf{g}_{agg}(\mathbf{z})\|_2 d\mathbf{z} \\ &= \int \|w\mathbf{g}_1(\mathbf{x}) + (1-w)\mathbf{g}_2(\mathbf{x}) - w\mathbf{g}_1(\mathbf{z}) - (1-w)\mathbf{g}_2(\mathbf{z})\|_2 d\mathbf{z} \\ &= \int \|w\mathbf{g}_1(\mathbf{x}) - w\mathbf{g}_1(\mathbf{z}) + (1-w)\mathbf{g}_2(\mathbf{x}) - (1-w)\mathbf{g}_2(\mathbf{z})\|_2 d\mathbf{z} \\ &= \int \|w(\mathbf{g}_1(\mathbf{x}) - \mathbf{g}_1(\mathbf{z})) + (1-w)(\mathbf{g}_2(\mathbf{x}) - \mathbf{g}_2(\mathbf{z}))\|_2 d\mathbf{z} \\ &\leq \int \|w(\mathbf{g}_1(\mathbf{x}) - \mathbf{g}_1(\mathbf{z}))\|_2 + \|(1-w)(\mathbf{g}_2(\mathbf{x}) - \mathbf{g}_2(\mathbf{z}))\|_2 d\mathbf{z} \\ &\leq \int w\|\mathbf{g}_1(\mathbf{x}) - \mathbf{g}_1(\mathbf{z})\|_2 + (1-w)\|\mathbf{g}_2(\mathbf{x}) - \mathbf{g}_2(\mathbf{z})\|_2 d\mathbf{z} \\ &\leq \int wD(\mathbf{g}_1(\mathbf{x}), \mathbf{g}_1(\mathbf{z})) + (1-w)D(\mathbf{g}_2(\mathbf{x}), \mathbf{g}_2(\mathbf{z})) d\mathbf{z} \\ &\leq w \int D(\mathbf{g}_1(\mathbf{x}), \mathbf{g}_1(\mathbf{z})) d\mathbf{z} + (1-w) \int D(\mathbf{g}_2(\mathbf{x}), \mathbf{g}_2(\mathbf{z})) d\mathbf{z} \\ &\leq w\mu_A(\mathbf{g}_1) + (1-w)\mu_A(\mathbf{g}_2) \end{aligned}$$

□

B.2 Proof of Proposition 2

Proof. To prove this, we just need to show that the sum of the squared distances is minimized by the mean of a set of explanation vectors:

$$\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \mathbf{g}_i(\mathbf{f}, \mathbf{x})$$

Recall we have a set of candidate explanation functions $\mathcal{G}_m = \{\mathbf{g}_1, \dots, \mathbf{g}_m\}$. Fix a point of interest \mathbf{x} . Since d is the ℓ_2 distance and $p = 2$, we define a loss function as follows:

$$L(\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x})) = \sum_{i=1}^m \|\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x}) - \mathbf{g}_i(\mathbf{f}, \mathbf{x})\|_2^2$$

We then compute the partial derivatives with respect to each feature of our aggregate explanation $\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x})_j$.

$$\frac{\partial L}{\partial \mathbf{g}_{agg}(\mathbf{f}, \mathbf{x})_j} = 2m\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x})_j - 2 \sum_{i=1}^m \mathbf{g}_i(\mathbf{f}, \mathbf{x})_j$$

$$\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x})_j = \frac{\sum_{i=1}^m \mathbf{g}_i(\mathbf{f}, \mathbf{x})_j}{m}$$

$$\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x}) = \begin{bmatrix} \frac{\sum_{i=1}^m \mathbf{g}_i(\mathbf{f}, \mathbf{x})_1}{m} \\ \vdots \\ \frac{\sum_{i=1}^m \mathbf{g}_i(\mathbf{f}, \mathbf{x})_d}{m} \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m \mathbf{g}_i(\mathbf{f}, \mathbf{x})$$

□

B.3 Proof of Proposition 3

Proof. To prove this, we just need to show that the sum of the absolute distances is minimized by the median of a set of explanation vectors:

$$\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x}) = \text{med}\{\mathbf{g}_i(\mathbf{f}, \mathbf{x})\}$$

Recall we have a set of candidate explanation functions $\mathcal{G}_m = \{\mathbf{g}_1, \dots, \mathbf{g}_m\}$. Fix a point of interest \mathbf{x} . Since d is the l_1 distance and $p = 1$, we define a loss function as follows:

$$L(\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x})) = \sum_{i=1}^m |\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x}) - \mathbf{g}_i(\mathbf{f}, \mathbf{x})|$$

Taking the partial derivative of the above with respect to each feature of our aggregate explanation $\mathbf{g}_{agg}(\mathbf{x})_j$ yields.

$$\frac{\partial L}{\partial \mathbf{g}_{agg}(\mathbf{f}, \mathbf{x})_j} = \sum_{i=1}^m \text{sign}(\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x})_j - \mathbf{g}_i(\mathbf{f}, \mathbf{x})_j)$$

Now the above partial derivative only equals zero when the number of positive and negative items are the same. The median is the only value where the number of positive items (those greater than the median) and the number of negative items (those less than the median) are equal. Thus, the median value for each feature j would minimize the sum of absolute deviations loss we crafted above (i.e. $\mathbf{g}_{agg}(\mathbf{f}, \mathbf{x})_j = \text{med}\{\mathbf{g}_1(\mathbf{f}, \mathbf{x})_j, \mathbf{g}_2(\mathbf{f}, \mathbf{x})_j, \dots, \mathbf{g}_m(\mathbf{f}, \mathbf{x})_j\}$). \square

B.4 Alternative Proof of Theorem 5

Proof. We want to show that $\mathbf{g}_{AVA}(\mathbf{f}, \mathbf{x}_{\text{test}}) = \Phi_{\mathbf{x}_{\text{test}}}$ is indeed a vector of Shapley values. Let $\mathbf{g}_{\text{SHAP}}(\mathbf{f}, \mathbf{z}) = \Phi_{\mathbf{z}}$ be the vector of Shapley value contributions for a point $\mathbf{z} \in \mathcal{N}_k$. By [Lundberg and Lee, 2017], we know that $\mathbf{g}_{\text{SHAP}}(\mathbf{f}, \mathbf{z})_i = \phi_i(v_{\mathbf{z}})$ is a unique Shapley value for the characteristic function $v_{\mathbf{z}}$. By linearity of Shapley values [Shapley, 1953], we know that:

$$\phi_i(v_{\mathbf{z}_1} + v_{\mathbf{z}_2}) = \phi_i(v_{\mathbf{z}_1}) + \phi_i(v_{\mathbf{z}_2}) \quad (5)$$

This means that the $\Phi_{\mathbf{z}_1} + \Phi_{\mathbf{z}_2}$ will yield a unique Shapley value contribution vector for the characteristic function $v_{\mathbf{z}_1} + v_{\mathbf{z}_2}$. By linearity (also called additivity), we also know that, for any scalar α :

$$\alpha \phi_i(v_{\mathbf{z}}) = \phi_i(\alpha v_{\mathbf{z}}) \quad (6)$$

This means that the $\alpha \Phi_{\mathbf{z}}$ will yield a unique Shapley value contribution vector for the characteristic function $\alpha v_{\mathbf{z}}$. Now, to show $\Phi_{\mathbf{x}_{\text{test}}}$ is a vector of Shapley values, it suffices to show that any $\phi_i(v_{\text{AVA}}) \in \Phi_{\mathbf{x}_{\text{test}}}$ is a Shapley value. As such, we define v_{AVA} to be the characteristic function of $\mathbf{g}_{\text{AVA}}(\mathbf{f}, \mathbf{x})$, where we find the average weighted importance score of the neighbors of \mathbf{x}_{test} .

$$\begin{aligned} v_{\text{AVA}}(S) &= \sum_{\mathbf{z} \in \mathcal{N}_k(\mathbf{x}_{\text{test}})} \frac{v_{\mathbf{z}}(S)}{\rho(\mathbf{x}_{\text{test}}, \mathbf{z})} \\ &= \sum_{\mathbf{z} \in \mathcal{N}_k(\mathbf{x}_{\text{test}})} \frac{1}{\rho(\mathbf{x}_{\text{test}}, \mathbf{z})} \mathbb{E}_Y \left[-\log \frac{1}{\mathbb{P}_{\mathbf{f}}(Y|z_s)} \Big| \mathbf{z} \right] \end{aligned} \quad (7)$$

By Equations 5, 6, and 7, we can see that $\phi_i(v_{\text{AVA}})$ is a Shapley value.

$$\begin{aligned} \mathbf{g}_{\text{AVA}}(\mathbf{f}, \mathbf{x}_{\text{test}})_i &= \phi_i(v_{\text{AVA}}) \\ &= \sum_{\mathbf{z} \in \mathcal{N}_k(\mathbf{x}_{\text{test}})} \frac{\mathbf{g}_{\text{SHAP}}(\mathbf{f}, \mathbf{z})_i}{\rho(\mathbf{x}_{\text{test}}, \mathbf{z})} \\ &= \sum_{\mathbf{z} \in \mathcal{N}_k(\mathbf{x}_{\text{test}})} \frac{\phi_i(v_{\mathbf{z}})}{\rho(\mathbf{x}_{\text{test}}, \mathbf{z})} \end{aligned} \quad (8)$$

\square

C Details on Lowering Complexity

Given a fixed input \mathbf{x} and an explanation function \mathbf{g}_i , the complexity can be rewritten as:

$$\mu_C(\mathbf{f}, \mathbf{g}_i; \mathbf{x}) = - \sum_{k=1}^d \frac{|\mathbf{g}_i(\mathbf{f}, \mathbf{x})_k|}{\sum_{j \in [d]} |\mathbf{g}_i(\mathbf{f}, \mathbf{x})_j|} \ln \left(\frac{|\mathbf{g}_i(\mathbf{f}, \mathbf{x})_k|}{\sum_{j \in [d]} |\mathbf{g}_i(\mathbf{f}, \mathbf{x})_j|} \right)$$

This will help us determine how a small perturbation of the k th component of $\mathbf{g}_i(\mathbf{f}, \mathbf{x})$ will affect the complexity of \mathbf{g}_i , which, in turn, will help find a lower complexity explanation. Note $\mathbf{g}_i(\mathbf{f}, \mathbf{x})_k$ is the k th component of $\mathbf{g}_i(\mathbf{f}, \mathbf{x})$. The partial derivative of $\mu_C(\mathbf{f}, \mathbf{g}_i; \mathbf{x})$ with respect to the k th component of $\mathbf{g}_i(\mathbf{f}, \mathbf{x})$ is:

$$\begin{aligned} \frac{\partial \mu_C(\mathbf{f}, \mathbf{g}_i; \mathbf{x})}{\partial \mathbf{g}_i(\mathbf{f}, \mathbf{x})_k} &= -(1 + \ln(a)) \frac{\sum_{l=1}^d |\mathbf{g}_i(\mathbf{f}, \mathbf{x})_l|}{\left(\sum_{j \in [d]} |\mathbf{g}_i(\mathbf{f}, \mathbf{x})_j| \right)^2} \\ &\quad + \sum_{\substack{l=1 \\ l \neq k}}^d (1 + \ln(b)) \frac{|\mathbf{g}_i(\mathbf{f}, \mathbf{x})_l|}{\left(\sum_{j \in [d]} |\mathbf{g}_i(\mathbf{f}, \mathbf{x})_j| \right)^2} \end{aligned}$$

where $a = \frac{|\mathbf{g}_i(\mathbf{f}, \mathbf{x})_k|}{\sum_{j \in [d]} |\mathbf{g}_i(\mathbf{f}, \mathbf{x})_j|}$ and $b = \frac{|\mathbf{g}_i(\mathbf{f}, \mathbf{x})_l|}{\sum_{j \in [d]} |\mathbf{g}_i(\mathbf{f}, \mathbf{x})_j|}$.

We now provide an additional discussion and comparison of the two algorithms for lowering complexity.

We presented two algorithms for finding a \mathbf{g}_{agg} with lower complexity: a gradient descent approach (Algorithm 1) and a region shrinking approach (Algorithm 2). Algorithm 1 relies on a greedy choice of selecting one of the j directions to move in. This algorithm works best for regions that are smooth and with decreasing complexity around \mathbf{g}_i and \mathbf{g}_{avg} . Since Algorithm 1 does not backtrack and moves component-wise, it can avoid areas of higher complexity, but can take a sub-optimal step. For example, consider when $d = 2$. During a walk, Algorithm 1 may start at \mathbf{g}_i , move in the x direction, but then get stuck as complexity in y direction increases. However, had we moved in the y direction first and then in the x direction, then we may have found a minimum. The choice of component plagues this approach. On the other hand, Algorithm 2 solves the issue of getting stuck because of regions of high complexity present in Algorithm 1. Since Algorithm 2 shrinks the region by choosing points in the convex combination, it can avoid the areas of high complexity. Since Algorithm 2 uses the line segments between the points chosen, it may be difficult to obtain the global minima, which

Algorithm 1 Gradient-Descent Style Approach to finding $g_{\text{agg}}(\mathbf{f}, \mathbf{x})$ with lower complexity

Require: $\alpha, \mathbf{g}_i(\mathbf{f}, \mathbf{x}), i = 1, \dots, m$, fixed \mathbf{x}

▷ Calculate the complexity of each $\mathbf{g}_i(\mathbf{f}, \mathbf{x})$

for $i = 1, \dots, m$ **do**

$$E_{\mathbf{g}_i(\mathbf{x})} \leftarrow \mu_{\mathbf{C}}(\mathbf{f}, \mathbf{g}_i; \mathbf{x}) \leftarrow \sum_{k=1}^d \frac{|\mathbf{g}_i(\mathbf{f}, \mathbf{x})_k|}{\sum_{j \in [d]} |\mathbf{g}_i(\mathbf{f}, \mathbf{x})_j|} \ln \left(\frac{|\mathbf{g}_i(\mathbf{f}, \mathbf{x})_k|}{\sum_{j \in [d]} |\mathbf{g}_i(\mathbf{f}, \mathbf{x})_j|} \right)$$

end for

$$\mathbf{g}_{\text{avg}}(\mathbf{f}, \mathbf{x}) \leftarrow \frac{1}{m} \sum_{i=1}^m \mathbf{g}_i(\mathbf{f}, \mathbf{x})$$

for $i = 1, \dots, m$ **do**

▷ Move in the direction of $\mathbf{g}_{\text{avg}}(\mathbf{f}, \mathbf{x})$ from $\mathbf{g}_i(\mathbf{f}, \mathbf{x})$ as long as the complexity decreases

$$\mathbf{t}_i \leftarrow \mathbf{g}_i(\mathbf{f}, \mathbf{x})$$

while Complexity of \mathbf{t}_i is decreasing and $\mathbf{t}_i \neq \mathbf{g}_{\text{avg}}(\mathbf{f}, \mathbf{x})$ **do**

for $j = 1, \dots, d$ **do**

$$\text{Calculate } \frac{\partial E_{\mathbf{t}_i}}{\partial \mathbf{t}_{ij}}$$

if Complexity decreases by moving in the j direction towards $\mathbf{g}_{\text{avg}}(\mathbf{f}, \mathbf{x})$ **then**

▷ Update \mathbf{t}_{ij}

$$\mathbf{t}_{ij} \leftarrow \mathbf{t}_{ij} + \alpha \frac{\partial E_{\mathbf{t}_i}}{\partial \mathbf{t}_{ij}}$$

end if

end for

end while

▷ Move in the direction of $\mathbf{g}_i(\mathbf{f}, \mathbf{x})$ from $\mathbf{g}_{\text{avg}}(\mathbf{f}, \mathbf{x})$ as long as the complexity decreases

$$\mathbf{q}_i \leftarrow \mathbf{g}_{\text{avg}}(\mathbf{f}, \mathbf{x})$$

while Complexity of \mathbf{q}_i is decreasing and $\mathbf{q}_i \neq \mathbf{g}_i(\mathbf{x})$ **do**

for $j = 1, \dots, d$ **do**

$$\text{Calculate } \frac{\partial E_{\mathbf{q}_i}}{\partial \mathbf{q}_{ij}}$$

if Complexity decreases by moving in the j direction towards $\mathbf{g}_i(\mathbf{x})$ **then**

▷ Update \mathbf{q}_{ij}

$$\mathbf{q}_{ij} \leftarrow \mathbf{q}_{ij} + \alpha \frac{\partial E_{\mathbf{q}_i}}{\partial \mathbf{q}_{ij}}$$

end if

end for

end while

▷ Take the $\mathbf{t}_i, \mathbf{q}_i$ that minimizes the complexity

$$\mathbf{b}_i = \min_{\mathbf{x}=\{\mathbf{q}_i, \mathbf{t}_i\}} E_{\mathbf{x}}$$

end for

▷ Take the \mathbf{b}_i that minimizes the complexity

$$\mathbf{g}_{\text{agg}}(\mathbf{f}, \mathbf{x}) = \min_{\mathbf{b}_i} E_{\mathbf{b}_i}$$

Algorithm 2 Region Shrinking Approach to finding $g_{\text{agg}}(\mathbf{f}, \mathbf{x})$ with lower complexity

Require: $\mathbf{g}_i(\mathbf{f}, \mathbf{x}), i = 1, \dots, m$, fixed \mathbf{x}

$$t \leftarrow 0$$

▷ Add all the \mathbf{g}_i into set S

$$S \leftarrow \{\mathbf{g}_i(\mathbf{f}, \mathbf{x}), i = 1, \dots, m\}$$

repeat

▷ Repeat K times

▷ Initialize S'

$$S' \leftarrow \emptyset$$

for every 2 points in $S: P_1, P_2$ **do**

Find point P with the minimum entropy in the convex combination of P_1, P_2

Add point P to S'

end for

▷ Update values

Choose the N minimum entropy points in S' to form S

$$t \leftarrow t + 1$$

until $t = K$

▷ Take the element in set S that minimizes the entropy

$$\mathbf{g}_{\text{agg}}(\mathbf{f}, \mathbf{x}) = \min_{k \in S} E_k$$

may not occur on the line segment. A combination of the two approaches can be used. First, Algorithm 2 can be used to shrink the region being considered into a set, S , of points with low complexity. This can avoid getting stuck in areas of high complexity, like in Algorithm 1. Then, Algorithm 1 can be used to move around the points in set S in order to find the global minima that may not occur on the line segments considered in Algorithm 2. It can refine the points in set S to obtain a lower complexity. In sum, we can shrink the region considered into several candidate sets and then refine the points in each set by perturbing and performing greedy walks around them to find g_{agg} with a low complexity.

D Experimental Setup

We provide additional details on the datasets used and their respective models from our experiments.

- Iris [Dua and Graff, 2017]: The iris dataset consists of 150 datapoints: 50 per class and 4 features per datapoint. We use a one layer multilayer perceptron trained to 96% accuracy as our \mathbf{f} .
- Adult [Dua and Graff, 2017]: Each of the 48842 datapoints has 38 features and falls in one of two classes. Note we label encode categorical attributes. We use a one layer MLP (40 hidden nodes with leaky-relu activation) trained to an accuracy of 82%.
- Mimic-III [Johnson *et al.*, 2016]: The MIMIC-III (Medical Information Mart for Intensive Care III) is a large electronic health record dataset comprised of health related data of over 40,000 patients who were admitted to the the critical care units of Beth Israel Deaconess Medical Center between the years 2001 and 2012. MIMIC-III consists of demographics, vital sign measurements, lab test results, medications, procedures,

caregiver notes, imaging reports, and mortality of the ICU patients. Using MIMIC-III dataset, we extracted seventeen real-valued features deemed critical in the sepsis diagnosis task as per [Purushotham *et al.*, 2018]. These are the processed features we extracted for every sepsis diagnosis (a binary variable indicating the presence of sepsis): Glasgow Coma Scale, Systolic Blood Pressure, Heart Rate, Body Temperature, Pao2 / Fio2 ratio, Urine Output, Serum Urea Nitrogen Level, White Blood Cells Count, Serum Bicarbonate Level, Sodium Level, Potassium Level, Bilirubin Level, Age, Acquired Immunodeficiency Syndrome, Hematologic Malignancy, Metastatic Cancer, Admission Type. We used two layers of 16 hidden nodes each and leaky-reLU activation to get an accuracy of 91% on the sepsis prediction task.

- MNIST with CNN [LeCun *et al.*, 1998]: We use a CNN trained to 90% accuracy with the following architecture: one layer with $32 \ 5 \times 5$ filters and ReLU activation; max pooling layer with a 2×2 filter and stride of 2; convolutional layer with $64 \ 5 \times 5$ filters and ReLU activation; max pooling layer with a 2×2 filter and stride of 2; final dense layer with 10 output neurons. We used the MNIST dataset with 60,000 28x28 grayscale images of the 10 digits, along with a test set of 10,000 images.

Note that we fix a dataset-model pairing for all experiments. In practice, when calculating *average sensitivity*, we use the following formulation:

$$\mu_A(\mathbf{f}, \mathbf{g}, \mathbf{x}) = \frac{1}{|\mathcal{N}_r|} \sum_{\mathbf{z} \in \mathcal{N}_r} \frac{D(\mathbf{g}(\mathbf{f}, \mathbf{x}), \mathbf{g}(\mathbf{f}, \mathbf{z}))}{\rho(\mathbf{x}, \mathbf{z})}$$

Effectively, we want to ensure that the distance between an explanation of \mathbf{x} and an explanation of \mathbf{z} , a point in the neighborhood of \mathbf{x} , is proportional to the distance between \mathbf{x} and \mathbf{z} . Some recent work has shown that *average sensitivity* can be lowered with simple smoothing tricks to explanation functions or with adversarial training of the predictor itself.