

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.

Evaluation Criteria

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. 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 f, g , and \mathbf{x} , and outputs a scalar: $\mu(f, g; \mathbf{x})$.

Desideratum: Low Sensitivity. If inputs are near each other and their model outputs are similar, then their explanations should be close to each other. Let $\mathcal{N}_r = \{z \in \mathcal{D}_{\mathbf{x}} \mid \rho(\mathbf{x}, z) \leq r, f(\mathbf{x}) = f(z)\}$ be a neighborhood of datapoints within a radius r of \mathbf{x} .

$$\begin{aligned} \text{Max Sensitivity: } \mu_M(f, g, r; \mathbf{x}) &= \max_{z \in \mathcal{N}_r} D(g(f, \mathbf{x}), g(f, z)) \\ \text{Average Sensitivity: } \mu_A(f, g, r; \mathbf{x}) &= \int_{z \in \mathcal{N}_r} D(g(f, \mathbf{x}), g(f, z)) \mathbb{P}_{\mathbf{x}}(z) dz \end{aligned}$$

Desideratum: High Faithfulness. 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 $|S| = d$, $\mathbf{x}_{[\mathbf{x}_s = \bar{\mathbf{x}}_s]} = \bar{\mathbf{x}}$.

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

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:

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

Let $\mathbb{P}_g(i)$ denote the fractional contribution of feature \mathbf{x}_i to the total magnitude of the attribution. If every feature had equal attribution, the explanation would be complex even if faithful. The simplest explanation would be concentrated on one feature. We define complexity as the entropy of \mathbb{P}_g .

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

Other Aggregation Methods

Given $f, \mathcal{G}_m = \{g_1, \dots, g_m\}$, μ , and a set of inputs $\mathcal{D}_{\mathbf{x}}$, we want to find an aggregate explanation function g_{agg} that satisfies μ at least as well as any $g_i \in \mathcal{G}_m$. Let $h(\cdot)$ represent some function that combines m explanations into a consensus $g_{\text{agg}} = h(\mathcal{G}_m)$.

Convex Combination. Suppose we have two different explanation functions g_1 and g_2 and have chosen a criterion μ to evaluate a g . Consider an aggregate explanation, $g_{\text{agg}} = h(g_1, g_2)$. A potential $h(\cdot)$ is a convex combination where $g_{\text{agg}} = h(g_1, g_2) = w g_1 + (1-w) g_2 = w^T \mathcal{G}_m$.

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:

$$g_{\text{agg}} \in \arg \min_{g \in \mathcal{G}} \mathbb{E}_{g_i \in \mathcal{G}_m} [D(g, g_i)^p] = \arg \min_{g \in \mathcal{G}} \sum_{i=1}^m D(g, g_i)^p$$

where p is a positive constant. Assuming real-valued attributions where $\mathcal{G} \subseteq \mathbb{R}^d$, when D is ℓ_2 and $p = 2$, the aggregate explanation is the feature-wise sample mean; when D is ℓ_1 and $p = 1$, the aggregate is the feature-wise sample median [1]. 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 is the Kemeny-Young rule.

Future Work. We could leverage multi-objective optimization to find g_{agg} . For example, we could learn a less sensitive and less complex explanation function by optimizing:

$$g_{\text{agg}} \in \arg \min_{g \in \mathcal{G}} \mathbb{E}_{\mathbf{x} \in \mathcal{D}} [\mu_A(f, g, r; \mathbf{x}) + \lambda \mu_C(f, g; \mathbf{x})]$$

Algorithms for Aggregation

We can aggregate explanations from various g_i to optimize a particular evaluation criterion.

Lowering Sensitivity. Termed Aggregate Valuation of Antecedents (AVA), we derive an explanation function that explains a data point in terms of the explanations of its neighbors. To obtain an explanation $g_{\text{AVA}}(f, \mathbf{x}_{\text{test}})$ for a point of interest \mathbf{x}_{test} , we first find the k nearest neighbors of \mathbf{x}_{test} under ρ denoted by $\mathcal{N}_k(\mathbf{x}_{\text{test}}, \mathcal{D})$.

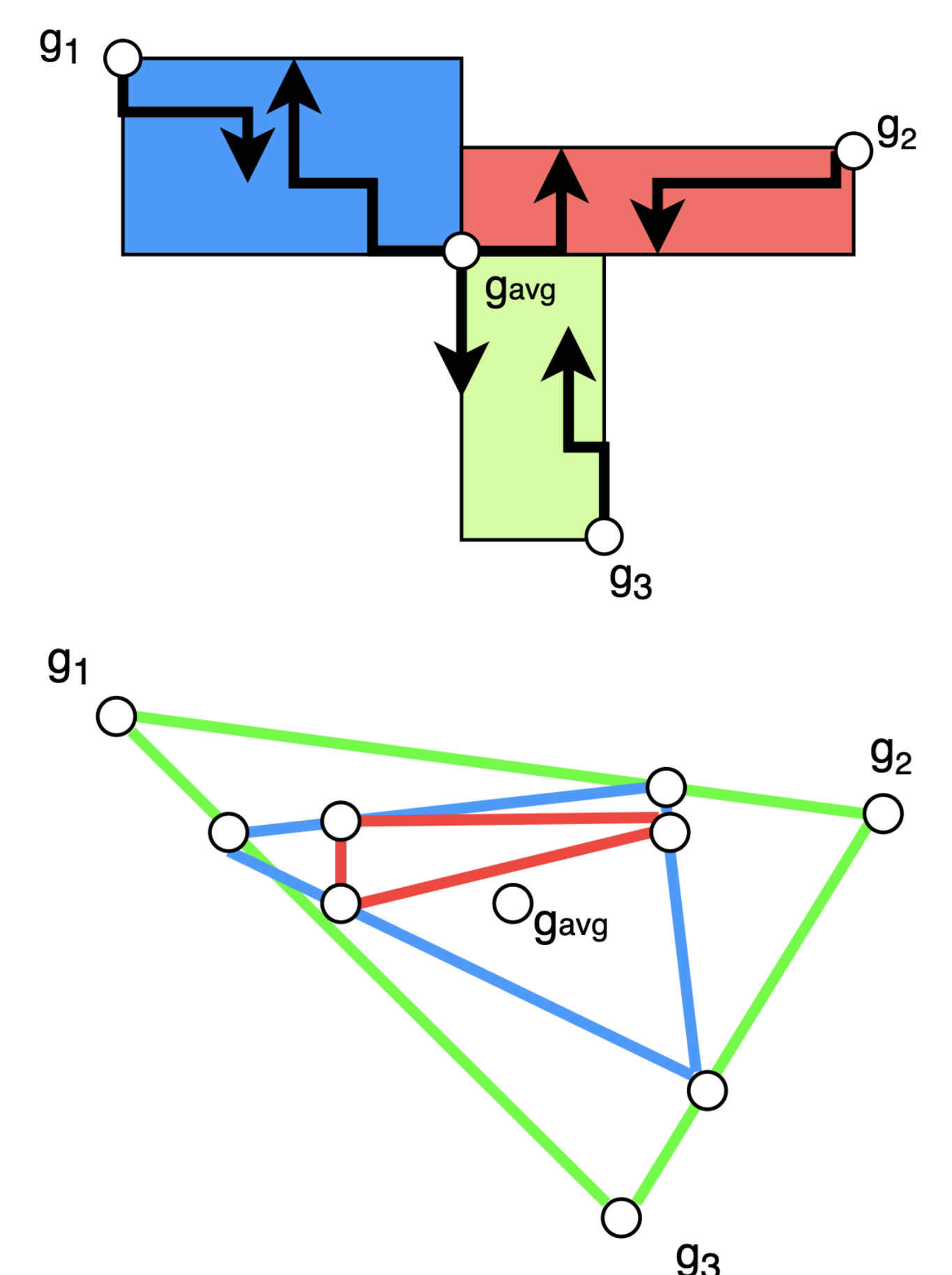
$$\mathcal{N}_k(\mathbf{x}_{\text{test}}, \mathcal{D}) = \arg \min_{\mathcal{N} \subset \mathcal{D}, |\mathcal{N}|=k} \sum_{z \in \mathcal{N}} \rho(\mathbf{x}_{\text{test}}, z)$$

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

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

Theorem 1. $g_{\text{AVA}}(f, \mathbf{x}_{\text{test}})$ is a Shapley value explanation.

Lowering Complexity. We devise iterative algorithms for aggregating explanation functions to obtain $g_{\text{agg}}(f, \mathbf{x})$ with lower complexity whilst combining m candidate explanation functions $\mathcal{G}_m = \{g_1, \dots, g_m\}$. We desire a $g_{\text{agg}}(f, \mathbf{x})$ that contains information from all candidate explanations $g_i(f, \mathbf{x})$ yet has entropy less than or equal to that of each explanation $g_i(f, \mathbf{x})$.



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

References

- [1] James O Berger. *Statistical decision theory and Bayesian analysis*. Springer Science & Business Media, 2013.