Explaining Tree Models

Bernhard Jaeger

University of Tübingen

Tübingen, Germany
bernhard.jaeger@student.uni-tuebingen.de

Abstract—Understanding the decisions made by machine learning models can be challenging due to their black box nature. We investigate the additive feature attribution method SHAP (SHapley Additive exPlanations) that addresses this problem [2]. While SHAP can be used to explain the decisions of any model, the generic KernelSHAP algorithm for it can be intractable for larger datasets. For tree models there exists a specialized algorithm (Tree SHAP) that can compute exact SHAP values in polynomial time [1]. We use the SHAP framework to investigate the relationship between Decision Trees and Hoeffding Trees [3]. Theory tells us that a Hoeffding Tree should converge towards a Decision Tree when trained. In our experiments however we could not observe convergence. The reason for this remains unclear.

Index Terms—Tree SHAP, SHAP, Decision Tree, Hoeffding Tree

I. Introduction

This . This and [3] document is a model and instructions for LATEX. Please observe the conference page limits.

II. ADDITIVE FEATURE ATTRIBUTION METHODS

Modern models of Machine Learning have become increasingly complex making it hard to explain their decisions. To address this challenge we can go one layer of abstraction further and make a model of the model which we want to be simple enough to be interpretable again. In [2] they coined the term explanation model for this approach of using an interpretable approximation of the original model. When viewed trough this lens multiple of the current explanation methods can be categorized in the class of **Additive Feature Attribution methods**. In [2] they defined this as methods having an explanation model that is a linear function of binary variables:

$$g(z') = \phi_0 + \sum_{i=1}^{M} \phi_i z_i'$$
 (1)

Where g is the explanation model, $z' \in \{0,1\}^M$, $\phi \in \mathbb{R}$. Local methods like this first approximate their model f(x) using simplified inputs x' where $x = h_x(x')$.

The objective is then to ensure $g(z') \approx f(h_x(z'))$ whenever $z' \approx x'$.

 ϕ_i is then called the **effect** of feature i and summing up all effects approximates the output of the model f(x) we wish to explain. Note that these *explanation models* are now easily interpretable. We can for a prediction f(x) given input x say what influence or effect each component x_i had on the final

decision.

Methods falling in the category of Additive Feature Attribution methods include LIME [4], DeepLIFT [5], layerwise relevance propagation [6], Quantitive Input Influence [7], Shapley sampling values [8] and SHAP [2].

Additive Feature Attribution methods can have 3 desirable theoretical properties [2]:

Local Accuracy:

$$f(x) = g(x')whenx = h_x(x')and\phi_0 = f(h_x(0))$$
 (2)

The intuition being that the explanation model should match the original model for the simplified input (notice x'=z' here).

Missingness:

$$x_i' = 0 \to \phi_i = 0 \tag{3}$$

The intuition is that feature missing in the original input may have no impact even if they are present in the simplified version. (Notice that even if x' = 0, z' does not necessarily have to be 0).

Consistency:

 $f_x(z') := f(h_x(z'))$, and let $z' \setminus i$ be $z'_i = 0$ For any two models f and f', if:

$$f_x'(z') - f_x'(z' \setminus i) \ge f_x(z') - f_x(z' \setminus i) \tag{4}$$

for all inputs $z' \in 0, 1^M$.

Meaning that the model changes in a way that some simplified input's contribution becomes larger.

Then: $\phi_i(f', x) \ge \phi_i(f, x)$

Meaning that it's inputs attribution should not become smaller.

Now it turns out that the only possible features for an Additive Feature attribution method that satisfies all 3 properties are the Shapley values defined as:

$$\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$$
(5)

 $z' \subseteq x'$ are all z' vectors where the non-zero entries are a subset of the non-zero entries in x'.

|z'| represents the number of non-zero entries in z'.

This means that at least in theory Additive Feature Attribution methods using Shapley Values as effects are to be preferred over other methods. In practice methods using shapely values can be intractable for larger sets.

In the next section we will introduce the SHAP framework which tries to unify these measures of feature importance.

III. SHAP VALUES

SHAP (SHapley Additive exPlanations) values [2] are obtained by solving equation 5 using a conditional expectation function of the original model:

$$f(h_x(z')) = E[f(z)|z_S]$$
(6)

S is the set of non-zero indexes in z'. The mapping function is therefore implicitly defined as:

$$h_x(z') = z_S \tag{7}$$

 z_S has missing values for features not in the set S. One problem with this approach is that in general models don't can't necessarily handle features that are just missing.

To address this issue $f(z_S)$ can be approximated with $E[f(z)|z_S].$

The intuition being that SHAP values calculate the attribution (or effect) of each feature using the change in expectation of the model prediction when you condition on that feature. The feature $\phi_0 = E[f(z)]$ is the base value or the prediction the model would make if it had no observation. The following ϕ_i are then calculated by conditioning on some input features and averaging over all possible orderings of them. As hinted already earlier the exact computation of SHAP values can be challenging as it requires to calculate combinatory many expectations. Assuming feature independence (equation 10) and model linearity (equation 11) can make the computation much easier.

$$f(h_x(z')) = E[f(z)|z_S]$$
(8)

$$=E_{z_{\overline{S}}|z_S}[f(z)] \tag{9}$$

$$\approx E_{z_{\overline{S}}}[f(z)]$$
 (10)

$$\approx f([z_S, E[z_{\overline{S}}]])$$
 (11)

It allows you to These are however strong assumptions as most modern models are non-linear and features are in general not independent of each other (e.g. a pixel value only has a meaningful interpretation when put into the context of the surrounding pixels).

In the following subsections we will discuss algorithms for computing SHAP values.

A. Kernel SHAP (Linear Lime + Shapely Value)

Kernel SHAP [2] is the only method presented here that can be applied to any model. Linear Lime is a method that minimizes the following objective function:

$$\xi = \operatorname*{argmin}_{g \in G} L(f, g, \pi_{x'}) + \Omega(g)$$
 (12)

Lime is a local linear explanation model that locally approximates the model around a given prediction. Kernel SHAP now chooses Ω , π and L in a way that recovers the Shapely Values (and therefore guarantees the properties missingness, local accuracy and consistency). This is possible because LIME is a additive feature attribution method. This lead to the following equations:

$$\Omega(g) = 0 \tag{13}$$

$$\pi_{x'}(z') = \frac{(M-1)}{(Mchoose|z'|)|z'|(M-|z'|)1}$$
(14)

$$\pi_{x'}(z') = \frac{(M-1)}{(Mchoose|z'|)|z'|(M-|z'|)1}$$

$$L(f, g, \pi_{x'}) = \sum_{z' \in Z} [f(h_x(z')) - g(z')]^2 \pi_{x'}(z')$$
(15)

Because the loss is quadratic and g(z') is assumed to be linear we can solve equation 12 using linear regression (and calculate the shapely values using weighted linear regression). The input mapping used by LIME is equivalent to equation 11 and we can therefore use Kernel SHAP on every model. The simplified input mapping $h_x(z')$ that LIME uses is equivalent to equation 11 which leads to estimation of the SHAP values.

To summarize Kernel SHAP is a model agnostic algorithm that under the (strong) assumptions model linearity and feature independence approximates the SHAP values. The time complexity of the algorithm is exponential in the number of input features M. In the original paper the time complexity was stated as $O(2^M + M^3)$ [2]. The ALIBI implementation of the algorithm claims $O(M2^M)$ as time complexity. [9] An additional problem that arises in practice is that most models can not handle arbitrary patterns of missing inputs at inference time (which is a requirement to compute the Shapley values). Implementation of the algorithm therefore make further approximations by replacing missing features with random values from a background dataset provided by the user (e.g. the training set). Given these severe flaws one can question whether the explanations provided by Kernel

Better algorithms for calculating SHAP values for specific model types exist and will be covered in the following subsections:

SHAP can be trusted in a realistic setting.

B. Linear SHAP

IV. DECISION AND HOEFFDING TREES

A. Abbreviations and Acronyms

Define abbreviations and acronyms the first time they are used in the text, even after they have been defined in the abstract. Abbreviations such as IEEE, SI, MKS, CGS, ac, dc, and rms do not have to be defined. Do not use abbreviations in the title or heads unless they are unavoidable.

B. Equations

Number equations consecutively. To make your equations more compact, you may use the solidus (/), the exp function, or appropriate exponents. Italicize Roman symbols for quantities and variables, but not Greek symbols. Use a long dash rather than a hyphen for a minus sign. Punctuate equations with commas or periods when they are part of a sentence, as in:

$$a + b = \gamma \tag{16}$$

C. LATEX-Specific Advice

Please use "soft" (e.g., \eqref{Eq}) cross references instead of "hard" references (e.g., (1)). That will make it possible to combine sections, add equations, or change the order of figures or citations without having to go through the file line by line.

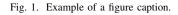
D. Figures and Tables

a) Positioning Figures and Tables: Place figures and tables at the top and bottom of columns. Avoid placing them in the middle of columns. Large figures and tables may span across both columns. Figure captions should be below the figures; table heads should appear above the tables. Insert figures and tables after they are cited in the text. Use the abbreviation "Fig. 1", even at the beginning of a sentence.

TABLE I TABLE TYPE STYLES

Table	Table Column Head		
Head	Table column subhead	Subhead	Subhead
copy	More table copy ^a		
^a Sample of a Table footnote.			

bumple of a facto feetilete



V. METHODOLOGY VI. EVALUATION VII. CONCLUSION

VIII. REPRODUCIBILITY CONSIDERATIONS

A. Software Libraries:

The software library versions we used (read out using pip –freeze) are documented in table II.

All experiments were conducted on a single personal computer. The hardware specs can be found in table III.

TABLE II SOFTWARE LIBRARY VERSIONS

Software library	Version
Python	3.7.4
anaconda-client	1.7.2
anaconda-navigator	1.9.7
anaconda-project	0.8.3
conda	4.7.12
conda-build	3.18.9
conda-package-handling	1.6.
0 conda-verify	3.4.2
ipykernel	5.1.2
ipython	7.8.0
ipython-genutils	0.2.0
jupyter	1.0.0
jupyter-client	5.3.3
jupyter-console	6.0.0
jupyter-core	4.5.0
jupyterlab	1.1.4
jupyterlab-server	1.0.6
matplotlib	3.1.1
notebook	6.0.1
numpy	1.16.5
pandas	0.25.1
scikit-image	0.15.0
scikit-learn	0.21.3
scikit-multiflow	0.5.0
scipy	1.3.1
shap	0.35.0
sklearn	0.0
torch	1.5.0+cu101
torchvision	0.6.0+cu101

TABLE III
HARDWARE SPECIFICATIONS:

Component	Model	
Processor	Intel Core i7-6700 CPU	
RAM	16 GB Corsair Vengeance LPX DDR4-2400	
GPU	NVIDIA GeForce GTX 980	
Motherboard	ASUS Z170 Pro Gaming	
Operating system	Microsoft Windows 10 Education N	
Hard Drive	SSD Crucial MX200	

ACKNOWLEDGMENT

The preferred spelling of the word "acknowledgment" in America is without an "e" after the "g". Avoid the stilted expression "one of us (R. B. G.) thanks ...". Instead, try "R. B. G. thanks...". Put sponsor acknowledgments in the unnumbered footnote on the first page.

REFERENCES

- Scott M. Lundberg, Gabriel G. Erion, Su-In Lee: Consistent Individualized Feature Attribution for Tree Ensembles. CoRR abs/1802.03888 (2018).
- [2] Scott M. Lundberg, Su-In Lee: A Unified Approach to Interpreting Model Predictions. NIPS 2017: 4765-4774.
- [3] Pedro M. Domingos, Geoff Hulten: Mining high-speed data streams. KDD 2000: 71-80
- [4] Marco Túlio Ribeiro, Sameer Singh, Carlos Guestrin: "Why Should I Trust You?": Explaining the Predictions of Any Classifier. KDD 2016: 1135-1144.
- [5] Avanti Shrikumar, Peyton Greenside, Anshul Kundaje: Learning Important Features Through Propagating Activation Differences. ICML 2017: 3145-3153.

- [6] Alexander Binder, Grégoire Montavon, Sebastian Bach, Klaus-Robert Müller, Wojciech Samek: Layer-wise Relevance Propagation for Neural Networks with Local Renormalization Layers. CoRR abs/1604.00825 (2016).
- [7] Anupam Datta, Shayak Sen, Yair Zick: Algorithmic Transparency via Quantitative Input Influence: Theory and Experiments with Learning Systems. IEEE Symposium on Security and Privacy 2016: 598-617.
- [8] Erik Strumbelj, Igor Kononenko: Explaining prediction models and individual predictions with feature contributions. Knowl. Inf. Syst. 41(3): 647-665 (2014).
- [9] https://docs.seldon.io/projects/alibi/en/stable/methods/KernelSHAP.html