
A Discussion of Model Explanation Tools with Practical Recommendations

Patrick Hall
H2O.ai, Mountain View, CA

PHALL@H2O.AI

Abstract

This paper discusses several explanatory methods that go beyond the error measurements and plots traditionally used to assess machine learning models. The approaches, decision tree surrogate models, individual conditional expectation (ICE) plots, local interpretable model-agnostic explanations (LIME), partial dependence plots, and Shapley explanations, vary in terms of scope, fidelity, and suitable application domain. Along with descriptions of these methods, practical guidance for usage and in-depth software examples are also presented.

1. Introduction

Interpretability of statistical and machine learning models is a multifaceted, complex, and evolving subject. This paper focuses mostly on just one aspect of model interpretability: explaining the mechanisms and predictions of models trained using supervised decision tree ensemble algorithms, like gradient boosting machines (GBMs) and random forests. Others have defined key terms and put forward general motivations for better interpretability of machine learning models (Lipton, 2016), (Doshi-Velez & Kim, 2017), (Gilpin et al., 2018), (Guidotti et al., 2018). Following Doshi-Velez and Kim (2017), this discussion uses “the ability to explain or to present in understandable terms to a human,” as the definition of *interpretable*. “When you can no longer keep asking why,” will serve as the working definition for a *good explanation* of model mechanisms or predictions (Gilpin et al., 2018).

As in Figure 1, the presented explanatory methods help practitioners make random forests, GBMs, and other types of popular supervised machine learning models more interpretable by enabling post-hoc explanations that are suitable for:

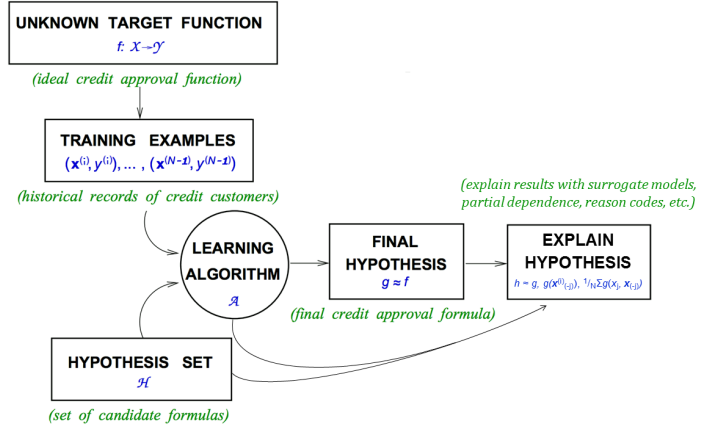


Figure 1. An augmented learning problem diagram in which several techniques create explanations for a credit scoring model. Adapted from **Learning From Data** (Abu-Mostafa et al., 2012).

- Facilitating regulatory compliance.
- Understanding or debugging model mechanisms and predictions.
- Preventing or debugging accidental or intentional discrimination in model predictions.
- Preventing or debugging malicious hacking of models or adversarial attacks on models.

Detailed discussions of the explanatory methods begin below by defining notation. Then sections 3 – 6 discuss explanatory methods and present recommendations for each method. Section 8 presents some general interpretability recommendations for practitioners. Section 9 discusses several additional interpretability subjects that are likely important for practitioners, and finally, section 10 highlights a few software resources that accompany this paper.

2. Notation

To facilitate technical descriptions of explanatory techniques, notation for input and output spaces, datasets, and models is defined.

2.1. Spaces

- Input features come from a set \mathcal{X} contained in a P -dimensional input space, $\mathcal{X} \subset \mathbb{R}^P$.
- Known labels corresponding to instances of \mathcal{X} come from the set \mathcal{Y} and are contained in a C -dimensional label space, $\mathcal{Y} \subset \mathbb{R}^C$.
- Learned output responses come from a set $\hat{\mathcal{Y}}$.

2.2. Datasets

- The input dataset \mathbf{X} is composed of observed instances of the set \mathcal{X} with a corresponding dataset of labels \mathbf{Y} , observed instances of the set \mathcal{Y} .
- Each i -th observation of \mathbf{X} is denoted as $\mathbf{x}^{(i)} = [x_0^{(i)}, x_1^{(i)}, \dots, x_{P-1}^{(i)}]$, with corresponding i -th labels in \mathbf{Y} , $\mathbf{y}^{(i)} = [y_0^{(i)}, y_1^{(i)}, \dots, y_{C-1}^{(i)}]$.
- \mathbf{X} and \mathbf{Y} consists of N tuples of observations: $[(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}), (\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(N-1)}, \mathbf{y}^{(N-1)})]$.
- Each j -th input column vector of \mathbf{X} is denoted as $\mathbf{X}_j = [x_j^{(0)}, x_j^{(1)}, \dots, x_j^{(N-1)}]^T$.

2.3. Models

- A type of machine learning model g , selected from a hypothesis set \mathcal{H} , is trained to represent an unknown target function f observed as \mathbf{X} with labels \mathbf{Y} using a training algorithm \mathcal{A} : $\mathbf{X}, \mathbf{Y} \xrightarrow{\mathcal{A}} g$.
- g generates learned output responses on the input dataset $g(\mathbf{X}) = \hat{\mathbf{Y}}$, and on the general input space $g(\mathcal{X}) = \hat{\mathcal{Y}}$.
- The model to be explained is denoted as g .

3. Surrogate Decision Trees

The phrase *surrogate model* is used here to refer to a simple model, h , of a complex model, g . This type of model is referred to by various other names, such as *proxy* or *shadow* models and the process of training surrogate models is often referred to as *model extraction* (Craven & Shavlik, 1996), (Williams et al., 2017), (Bastani et al., 2017).

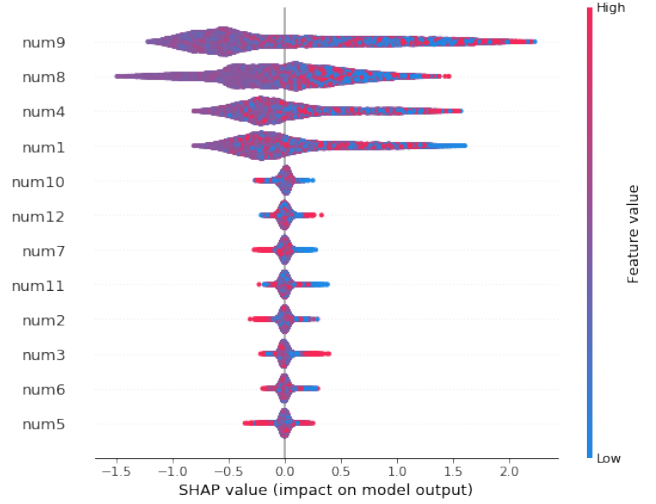


Figure 2. Shapley summary plot for known target function $f = \text{num}_1 * \text{num}_4 + |\text{num}_8| * \text{num}_9^2 + e$, and for a machine-learned GBM response function g_{GBM} .

3.1. Description

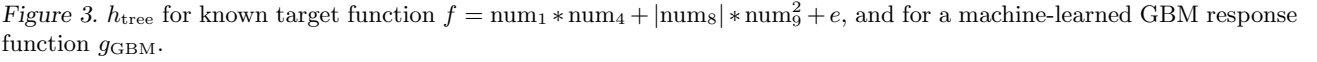
Given a learned function g and set of learned output responses $g(\mathbf{X}) = \hat{\mathbf{Y}}$, and a tree splitting and pruning approach \mathcal{A} , a surrogate decision tree, h_{tree} , can be extracted:

$$\mathbf{X}, g(\mathbf{X}) \xrightarrow{\mathcal{A}} h_{\text{tree}} \quad (1)$$

such that $h_{\text{tree}}(\mathbf{X}) \approx g(\mathbf{X})$.

Decision trees can be represented as directed graphs where the relative positions of input variables can provide insight into their importance and interactions (Breiman et al., 1984). This makes decision trees useful surrogate models. Input variables that appear high and often in the directed graph representation of h_{tree} are assumed to have high importance in g . Input variables directly above or below one-another in h_{tree} are assumed to have strong interactions in g . These relative relationships between input variables in h_{tree} can be used to verify and debug the variable importance, interactions, and predictions of g .

Figures 2 and 3 use simulated data to empirically demonstrate the desired relationships between input variable importance and interactions in the input space \mathbf{X} , the label space $f(\mathbf{X}) = \mathbf{Y}$, a GBM model to be explained g_{GBM} , and a decision tree surrogate h_{tree} . Data with a known target function depending on four input variables with interactions:



and with eight noise variables is simulated. g_{GBM} is trained: $\mathbf{X}, \mathbf{f}(\mathbf{X}) \xrightarrow{A} g_{\text{GBM}}$ such that $g_{\text{GBM}} \approx f$. Then h_{tree} is extracted by $\mathbf{X}, \mathbf{g}_{\text{GBM}}(\mathbf{X}) \xrightarrow{A} h_{\text{tree}}$, such that $h_{\text{tree}}(\mathbf{X}) \approx g_{\text{GBM}}(\mathbf{X}) \approx f(\mathbf{X})$.

Figure 3 is a directed graph representation of h_{tree} that prominently displays the importance of input variables `num9` and `num8` and also visually highlights the interactions between the two inputs. URLs to the data and software used to generate Figures 2 and 3 are available in section 10.

- A shallow-depth h_{tree} displays a global, low-fidelity, high-interpretability flow chart of important features and interactions in g . Because there are few theoretical guarantees that h_{tree} truly represents g , always use error measures to assess the trustworthiness of h_{tree} .
- Prescribed methods (Craven & Shavlik, 1996); (Bas-

- Hu et al. (2018) use local linear surrogate models, h_{GLM} , in h_{tree} leaf nodes to increase overall surrogate model fidelity while also retaining a high degree of interpretability.

4.1. Description

Following Friedman et al. (2001) a single feature $X_j \in \mathbf{X}$ and its complement set $X_{(-j)} \in \mathbf{X}$ (where $X_j \cup X_{(-j)} = \mathbf{X}$) is considered. $\text{PD}(X_j, g)$ for a given feature X_j is estimated as the average output of the learned function g when all the components of X_j are set to a constant $x \in \mathcal{X}$ and $X_{(-j)}$ is left untouched. $\text{ICE}(X_j, \mathbf{x}^{(i)}, g)$ for a given observation $\mathbf{x}^{(i)}$ and fea-

ture X_j is estimated as the output of the learned function g when $x_j^{(i)}$ is set to a constant $x \in \mathcal{X}$ and all $\mathbf{x}^{(i)} \in X_{(-j)}$ are left untouched. PD and ICE curves are usually plotted over some set of interesting constants $x \in \mathcal{X}$.

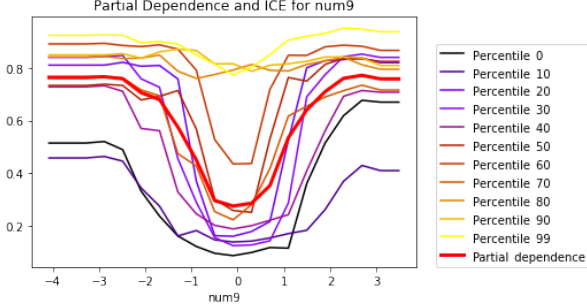


Figure 4. Partial dependence and ICE curves for $X_j = \text{num}_9$, for known signal generating function $f = \text{num}_1 * \text{num}_4 + |\text{num}_8| * \text{num}_9^2 + e$, and for machine-learned GBM response function g_{GBM} .

As in section 3, simulated data is used to highlight desirable characteristics of partial dependence and ICE plots. In Figure 4 partial dependence and ICE at the minimum, maximum, and each decile of g_{GBM} are plotted. The known quadratic behavior of num_9 is plainly visible, except for high value predictions, the 80th percentiles of $g_{\text{GBM}}(\mathbf{X})$ and above, and for $-1 < \text{num}_9 < 1$. When partial dependence and ICE curves diverge, this often points to an interaction that is being averaged out of the partial dependence. Given the form of 2, there is a known interaction between num_9 and num_8 and in Figure 3, every path through h_{tree} that includes a num_8 yields a higher prediction

the highest predictions from h_{tree} occurs when a decision path contains a split based

4.2. Recommendations

5. Local Interpretable Model-agnostic Explanations (LIME)

(Ribeiro et al., 2016) defines LIME for some observation $\mathbf{x} \in \mathcal{X}$:

$$\arg \max_{h \in \mathcal{H}} \mathcal{L}(g, h, \pi_{\mathbf{x}}) + \Omega(h) \quad (3)$$

Here g is the function to be explained, h is an interpretable surrogate model of g , often a linear model h_{GLM} , $\pi_{\mathbf{x}}$ is a weighting function over the domain of g , and $\Omega(h)$ limits the complexity of h .

Typically, h_{GLM} is constructed such that

$$\mathbf{X}^{(*)}, g(\mathbf{X}^{(*)}) \xrightarrow{\mathcal{A}_{\text{surrogate}}} h_{\text{GLM}} \quad (4)$$

where $\mathbf{X}^{(*)}$ is a generated sample, $\pi_{\mathbf{x}}$ weighs $\mathbf{X}^{(*)}$ samples by their Euclidean similarity to \mathbf{x} , local feature importance is estimated using $\beta_j x_j$, and L_1 regularization is used to induce a simplified, sparse h_{GLM} .

- LIME is ideal for creating low-fidelity, highly interpretable explanations for non-DT models and for neural network models trained on unstructured data, e.g. deep learning.
- Always use regression fit measures to assess the trustworthiness of LIME explanations.
- LIME can be difficult to deploy, but there are highly deployable variants. (Hu et al., 2018); (Hall et al., 2017)
- Local feature importance values are offsets from a local intercept.
 - Note that the intercept in LIME can account for the most important local phenomena.
 - Generated LIME samples can contain large proportions of out-of-range data that can lead to unrealistic intercept values.
- To increase the fidelity of LIME explanations, try LIME on discretized input features and on manually constructed interactions.
- Use cross-validation to construct standard deviations or even confidence intervals for local feature importance values.
- LIME can fail, particularly in the presence of extreme nonlinearity or high-degree interactions.

6. Tree Shap

Shapley explanations are a class of additive, consistent local feature importance measures with long-standing theoretical support, (Lundberg & Lee, 2017). For some observation $\mathbf{x} \in \mathcal{X}$, Shapley explanations take the form:

$$\phi_0 + \sum_{j=0}^{j=\mathcal{P}-1} \phi_j \mathbf{x}'_j \quad (5)$$

Here $\mathbf{x}' \in \{0, 1\}^{\mathcal{P}}$ is a binary representation of \mathbf{x} where 0 indicates missingness. Each ϕ_j is the local feature importance value associated with x_j .

- Calculating Shapley values directly is typically infeasible, but they can be estimated in different ways.
- Tree Shap is a specific implementation of Shapley explanations that leverages DT structures to disaggregate the contribution of each x_j to $g(\mathbf{x})$ in a DT or DT-based ensemble model. (Lundberg et al., 2018)
- Tree Shap is ideal for high-fidelity explanations of DT-based models, perhaps even in regulated applications.
- Local feature importance values are offsets from a global intercept.
- LIME can be constrained to become Shapley explanations, i.e. kernel shap.
- A similar, popular method known as *treeinterpreter* appears untrustworthy when applied to GBM models.

7. Use Case

8. General Recommendations

- Monotonically constrained XGBoost, Surrogate DT, PD and ICE plots, and Tree Shap are a direct and open source way to create an interpretable non-linear model.
- Global and local explanatory techniques are often necessary to explain a model.
- Use simpler low-fidelity or sparse explanations to understand more accurate and complex high-fidelity explanations.
- Seek consistent results across multiple explanatory techniques.
- Methods relying on generated data are sometimes unpalatable to users. They want to understand *their* data.
- Surrogate models can provide low-fidelity explanations for model mechanisms in original feature spaces if g is defined to include feature extraction or engineering.
- To increase adoption, production deployment of explanatory methods must be straightforward.

9. Suggested Reading

- xNN derivatives, neural net-specific methods.
- Accurate and interpretable classifiers.
- Fairness.

10. Software Resources

Comparison of Explanatory Techniques on Simulated Data:

https://github.com/h2oai/ml-resources/tree/master/lime_shap_treeint_compare

In-depth Explanatory Technique Examples:

https://github.com/jphall663/interpretable_machine_learning_with_python

"Awesome" Machine Learning Interpretability Resource List:

<https://github.com/jphall663/awesome-machine-learning-interpretability>

11. Acknowledgements

References

- Abu-Mostafa, Yaser S., Magdon-Ismael, Malik, and Lin, Hsuan-Tien. *Learning from Data*. AML-Book, New York, 2012. URL <https://work.caltech.edu/textbook.html>.
- Bastani, Osbert, Kim, Carolyn, and Bastani, Hamsa. Interpreting blackbox models via model extraction. *arXiv preprint arXiv:1705.08504*, 2017. URL <https://arxiv.org/pdf/1705.08504.pdf>.
- Breiman, Leo, Friedman, Jerome H., Olshen, Richard A., and Stone, Charles J. *Classification and Regression tTrees*. Routledge, 1984.
- Craven, Mark W. and Shavlik, Jude W. Extracting tree-structured representations of trained networks. *Advances in Neural Information Processing Systems*, 1996. URL <http://papers.nips.cc/paper/1152-extracting-tree-structured-representations-of-trained-networks.pdf>.
- Doshi-Velez, Finale and Kim, Been. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017. URL <https://arxiv.org/pdf/1702.08608.pdf>.

- Friedman, Jerome, Hastie, Trevor, and Tibshirani, Robert. *The Elements of Statistical Learning*. Springer, New York, 2001. URL <https://web.stanford.edu/~hastie/ElemStatLearn/printings/ESLII.print12.pdf>.
- Gilpin, Leilani H, Bau, David, Yuan, Ben Z., Bajwa, Ayesha, Specter, Michael, and Kagal, Lalana. Explaining explanations: An approach to evaluating interpretability of machine learning. *arXiv preprint arXiv:1806.00069*, 2018. URL <https://arxiv.org/pdf/1806.00069.pdf>.
- Goldstein, Alex, Kapelner, Adam, Bleich, Justin, and Pitkin, Emil. Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *Journal of Computational and Graphical Statistics*, 24(1), 2015.
- Guidotti, Riccardo, Monreale, Anna, Turini, Franco, Pedreschi, Dino, and Giannotti, Fosca. A survey of methods for explaining black box models. *arXiv preprint arXiv:1802.01933*, 2018. URL <https://arxiv.org/pdf/1802.01933.pdf>.
- Hall, Patrick, Gill, Navdeep, Kurka, Megan, and Phan, Wen. Machine learning interpretability with h2o driverless ai, 2017. URL <http://docs.h2o.ai/driverless-ai/latest-stable/docs/booklets/MLIBooklet.pdf>.
- Hu, Linwei, Chen, Jie, Nair, Vijayan N., and Sudjianto, Agus. Locally interpretable models and effects based on supervised partitioning (lime-sup). *arXiv preprint arXiv:1806.00663*, 2018. URL <https://arxiv.org/ftp/arxiv/papers/1806/1806.00663.pdf>.
- Lipton, Zachary C. The mythos of model interpretability. *arXiv preprint arXiv:1606.03490*, 2016. URL <https://arxiv.org/pdf/1606.03490.pdf>.
- Lundberg, Scott M and Lee, Su-In. A unified approach to interpreting model predictions. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 30*, pp. 4765–4774. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>.
- Lundberg, Scott M, Erion, Gabriel G, and Lee, Su-In. Consistent individualized feature attribution for tree ensembles. *arXiv preprint arXiv:1802.03888*, 2018. URL <https://arxiv.org/pdf/1706.06060.pdf>.
- Ribeiro, Marco Tulio, Singh, Sameer, and Guestrin, Carlos. Why should I trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1135–1144. ACM, 2016. URL <http://www.kdd.org/kdd2016/papers/files/rfp0573-ribeiroA.pdf>.
- Williams, Mike et al. *Interpretability*. Fast Forward Labs, 2017. URL <https://www.fastforwardlabs.com/research/FF06>.