

---

# Explanations for Multinomial Classifiers

Tips and Tricks for Practitioners

---

**Pramit Choudhary\***

Los Angeles, CA

pramit.choudhary@h2o.ai

**Navdeep Gill\***

Mountain View, CA

navdeep.gill@h2o.ai

**Patrick Hall†**

Washington, DC

patrick.hall@h2o.ai

## Abstract

## 1 Introduction

This short discussion bookends popular and practical texts on machine learning explanations by Chaudhary, Gill, Hall et al by specifically addressing the common and somewhat vexing problem of explaining the behavior and predictions of multinomial classifiers [7], [6], [2].

## 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 the 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  $\mathcal{Y}$  contained in a  $C$ -dimensional input space,  $\mathcal{Y} \subset \mathbb{R}^C$ .
- Learned output responses come from the set  $\hat{\mathcal{Y}}$ . For classification models the set  $\hat{\mathcal{Y}}$  typically contains a column vector for each unique class in  $\mathcal{Y}$ . In this text, the space  $\hat{\mathcal{Y}}$  is said to be contained in a  $C'$ -dimensional output space,  $\hat{\mathcal{Y}} \subset \mathbb{R}^{C'}$ .

### 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  $\mathbf{Y}$ ,  $\mathbf{y}^{(i)} = [y_0^{(i)}, y_1^{(i)}, \dots, y_{C-1}^{(i)}]$ , and corresponding predictions in  $\hat{\mathbf{Y}}$ ,  $\hat{\mathbf{y}}^{(i)} = [\hat{y}_0^{(i)}, \hat{y}_1^{(i)}, \dots, \hat{y}_{C'-1}^{(i)}]$ .
- $\mathbf{X}$  and  $\mathbf{Y}$  consist 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  $X_j = [x_j^{(0)}, x_j^{(1)}, \dots, x_j^{(N-1)}]^T$ .

---

\*H2O.ai

†H2O.ai and George Washington University

- A type of machine learning model  $g$ , selected from a hypothesis set  $\mathcal{H}$ , is trained to represent an unknown signal-generating function  $f$  observed as  $\mathbf{X}$  with labels  $\mathbf{Y}$  using a training algorithm  $\mathcal{A}$ :
- $g$  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  $g$ .

### 3.1 Decision Tree Surrogate

Given a learned function  $g$ , a set of learned output responses  $g(\mathbf{X}) = \hat{\mathbf{Y}}$ , and a tree splitting and pruning approach  $\mathcal{A}$ , a global – or over all  $\mathbf{X}$  – surrogate decision tree  $h_{\text{tree}}$  can be extracted such that  $h_{\text{tree}}(\mathbf{X}) \approx g(\mathbf{X})$ :

```

graph TD
    Node1["node #1  
rump ≤ 1.191  
giri = 0.667  
samples = 100.0%  
value = [0.333, 0.333, 0.333]  
class = 0.0"]
    Node2["node #2  
giri ≤ 0.5  
samples = 11.7%  
value = [0.154, 0.177, 0.666]  
class = 2.0"]
    Node3["node #3  
giri ≤ 0.543  
numf = 1.044  
giri = 0.662  
samples = 88.3%  
value = [0.362, 0.358, 0.28]  
class = 0.0"]
    Node4["node #4  
numf ≤ 0.789  
giri = 0.639  
samples = 73.2%  
value = [0.43, 0.39, 0.197]  
class = 0.0"]
    Node5["node #5  
numf ≤ 0.987  
giri = 0.619  
samples = 57.0%  
value = [0.403, 0.441, 0.156]  
class = 1.0"]
    Node6["node #6  
numf ≤ 0.029  
giri = 0.645  
samples = 16.2%  
value = [0.442, 0.234, 0.324]  
class = 0.0"]
    Node7["node #7  
giri ≤ 0.229  
samples = 7.9%  
value = [0.87, 0.112, 0.017]  
class = 0.0"]
    Node8["node #8  
giri ≤ 0.494  
samples = 8.3%  
value = [0.031, 0.351, 0.61]  
class = 2.0"]
    Node9["node #9  
numf ≤ 0.055  
giri = 0.62  
samples = 36.7%  
value = [0.286, 0.503, 0.21]  
class = 1.0"]
    Node10["node #10  
numf ≤ 0.065  
giri = 0.522  
samples = 20.3%  
value = [0.601, 0.336, 0.063]  
class = 0.0"]
    Node11["node #11  
numf ≤ 0.137  
giri = 0.515  
samples = 11.4%  
value = [0.589, 0.358, 0.042]  
class = 0.0"]
    Node12["node #12  
giri ≤ 0.308  
samples = 12.3%  
value = [0.869, 0.111, 0.021]  
class = 0.0"]
    Node13["node #13  
numf ≤ 0.785  
giri = 0.574  
samples = 25.3%  
value = [0.143, 0.567, 0.29]  
class = 1.0"]
    Node14["node #14  
numf ≤ 0.137  
giri = 0.515  
samples = 11.4%  
value = [0.589, 0.358, 0.042]  
class = 0.0"]
    Node15["node #15  
giri = 0.532  
samples = 9.3%  
value = [0.066, 0.351, 0.583]  
class = 0.0"]
    Node16["node #16  
giri = 0.495  
samples = 17.0%  
value = [0.192, 0.705, 0.103]  
class = 1.0"]
    Node17["node #17  
giri = 0.373  
samples = 7.8%  
value = [0.754, 0.205, 0.031]  
class = 0.0"]
    Node18["node #18  
giri = 0.248  
samples = 3.5%  
value = [0.062, 0.961, 0.077]  
class = 0.0"]

    Node1 -- True --> Node2
    Node1 -- False --> Node3
    Node2 -- True --> Node7
    Node2 -- False --> Node8
    Node3 -- True --> Node9
    Node3 -- False --> Node4
    Node4 -- True --> Node5
    Node4 -- False --> Node6
    Node5 -- True --> Node10
    Node5 -- False --> Node11
    Node6 -- True --> Node12
    Node6 -- False --> Node13
    Node7 -- True --> Node15
    Node7 -- False --> Node16
    Node8 -- True --> Node17
    Node8 -- False --> Node18
    Node9 -- True --> Node15
    Node9 -- False --> Node16
    Node10 -- True --> Node17
    Node10 -- False --> Node18
    Node11 -- True --> Node17
    Node11 -- False --> Node18
    Node12 -- True --> Node18
    Node12 -- False --> Node15
    Node13 -- True --> Node16
    Node13 -- False --> Node15
    Node14 -- True --> Node17
    Node14 -- False --> Node18
  
```

Prescribed methods for training  $h_{\text{tree}}$  do exist [3] [1]. In practice, straightforward cross-validation and pruning approaches are often sufficient. Moreover, comparing cross-validated training error to traditional training error can give an indication of the stability of the single decision tree  $h_{\text{tree}}$ .

Elegantly handles high cardinality targets.

What to do if very high cardinality:

- 2- or 3-D plot against most important variables
- 2- or 3-D plot against sparse, interpretable extracted features: NMF, Sparse PCA

What to do if very high cardinality:

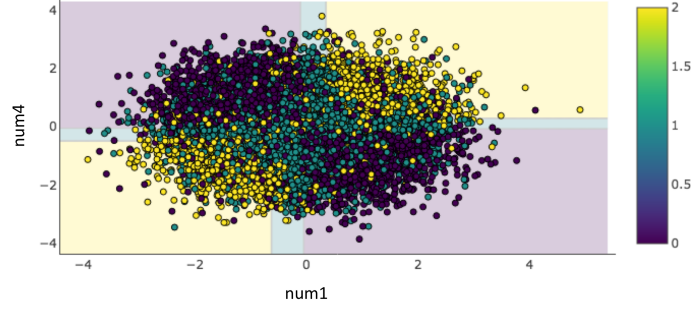


Figure 2:

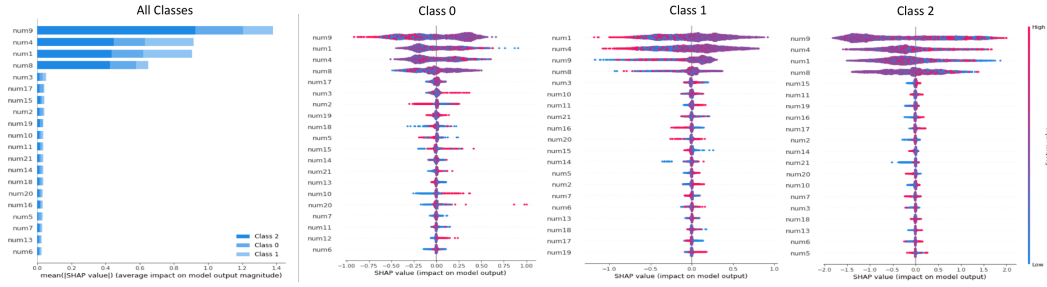


Figure 3:

- Examine top-K most frequent classes
- Examine top-K most accurate and inaccurate classes
- Examine classes with highest variance in  $\text{sum}(\text{absolute}(\text{shap}))$

### 3.4 Partial Dependence and ICE

Partial dependence (PD) plots are a widely-used method for describing the average predictions of a complex model  $g$  across some partition of data  $\mathbf{X}$  for some interesting input feature  $X_j$  [4]. Individual conditional expectation (ICE) plots are a newer method that describes the local behavior of  $g$  for a single instance  $\mathbf{x} \in \mathcal{X}$ . Partial dependence and ICE can be combined in the same plot to identify interactions modeled by  $g$  and to create a holistic portrait of the predictions of a complex model for some  $X_j$  [5].

What to do if very high cardinality:

- Examine top-K most frequent classes
- Examine top-K most accurate and inaccurate classes
- Examine classes with highest variance in partial dependence
- Examine classes with largest differences between partial dependence and ICE

### 3.5 Shapley Local Feature Importance

Shapley explanations, including tree shap and even certain implementations of LIME, are a class of additive, consistent local feature contribution measures with long-standing theoretical support [10]. Shapley explanations are the only possible locally accurate and consistent feature contribution values, meaning that Shapley explanation values for input features always sum to  $g(\mathbf{x})$  and that Shapley

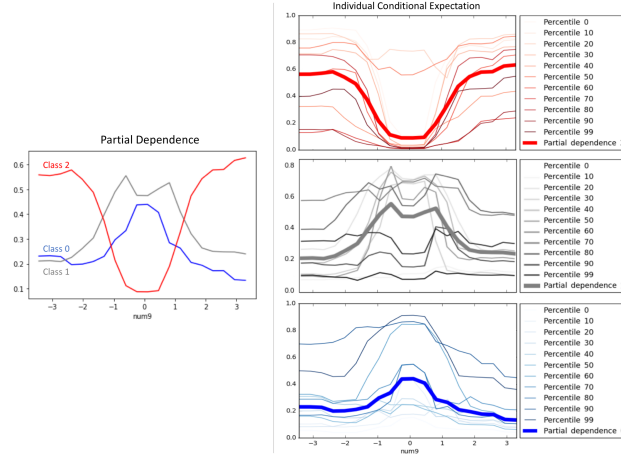


Figure 4:

explanation values can never decrease for some  $x_j$  when  $g$  is changed such that  $x_j$  truly makes a stronger contribution to  $g(\mathbf{x})$  [10].

$$g(\mathbf{x}) = \phi_0 + \sum_{j=0}^{j=\mathcal{P}-1} \phi_j \mathbf{z}_j \quad (2)$$

$$\phi_j = \sum_{S \subseteq \mathcal{P} \setminus \{j\}} \frac{|S|!(\mathcal{P} - |S| - 1)!}{\mathcal{P}!} [g_x(S \cup \{j\}) - g_x(S)] \quad (3)$$

Shapley values can be estimated in different ways. Tree shap is a specific implementation of Shapley explanations. It does not rely on surrogate models. Both tree shap and a related technique known as *treeinterpreter* rely instead on traversing internal tree structures to estimate the impact of each  $x_j$  for some  $g(\mathbf{x})$  of interest [9], [11].

What to do if very high cardinality:

## 4 Supplementary Materials

UCI credit card dataset [8].

<https://github.com/navdeep-G/interpretable-ml/tree/master/notebooks>

## 5 Conclusion

## 6 NIPS Style examples

**Paragraphs** There is also a `\paragraph` command available, which sets the heading in bold, flush left, and inline with the text, with the heading followed by 1 em of space.

### 6.1 Citations, figures, tables, references

These instructions apply to everyone.

## 6.2 Citations within the text

The `natbib` package will be loaded for you by default. Citations may be author/year or numeric, as long as you maintain internal consistency. As to the format of the references themselves, any style is acceptable as long as it is used consistently.

The documentation for `natbib` may be found at

<http://mirrors.ctan.org/macros/latex/contrib/natbib/natnotes.pdf>

Of note is the command `\citet`, which produces citations appropriate for use in inline text. For example,

```
\citet{hasselmo} investigated\dots
```

produces

Hasselmo, et al. (1995) investigated...

If you wish to load the `natbib` package with options, you may add the following before loading the `nips_2018` package:

```
\PassOptionsToPackage{options}{natbib}
```

If `natbib` clashes with another package you load, you can add the optional argument `nonatbib` when loading the style file:

```
\usepackage[nonatbib]{nips_2018}
```

As submission is double blind, refer to your own published work in the third person. That is, use “In the previous work of Jones et al. [4],” not “In our previous work [4].” If you cite your other papers that are not widely available (e.g., a journal paper under review), use anonymous author names in the citation, e.g., an author of the form “A. Anonymous.”

**New preprint option for 2018** If you wish to post a preprint of your work online, e.g., on arXiv, using the NIPS style, please use the `preprint` option. This will create a nonanonymized version of your work with the text “Preprint. Work in progress.” in the footer. This version may be distributed as you see fit. Please **do not** use the `final` option, which should **only** be used for papers accepted to NIPS.

At submission time, please omit the `final` and `preprint` options. This will anonymize your submission and add line numbers to aid review. Please *do not* refer to these line numbers in your paper as they will be removed during generation of camera-ready copies.

The file `nips_2018.tex` may be used as a “shell” for writing your paper. All you have to do is replace the author, title, abstract, and text of the paper with your own.

The formatting instructions contained in these style files are summarized in Sections ??, ??, and 6.1 below.

## 6.3 Footnotes

Footnotes should be used sparingly. If you do require a footnote, indicate footnotes with a number<sup>3</sup> in the text. Place the footnotes at the bottom of the page on which they appear. Precede the footnote with a horizontal rule of 2 inches (12 picas).

Note that footnotes are properly typeset *after* punctuation marks.<sup>4</sup>

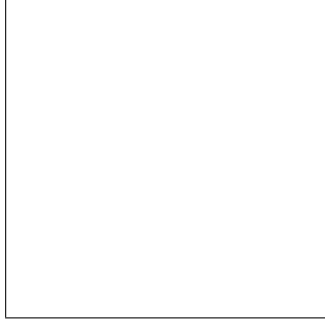


Figure 5: Sample figure caption.

Table 1: Sample table title

Part		
Name	Description	Size ( $\mu\text{m}$ )
Dendrite	Input terminal	$\sim 100$
Axon	Output terminal	$\sim 10$
Soma	Cell body	up to $10^6$

## 6.4 Figures

## 6.5 Tables

Place one line space before the table title, one line space after the table title, and one line space after the table. The table title must be lower case (except for first word and proper nouns); tables are numbered consecutively.

Note that publication-quality tables *do not contain vertical rules*. We strongly suggest the use of the `booktabs` package, which allows for typesetting high-quality, professional tables:

<https://www.ctan.org/pkg/booktabs>

This package was used to typeset Table 1.

The `\bbold` package almost always uses bitmap fonts. You should use the equivalent AMS Fonts:

```
\usepackage{amsfonts}
```

followed by, e.g., `\mathbb{R}`, `\mathbb{N}`, or `\mathbb{C}` for  $\mathbb{R}$ ,  $\mathbb{N}$  or  $\mathbb{C}$ . You can also use the following workaround for reals, natural and complex:

```
\newcommand{\RR}{\mathbb{R}} %real numbers
\newcommand{\Nat}{\mathbb{N}} %natural numbers
\newcommand{\CC}{\mathbb{C}} %complex numbers
```

Note that `amsfonts` is automatically loaded by the `amssymb` package.

## Acknowledgments

Use unnumbered third level headings for the acknowledgments. All acknowledgments go at the end of the paper. Do not include acknowledgments in the anonymized submission, only in the final paper.

---

<sup>3</sup>Sample of the first footnote.

<sup>4</sup>As in this example.

## References

- [1] Osbert Bastani, Carolyn Kim, and Hamsa Bastani. Interpreting Blackbox Models via Model Extraction. *arXiv preprint arXiv:1705.08504*, 2017. URL: <https://arxiv.org/pdf/1705.08504.pdf>.
- [2] Pramit Choudhary. Interpreting Predictive Models with Skater: Unboxing Model Opacity. *O'Reilly Ideas*, 2018. URL: <https://www.oreilly.com/ideas/interpreting-predictive-models-with-skater-unboxing-model-opacity>.
- [3] Mark W. Craven and Jude W. Shavlik. 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>.
- [4] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The Elements of Statistical Learning*. Springer, New York, 2001. URL: [https://web.stanford.edu/~hastie/ElemStatLearn/printings/ESLII\\_print12.pdf](https://web.stanford.edu/~hastie/ElemStatLearn/printings/ESLII_print12.pdf).
- [5] Alex Goldstein, Adam Kapelner, Justin Bleich, and Emil Pitkin. Peeking Inside the Black Box: Visualizing Statistical Learning with Plots of Individual Conditional Expectation. *Journal of Computational and Graphical Statistics*, 24(1), 2015. URL: <https://arxiv.org/pdf/1309.6392.pdf>.
- [6] Patrick Hall. On the Art and Science of Machine Learning Explanations. *arXiv preprint arXiv:1810.02909*, 2018. URL: <https://arxiv.org/pdf/1810.02909.pdf>.
- [7] Patrick Hall, Wen Phan, and Sri Satish Ambati. Ideas on Interpreting Machine Learning. *O'Reilly Ideas*, 2017. URL: <https://www.oreilly.com/ideas/ideas-on-interpreting-machine-learning>.
- [8] M. Lichman. UCI Machine Learning Repository, 2013. URL: <http://archive.ics.uci.edu/ml>.
- [9] Scott M. Lundberg, Gabriel G. Erion, and Su-In Lee. Consistent Individualized Feature Attribution for Tree Ensembles. *arXiv preprint arXiv:1802.03888*, 2018. URL: <https://arxiv.org/pdf/1706.06060.pdf>.
- [10] Scott M. Lundberg and Su-In Lee. A Unified Approach to Interpreting Model Predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc., 2017. URL: <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>.
- [11] Ando Saabas. Interpreting Random Forests, 2014. URL: <http://blog.datadive.net/interpreting-random-forests/>.