Front Matter
○
○
○

Surrogate DT
○○

PD and ICE
○○○

LIME
○○

Tree Shap
○○

Recommendations
○

Software
○

References

# Machine Learning Interpretability
## The Good, the Bad, and the Ugly

Patrick Hall
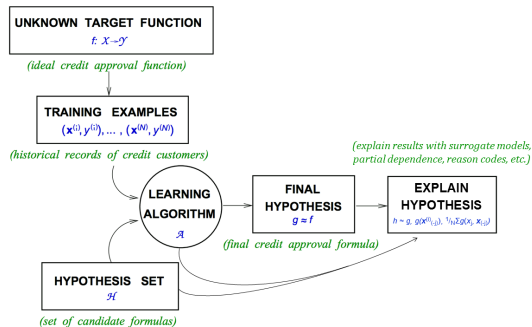
$H_2O$.ai

Aug. 1 2018

$H_2O$.ai

# Contents

H₂O.ai

# Obligatory Front Matter

- **What is interpretation?** "The ability to explain or to present in understandable terms to a human." Doshi-Velez and Kim, 2017
- **What is a good interpretation?** "When you can no longer keep asking why." Gilpin et al., 2018
- **Why should you care?**
  - Understanding of an impactful and quickly expanding set of technologies.
  - Addressing accidental or intentional discrimination.
  - Preventing malicious hacking and adversarial attacks.
  - Enabling regulatory compliance and increased financial margins.

$H_2O$ai

## Notation

- **Spaces.**
  - The input features come from a set $\mathcal{X}$ contained in a $P$-dimensional input space (i.e. $\mathcal{X} \subset \mathbb{R}^P$).
  - The output responses come from a set $\mathcal{Y}$ contained in a $C$-dimensional output space (i.e. $\mathcal{Y} \subset \mathbb{R}^C$).

- **Dataset.** A dataset $\mathbf{D}$ consists of $N$ tuples of observations: $[(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}), (\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \ldots, (\mathbf{x}^{(N-1)}, \mathbf{y}^{(N-1)})], \mathbf{x}^{(i)} \in \mathcal{X}, \mathbf{y}^{(i)} \in \mathcal{Y}$.

  - The input data $\mathbf{X}$ is composed of the set of row vectors $\mathbf{x}^{(i)}$.

    - let $\mathcal{P}$ be the set of features $\{X_0, X_1, \ldots, X_{P-1}\}$, where $X_j = \left[ x_j^{(0)}, x_j^{(1)}, \ldots, x_j^{(N-1)} \right]^T$.
    - then each $i$-th observation denoted as $\mathbf{x}^{(i)} = \left[ x_0^{(i)}, x_1^{(i)}, \ldots, x_{P-1}^{(i)} \right]$ is an instance of $\mathcal{P}$.

**H₂O**.ai

# Proposed Updates to the Learning Problem



The learning problem. Adapted from *Learning From Data*, Abu-Mostafa, Magdon-Ismail, and Lin, 2012.

Front Matter ○○○

Surrogate DT ●○

PD and ICE ○○○

LIME ○○

Tree Shap ○○

Recommendations ○

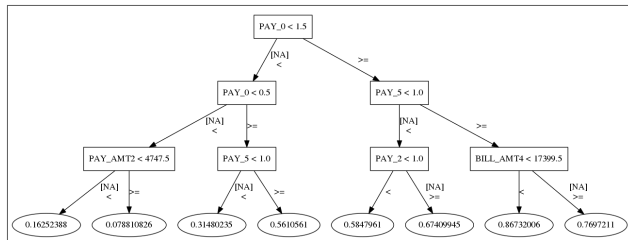Software ○

References

# Surrogate Decision Trees (DT)



Figure: $h_{\text{tree}}$ for Taiwanese credit card data Lichman, 2013, and for machine-learned GBM response function $g(X)$.

- Given a learned function $g$ and set of predictions, $g(\mathbf{X}) = \hat{\mathbf{Y}}$, a surrogate DT can be trained: $\mathbf{X}, \hat{\mathbf{Y}} \xrightarrow{\mathcal{A}_{\text{surrogate}}} h_{\text{tree}}$.

- $h_{\text{tree}}$ displays a low-fidelity flow chart of $g$'s decision making process, important features in $g$, and important interactions in $g$.

**H₂O**.ai

# Surrogate Decision Trees (DT)

- Always use error measures to assess the trustworthiness of $h_{tree}$.

- Prescribed methods (Craven and Shavlik, 1996; Bastani, Kim, and Bastani, 2017) for training $h_{tree}$ do exist. In practice, straightforward cross-validation approaches are typically sufficient.

- Comparing cross-validated error to standard training error can give an indication of the stability of the single tree model, $h_{tree}$.

- Hu et al., 2018 use local linear surrogate models, $h_{glm}$, in $h_{tree}$ leaf nodes to increase overall surrogate model accuracy while retaining a high degree of interpretability.

- $h_{tree}$ can provide low-fidelity explanations for model mechanisms in the original feature space if $g$ is defined to include feature extraction.

$H_2O$.ai

Front Matter
○
○
○

Surrogate DT
○○

PD and ICE
●○○

LIME
○○

Tree Shap
○○

Recommendations
○

Software
○

References

## Partial Dependence (PD) and Individual Conditional Expectation (ICE)

- Following Friedman, Hastie, and Tibshirani, 2001 a single feature $X_j \in \mathcal{P}$, a $P$-dimensional feature space, and its complement set $\mathcal{P}_{(-j)}$ (where $X_j \cup \mathcal{P}_{(-j)} = \mathcal{P}$) is considered.

- $PD(X_j, g)$ for a given feature $X_j$ is estimated as the average of the output of the learned function $g$, where all the components of $X_j$ are set to a constant $x_j^{(i)} \in X_j$, and $\mathcal{P}_{(-j)}$ is left untouched.

- $ICE(x_j^{(i)}, g)$ for a given row $\mathbf{x}^{(i)}$ and feature $X_j$ is estimated as the output of the learned function $g$ where $x_j^{(i)}$ is set to a constant $x_j^{(i)} \in X_j$ and $\mathbf{x}^{(i)} \in \mathcal{P}_{(-j)}$ are left untouched.

- PD and ICE are usually plotted over some set of interesting $x_j^{(i)} \in X_j$.

**H₂O**.ai

# Partial Dependence (PD) and Individual Conditional Expectation (ICE)



Figure: PD and ICE curves for $X_j = \text{num}_9$, for known signal generating function $f(X) = \text{num}_1 * \text{num}_4 + |\text{num}_8| * \text{num}_9^2 + e$, and for machine-learned GBM response function $g(X)$.

Overlaying PD and ICE curves is a succinct method for describing global and local prediction behavior and can be used to detect interactions. Goldstein et al., 2015

$H_2O$ ai

# Partial Dependence (PD) and Individual Conditional Expectation (ICE)



Figure: Surrogate DT, PD, and ICE curves for $X_j = \text{num}_9$, for known signal generating function $f(X) = \text{num}_1 * \text{num}_4 + |\text{num}_8| * \text{num}_9^2 + e$, and for machine-learned GBM response function $g(X)$.

Combining Surrogate DT models with PD and ICE curves is a convenient method for detecting, confirming, and understanding important interactions.

$\mathbf{H_2O}$.ai

Front Matter
○
○
○

Surrogate DT
○○

PD and ICE
○○○

LIME
●○

Tree Shap
○○

Recommendations
○

Software
○

References

# Local Interpretable Model-agnostic Explanations (LIME) - *Description*

# Local Interpretable Model-agnostic Explanations (LIME) - *Recommendations*

# Tree Shap - *Description*

Front Matter

Surrogate DT
○○

PD and ICE
○○○

LIME
○○

Tree Shap
○●

Recommendations
○

Software
○

References

# Tree Shap - *Recommendations*

## Closing Recommendations

- Monotonically constrained XGBoost, Surrogate DT, PD and ICE plots, and Tree Shap are a direct and open source way to create an interpretable nonlinear model.
- Global and local explanatory techniques are often necessary to explain a model.
- Simpler low-fidelity or sparse explanations should help in understanding more accurate and complex high-fidelity explanations.
- Seek consistency in results across multiple explanatory techniques.
- Methods that rely on generated data are sometimes unpalatable to users. They want to understand *their* data.
- Beware of uninterpretable features.
- Consider production deployment of explanatory techniques.

$H_2O$ ai

# Software Examples and Resources

**Comparison of Explanatory Techniques on Simulated Data:**
https://github.com/h2oai/mli-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

H$_2$O ai

# References I

Abu-Mostafa, Yaser S., Malik Magdon-Ismail, and Hsuan-Tien Lin (2012). **Learning from Data**. New York: AMLBook. URL: https://work.caltech.edu/textbook.html.

Bastani, Osbert, Carolyn Kim, and Hamsa Bastani (2017). "Interpreting blackbox models via model extraction." In: *arXiv preprint arXiv:1705.08504*. URL: https://arxiv.org/pdf/1705.08504.pdf.

Craven, Mark W. and Jude W. Shavlik (1996). "Extracting Tree-Structured Representations of Trained Networks." In: *Advances in Neural Information Processing Systems*. URL: http://papers.nips.cc/paper/1152-extracting-tree-structured-representations-of-trained-networks.pdf.

Doshi-Velez, Finale and Been Kim (2017). "Towards a rigorous science of interpretable machine learning." In: *arXiv preprint arXiv:1702.08608*. URL: https://arxiv.org/pdf/1702.08608.pdf.

Friedman, Jerome, Trevor Hastie, and Robert Tibshirani (2001). **The Elements of Statistical Learning**. New York: Springer. URL: https://web.stanford.edu/~hastie/ElemStatLearn/printings/ESLII_print12.pdf.

Gilpin, Leilani H et al. (2018). "Explaining Explanations: An Approach to Evaluating Interpretability of Machine Learning." In: *arXiv preprint arXiv:1806.00069*. URL: https://arxiv.org/pdf/1806.00069.pdf.

Goldstein, Alex et al. (2015). "Peeking Inside the Black Box: Visualizing Statistical Learning with Plots of Individual Conditional Expectation." In: *Journal of Computational and Graphical Statistics* 24.1.

$H_2O$ai

# References II

Hu, Linwei et al. (2018). "Locally Interpretable Models and Effects based on Supervised Partitioning (LIME-SUP)." In: *arXiv preprint arXiv:1806.00663*. URL: https://arxiv.org/ftp/arxiv/papers/1806/1806.00663.pdf.

Lichman, M. (2013). *UCI Machine Learning Repository*. URL: http://archive.ics.uci.edu/ml.

$H_2O$.ai