

# Guidelines for Responsible Use of Explainable Machine Learning

© Patrick Hall\*

H<sub>2</sub>O.ai

August 1, 2019

---

\*This material is shared under a [CC By 4.0 license](#) which allows for editing and redistribution, even for commercial purposes. However, any derivative work should attribute the author and H2O.ai.



# Contents

Introduction

Understanding and Trust

The Dark Side

Surrogates

High Stakes Applications



# What is explainable machine learning (ML)?

*“A collection of visual and/or interactive artifacts that provide a user with sufficient description of the model behavior to accurately perform tasks like evaluation, trusting, predicting, or improving the model.”*

— Sameer Singh, *UCI*

Variously defined along with aliases or similar concepts:

- “Towards a Rigorous Science of Interpretable Machine Learning” (Doshi-Velez and Kim [8])
- “Explaining Explanations” (Gilpin et al. [13])
- “A Survey Of Methods For Explaining Black Box Models” (Guidotti et al. [16])
- “The Mythos of Model Interpretability” (Lipton [23])
- *Interpretable Machine Learning* (Molnar [26])
- “Interpretable Machine Learning: Definitions, Methods, and Applications” (Murdoch et al. [27])
- “Challenges for Transparency” (Weller [39]).



# What is explainable ML?

What do *I* mean by explainable ML?

Mostly post-hoc techniques used to enhance *understanding* of trained model mechanisms and predictions, e.g. ...

- **Direct measures of global and local feature importance:**
  - Gradient-based feature attribution (Ancona et al. [2])
  - Shapley values (Lundberg and Lee [25])
- **Global and local surrogate models:**
  - Decision tree variants (Bastani, Pu, and Solar-Lezama [6], Craven and Shavlik [7])
  - Anchors (Ribeiro, Singh, and Guestrin [28])
  - Local interpretable model-agnostic explanations (LIME) (Ribeiro, Singh, and Guestrin [29])
- **Global and local visualizations of trained model predictions:**
  - Accumulated local effect (ALE) (Apley [4])
  - Partial dependence (Friedman, Hastie, and Tibshirani [11])
  - Individual conditional expectation (ICE) (Goldstein et al. [14])



# Why explainable ML?

Responsible Use of Explainable ML can enable:

- Human learning from machine learning
- Human appeal of automated decisions
- Regulatory compliance
- White-hat/red-team hacking

Misuse and Abuse of Explainable ML can enable:

- Model and data stealing (Tramèr et al. [36], Shokri et al. [34], Shokri, Strobel, and Zick [33])
- False justification for harmful black-boxes, e.g. “fairwashing” (Aïvodji et al. [1], Rudin [30])



## Proposed Guidelines for Responsible Use

Explainable ML is already in-use: numerous open source<sup>†</sup> and commercial packages<sup>‡</sup> are available today.

Best-practices are needed to prevent misuse and abuse. So, four basic guidelines are proposed here:

- Use explainable ML to enhance understanding.
- Learn how explainable ML is used for nefarious purposes.
- Augment surrogate models with direct explanations.
- Use highly transparent mechanisms for high stakes applications.

---

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

<sup>‡</sup>For instance Datarobot, H2O Driverless AI, SAS Visual Data Mining and Machine Learning, Zest AutoML



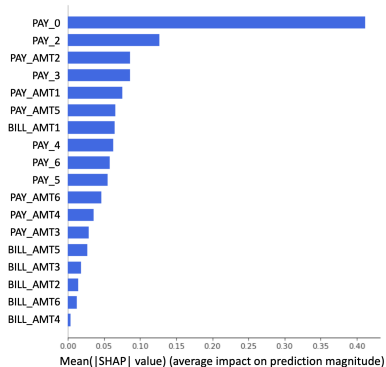
## Use Explainable ML to Enhance Understanding.

Explanations enhance understanding directly, and increase trust as a side-effect.

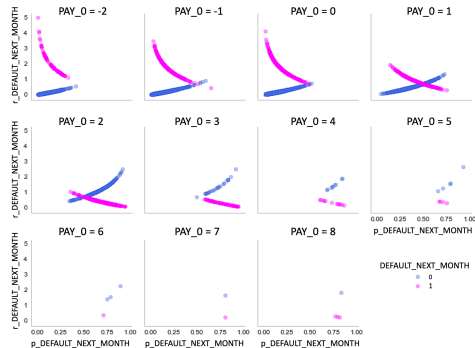
Models can be understood and not trusted, and trusted but not understood.

Explanations alone are neither necessary nor sufficient for trust.

# Understanding Without Trust



***g<sub>mono</sub>*** monotonically-constrained probability of default (PD) classifier trained on the UCI credit card dataset over-emphasizes the most important feature, a customer's most recent repayment status,  $PAY\_0$  [21].



***g<sub>mono</sub>*** also struggles to predict default for favorable statuses,  $-2 \leq PAY\_0 < 2$ , and often cannot predict on-time payment when recent payments are late,  $PAY\_0 \geq 2$ .





## Trust Without Understanding

Years before reliable explanation techniques were widely acknowledged and available, black-box predictive models, such as autoencoder and MLP neural networks, were used for fraud detection in the financial services industry [15]. When these models performed well, they were trusted.<sup>§</sup> However, they were not explainable or well-understood by contemporary standards.

---

<sup>§</sup>For example: [https://www.sas.com/en\\_ph/customers/hsbc.html](https://www.sas.com/en_ph/customers/hsbc.html),  
<https://www.kdnuggets.com/2011/03/sas-patent-fraud-detection.html>.

## Learn How Explainable ML is Used for Nefarious Purposes

When unintentionally misused, explainable ML can act as a faulty-safeguard for a potentially harmful black-box.

When intentionally abused, explainable ML can be used for:

- Hacking of data, models, or other intellectual property.
- *Fairwashing*, to mask the sociological biases of a discriminatory black-box.

# ML Hacking

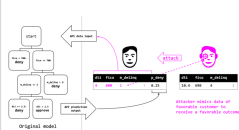
Many ML hacks use, or are exacerbated by, explainable ML techniques.

## Machine Learning Attack Cheatsheet

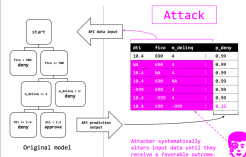
### Backdoors and Watermarks



### Impersonation



### Adversarial Examples



### Data Poisoning



### Model Inversion and Stealing



### Membership Inference



# White-hat Attacks

The flip-side of the dark side is community oversight of black-boxes.

Recent high profile analyses of commercial black-boxes, [Propublica and COMPAS](#) and [Gendershades and Rekognition](#), can be viewed as white-hat attacks (model stealing, adversarial examples) on proprietary black-boxes.<sup>||</sup>

---

<sup>||</sup>This presentation makes no claim on the quality of the analysis in Angwin et al. (2016), which has been criticized, but is simply stating that such cracking is possible [3], [10]. H<sub>2</sub>O.ai

## Explanation *is Not* a Front Line Fairness Tool

Use fairness tools, e.g. ...

- Disparate impact testing (Feldman et al. [9])
- Reweighting (Kamiran and Calders [18])
- Reject option based classification (Kamiran, Karim, and Zhang [19])
- Adversarial de-biasing (Zhang, Lemoine, and Mitchell [41])
- [aequitas](#), [ALF360](#), [Themis](#), [themis-ml](#)

... for fairness tasks: bias testing, bias remediation, and to establish trust.

Explanations can be used to understand and augment such results.

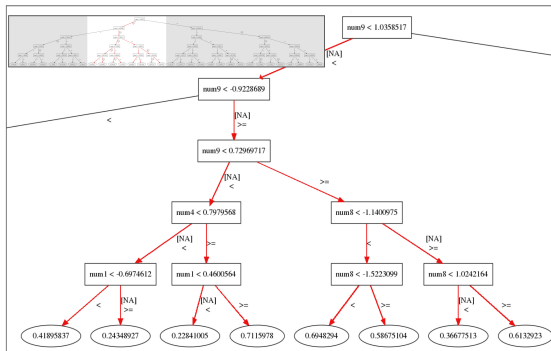
## Surrogate Models

Models of models, or surrogate models, can be helpful explanatory or modeling tools, but they are usually approximate, low-fidelity explainers.

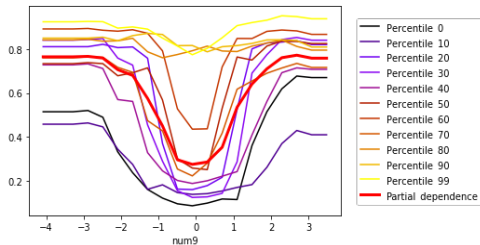
Much work in explainable ML has been directed toward improving the fidelity and usefulness of surrogate models (e.g., Bastani, Kim, and Bastani [5], Bastani, Pu, and Solar-Lezama [6], Craven and Shavlik [7], Hu et al. [17], Vaughan et al. [38])

***Many explainable ML techniques have nothing to do with surrogate models!***

# Augment Surrogate Models with Direct Explanations



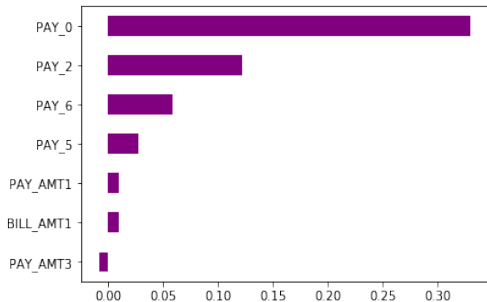
Naïve  $h_{\text{tree}}$ , a surrogate model, forms an approximate overall flowchart for the explained model,  $g_{\text{GBM}}$ .



Partial dependence and ICE curves generated *directly* from the explained model,  $g_{\text{GBM}}$ .

$h_{\text{tree}}$  displays known interactions in  $f = X_{\text{num}1} * X_{\text{num}4} + |X_{\text{num}8}| * X_{\text{num}9}^2$  for  $\sim -0.923 < X_{\text{num}9} < \sim 1.04$ . Modeling of the known interaction between  $X_{\text{num}9}$  and  $X_{\text{num}8}$  in  $f$  by  $g_{\text{GBM}}$  is confirmed by the divergence of partial dependence and ICE curves for  $\sim -1 < X_{\text{num}9} < \sim 1$ .

## Augment LIME with Direct Explanations



Locally-accurate Shapley contributions for a high risk individual's probability of default as predicted by a simple decision tree model,  $g_{tree}$ . See slide 19 for a directed graph representation of  $g_{tree}$ .

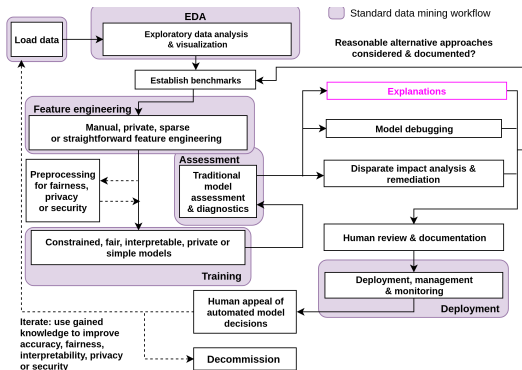
$h_{GLM}$ Feature	$h_{GLM}$ Coefficient
PAY_0 == 4	0.0009
PAY_2 == 3	0.0065
PAY_5 == 2	-0.0006
PAY_6 == 2	0.0036
BILL_AMT1	3.4339e-08
PAY_AMT1	4.8062e-07
PAY_AMT3	-5.867e-07

Coefficients for a local linear interpretable model,  $h_{GLM}$ , with an intercept of 0.77 and an  $R^2$  of 0.73., trained between the original inputs and predictions of  $g_{tree}$  for a segment of the UCI credit card dataset with late most recent repayment statuses,  $X_{PAY\_0} > 1$ .

Because  $h_{GLM}$  is relatively well-fit and has a logical intercept, it can be used along with Shapley values to reason about the modeled average behavior for risky customers and to differentiate the behavior of any one specific risky customer from their peers under the model.



# Use Highly Transparent Mechanisms for High Stakes Applications



A diagram of a proposed low risk ML workflow in which explanations (highlighted in **fuschia**) are used along with interpretable or white-box models, disparate impact analysis and remediation techniques, and other review and appeal mechanisms to create a fair, accountable, and transparent ML system.

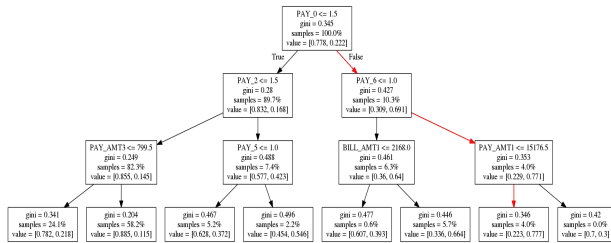
## Interpretable Models

Many types of novel interpretable models are freely available today, e.g.

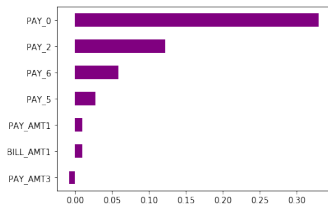
- Explainable boosting machine (EBM)
- Monotonic GBM in [h2o](#) or [XGBoost](#)
- [RuleFit](#) (Friedman, Popescu, et al. [12])
- [Super-sparse linear integer model](#) (SLIM) (Ustun and Rudin [37])
- [Scalable Bayesian rule list](#) (Yang, Rudin, and Seltzer [40])



# Explanations and Interpretable Models are Not Mutually Exclusive



Simple decision tree, *gtree*, trained on the UCI credit card data to predict default with validation AUC of 0.74. The decision policy for a high risk individual is highlighted in red.



Locally-accurate Shapley contributions for the highlighted individual's probability of default. See slide 16 for LIMEs for the high risk customers in *gtree*.

The Shapley values are helpful because they highlight the local importance of features not on the decision path, which could be underestimated by examining the decision policy alone.



# An Ode to the Shapley Value

1. **In the beginning:** A Value for N-Person Games, 1953
2. **Nobel-worthy contributions:** *The Shapley value: Essays in honor of Lloyd S. Shapley*, 1988
3. **Shapley regression:** Analysis of Regression in Game Theory Approach, 2001
4. **First reference in ML?** Fair Attribution of Functional Contribution in Artificial and Biological Networks, 2004
5. **Into the ML research mainstream, i.e. JMLR:** An Efficient Explanation of Individual Classifications using Game Theory, 2010
6. **Into the real-world data mining workflow ... finally:** Consistent Individualized Feature Attribution for Tree Ensembles, 2017. \*\*
7. **Unification:** A Unified Approach to Interpreting Model Predictions, 2017. ††

---

\*\* See [h2o](#), [LightGBM](#), or [XGBoost](#) for implementation.

†† See [shap](#) for implementation.



## Explanation and Fairness Techniques are Not Mutually Exclusive

	Adverse Impact Disparity	Accuracy Disparity	TPR Disparity	TNR Disparity	FPR Disparity	FNR Disparity
single	0.89	1.03	0.99	1.03	0.85	1.01
divorced	1.01	0.93	0.81	0.96	1.25	1.22
other	0.26	1.12	0.62	1.17	0	1.44

Basic group disparity metrics across different marital statuses for monotonically constrained GBM model,  $g_{\text{mono}}$ , trained on the UCI credit card dataset. See slide 8 for global Shapley feature importance for  $g_{\text{mono}}$  and slide 13 for an important note about explanation as fairness techniques.

Many fairness techniques are freely available today: [aequitas](#), [AlF360](#), [Themis](#), [themis-ml](#).

Traditional disparate impact testing tools are best-suited for constrained models because average group metrics cannot reliably identify local instances of discrimination that can occur when using complex, unconstrained models.

# References

## This presentation:

[https://www.github.com/jphall663kdd\\_2019](https://www.github.com/jphall663kdd_2019)

## Code examples for this presentation:

[https://www.github.com/jphall663/interpretable\\_machine\\_learning\\_with\\_python](https://www.github.com/jphall663/interpretable_machine_learning_with_python)

[https://www.github.com/jphall663/responsible\\_xai](https://www.github.com/jphall663/responsible_xai)

## Associated texts:

<https://arxiv.org/pdf/1810.02909.pdf>

<https://arxiv.org/pdf/1906.03533.pdf>



## References

- [1] Ulrich Aïvodji et al. “Fairwashing: the Risk of Rationalization.” In: *arXiv preprint arXiv:1901.09749* (2019). URL: <https://arxiv.org/pdf/1901.09749.pdf>.
- [2] Marco Ancona et al. “Towards Better Understanding of Gradient-based Attribution Methods for Deep Neural Networks.” In: *6th International Conference on Learning Representations (ICLR 2018)*. URL: [https://www.research-collection.ethz.ch/bitstream/handle/20.500.11850/249929/Flow\\_ICLR\\_2018.pdf](https://www.research-collection.ethz.ch/bitstream/handle/20.500.11850/249929/Flow_ICLR_2018.pdf). 2018.
- [3] Julia Angwin et al. “Machine Bias: There’s Software Used Across the Country to Predict Future Criminals. And It’s Biased Against Blacks.” In: *ProPublica* (2016). URL: <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.
- [4] Daniel W. Apley. “Visualizing the Effects of Predictor Variables in Black Box Supervised Learning Models.” In: *arXiv preprint arXiv:1612.08468* (2016). URL: <https://arxiv.org/pdf/1612.08468.pdf>.
- [5] Osbert Bastani, Carolyn Kim, and Hamsa Bastani. “Interpreting Blackbox Models via Model Extraction.” In: *arXiv preprint arXiv:1705.08504* (2017). URL: <https://arxiv.org/pdf/1705.08504.pdf>.



## References

- [6] Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. “Verifiable Reinforcement Learning Via Policy Extraction.” In: *Advances in Neural Information Processing Systems*. URL: <http://papers.nips.cc/paper/7516-verifiable-reinforcement-learning-via-policy-extraction.pdf>. 2018, pp. 2494–2504.
- [7] Mark W. Craven and Jude W. Shavlik. “Extracting Tree-Structured Representations of Trained Networks.” In: *Advances in Neural Information Processing Systems* (1996). URL: <http://papers.nips.cc/paper/1152-extracting-tree-structured-representations-of-trained-networks.pdf>.
- [8] Finale Doshi-Velez and Been Kim. “Towards a Rigorous Science of Interpretable Machine Learning.” In: *arXiv preprint arXiv:1702.08608* (2017). URL: <https://arxiv.org/pdf/1702.08608.pdf>.
- [9] Michael Feldman et al. “Certifying and Removing Disparate Impact.” In: *Proceedings of the 21<sup>st</sup> ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. URL: <https://arxiv.org/pdf/1412.3756.pdf>. ACM. 2015, pp. 259–268.





## References

- [10] Anthony W. Flores, Kristin Bechtel, and Christopher T. Lowenkamp. “False Positives, False Negatives, and False Analyses: A Rejoinder to Machine Bias: There’s Software Used across the Country to Predict Future Criminals. And It’s Biased against Blacks.” In: *Fed. Probation* 80 (2016). URL: <https://bit.ly/2Gesf9Y>, p. 38.
- [11] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. ***The Elements of Statistical Learning***. URL: [https://web.stanford.edu/~hastie/ElemStatLearn/printings/ESLII\\_print12.pdf](https://web.stanford.edu/~hastie/ElemStatLearn/printings/ESLII_print12.pdf). New York: Springer, 2001.
- [12] Jerome H. Friedman, Bogdan E Popescu, et al. “Predictive Learning via Rule Ensembles.” In: *The Annals of Applied Statistics* 2.3 (2008). URL: [https://projecteuclid.org/download/pdfview\\_1/euclid.aos/1223908046](https://projecteuclid.org/download/pdfview_1/euclid.aos/1223908046), pp. 916–954.
- [13] Leilani H. Gilpin et al. “Explaining Explanations: An Approach to Evaluating Interpretability of Machine Learning.” In: *arXiv preprint arXiv:1806.00069* (2018). URL: <https://arxiv.org/pdf/1806.00069.pdf>.
- [14] Alex Goldstein et al. “Peeking Inside the Black Box: Visualizing Statistical Learning with Plots of Individual Conditional Expectation.” In: *Journal of Computational and Graphical Statistics* 24.1 (2015). URL: <https://arxiv.org/pdf/1309.6392.pdf>.



## References

- [15] Krishna M. Gopinathan et al. *Fraud Detection using Predictive Modeling*. US Patent 5,819,226. URL: <https://patents.google.com/patent/US5819226A>. 1998.
- [16] Riccardo Guidotti et al. “A Survey of Methods for Explaining Black Box Models.” In: *ACM Computing Surveys (CSUR)* 51.5 (2018). URL: <https://arxiv.org/pdf/1802.01933.pdf>, p. 93.
- [17] Linwei Hu et al. “Locally Interpretable Models and Effects Based on Supervised Partitioning (LIME-SUP).” In: *arXiv preprint arXiv:1806.00663* (2018). URL: <https://arxiv.org/ftp/arxiv/papers/1806/1806.00663.pdf>.
- [18] Faisal Kamiran and Toon Calders. “Data Preprocessing Techniques for Classification Without Discrimination.” In: *Knowledge and Information Systems* 33.1 (2012). URL: <https://link.springer.com/content/pdf/10.1007/s10115-011-0463-8.pdf>, pp. 1–33.
- [19] Faisal Kamiran, Asim Karim, and Xiangliang Zhang. “Decision Theory for Discrimination-aware Classification.” In: *2012 IEEE 12th International Conference on Data Mining*. URL: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.722.3030&rep=rep1&type=pdf>. IEEE. 2012, pp. 924–929.



## References

- [20] Alon Keinan et al. “Fair Attribution of Functional Contribution in Artificial and Biological Networks.” In: *Neural Computation* 16.9 (2004). URL: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.436.6801&rep=rep1&type=pdf>, pp. 1887–1915.
- [21] M. Lichman. *UCI Machine Learning Repository*. URL: <http://archive.ics.uci.edu/ml>. 2013.
- [22] Stan Lipovetsky and Michael Conklin. “Analysis of Regression in Game Theory Approach.” In: *Applied Stochastic Models in Business and Industry* 17.4 (2001), pp. 319–330.
- [23] Zachary C. Lipton. “The Mythos of Model Interpretability.” In: *arXiv preprint arXiv:1606.03490* (2016). URL: <https://arxiv.org/pdf/1606.03490.pdf>.
- [24] Scott M. Lundberg, Gabriel G. Erion, and Su-In Lee. “Consistent Individualized Feature Attribution for Tree Ensembles.” In: *Proceedings of the 2017 ICML Workshop on Human Interpretability in Machine Learning (WHI 2017)*. Ed. by Been Kim et al. URL: <https://openreview.net/pdf?id=ByTKSo-m->. ICML WHI 2017, 2017, pp. 15–21.

## References

- [25] Scott M. Lundberg and Su-In Lee. “A Unified Approach to Interpreting Model Predictions.” In: *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon et al. URL: <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>. Curran Associates, Inc., 2017, pp. 4765–4774.
- [26] Christoph Molnar. ***Interpretable Machine Learning***. URL: <https://christophm.github.io/interpretable-ml-book/>. christophm.github.io, 2018.
- [27] W. James Murdoch et al. “Interpretable Machine Learning: Definitions, Methods, and Applications.” In: *arXiv preprint arXiv:1901.04592* (2019). URL: <https://arxiv.org/pdf/1901.04592.pdf>.
- [28] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “Anchors: High-Precision Model-agnostic Explanations.” In: *AAAI Conference on Artificial Intelligence*. URL: <https://homes.cs.washington.edu/~marcotcr/aaai18.pdf>. 2018.
- [29] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “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*. URL: <http://www.kdd.org/kdd2016/papers/files/rfp0573-ribeiroA.pdf>. ACM. 2016, pp. 1135–1144.



## References

- [30] Cynthia Rudin. “Please Stop Explaining Black Box Models for High Stakes Decisions.” In: *arXiv preprint arXiv:1811.10154* (2018). URL: <https://arxiv.org/pdf/1811.10154.pdf>.
- [31] Lloyd S. Shapley. “A Value for N-Person Games.” In: *Contributions to the Theory of Games* 2.28 (1953). URL: <http://www.library.fa.ru/files/Roth2.pdf#page=39>, pp. 307–317.
- [32] Lloyd S. Shapley, Alvin E. Roth, et al. *The Shapley value: Essays in honor of Lloyd S. Shapley*. URL: <http://www.library.fa.ru/files/Roth2.pdf>. Cambridge University Press, 1988.
- [33] Reza Shokri, Martin Strobel, and Yair Zick. “Privacy Risks of Explaining Machine Learning Models.” In: *arXiv preprint arXiv:1907.00164* (2019). URL: <https://arxiv.org/pdf/1907.00164.pdf>.
- [34] Reza Shokri et al. “Membership Inference Attacks Against Machine Learning Models.” In: *2017 IEEE Symposium on Security and Privacy (SP)*. URL: <https://arxiv.org/pdf/1610.05820.pdf>. IEEE. 2017, pp. 3–18.
- [35] Erik Strumbelj and Igor Kononenko. “An Efficient Explanation of Individual Classifications using Game Theory.” In: *Journal of Machine Learning Research* 11.Jan (2010). URL: <http://www.jmlr.org/papers/volume11/strumbelj10a/strumbelj10a.pdf>, pp. 1–18.



## References

- [36] Florian Tramèr et al. “Stealing Machine Learning Models via Prediction APIs.” In: *25th {USENIX} Security Symposium ({USENIX} Security 16)*. URL: [https://www.usenix.org/system/files/conference/usenixsecurity16/sec16\\_paper\\_tramer.pdf](https://www.usenix.org/system/files/conference/usenixsecurity16/sec16_paper_tramer.pdf). 2016, pp. 601–618.
- [37] Berk Ustun and Cynthia Rudin. “Supersparse Linear Integer Models for Optimized Medical Scoring Systems.” In: *Machine Learning* 102.3 (2016). URL: <https://users.cs.duke.edu/~cynthia/docs/UstunTrRuAAAI13.pdf>, pp. 349–391.
- [38] Joel Vaughan et al. “Explainable Neural Networks Based on Additive Index Models.” In: *arXiv preprint arXiv:1806.01933* (2018). URL: <https://arxiv.org/pdf/1806.01933.pdf>.
- [39] Adrian Weller. “Challenges for Transparency.” In: *arXiv preprint arXiv:1708.01870* (2017). URL: <https://arxiv.org/pdf/1708.01870.pdf>.
- [40] Hongyu Yang, Cynthia Rudin, and Margo Seltzer. “Scalable Bayesian Rule Lists.” In: *Proceedings of the 34th International Conference on Machine Learning (ICML)*. URL: <https://arxiv.org/pdf/1602.08610.pdf>. 2017.

## References

- [41] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. “Mitigating Unwanted Biases with Adversarial Learning.” In: *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*. URL: <https://arxiv.org/pdf/1801.07593.pdf>. ACM. 2018, pp. 335–340.