

Training Characteristic Functions with Reinforcement Learning

Stephan Wäldchen, Felix Huber, Sebastian Pokutta

Zuse Institut Berlin

39th International Conference on Machine Learning

July 17 – 23, 2022

How do we interpret black-box Functions?

- | Core idea: Probe black-box function with different inputs

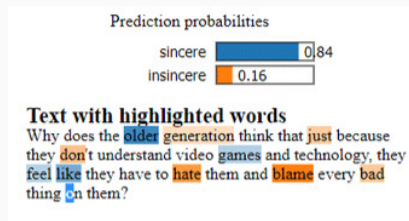


How do we interpret black-box Functions?

- | Core idea: Probe black-box function with different inputs
- | Local interpretation (for specific input \mathbf{x}): vary features



- | Core idea: Probe black-box function with different inputs
- | Local interpretation (for specific input \mathbf{x}): vary features

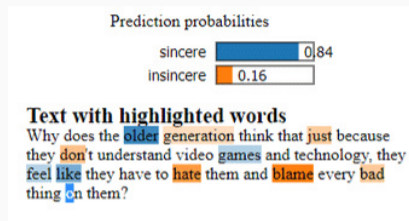


LIME saliency¹

1) "Explainable ai: A review of machine learning interpretability methods", Linardatos et al. [7]

How do we interpret black-box Functions?

- | Core idea: Probe black-box function with different inputs
- | Local interpretation (for specific input \mathbf{x}): vary features
- | ML: Saliency & Coop. Game Theory: Surplus Attribution

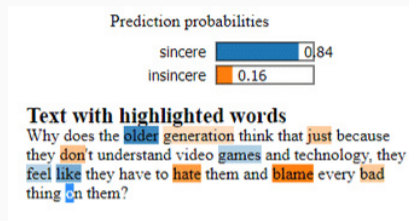


LIME saliency¹

1) "Explainable ai: A review of machine learning interpretability methods", Linardatos et al. [7]

How do we interpret black-box Functions?

- | Core idea: Probe black-box function with different inputs
- | Local interpretation (for specific input \mathbf{x}): vary features
- | ML: Saliency & Coop. Game Theory: Surplus Attribution
- | Characteristic function: $\mathcal{F} : 2^{[d]} \rightarrow \mathbb{R}$



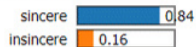
LIME saliency¹

1) "Explainable ai: A review of machine learning interpretability methods", Linardatos et al. [7]

How do we interpret black-box Functions?

- Core idea: Probe black-box function with different inputs
- Local interpretation (for specific input x): vary features
- ML: Saliency & Coop. Game Theory: Surplus Attribution
- Characteristic function: $v: 2^d \rightarrow \mathbb{R}$

Prediction probabilities



Text with highlighted words

Why does the **older** generation think that **just** because they **don't** understand video **games** and technology, they **feel like** they have to **hate** them and **blame** every **bad** thing **on** them?

LIME saliency¹



How do we interpret black-box Functions?

- | Core idea: Probe black-box function with different inputs
- | Local interpretation (for specific input \mathbf{x}): vary features
- | ML: Saliency $\&$ Coop. Game Theory: Surplus Attribution
- | Characteristic function: $\mathcal{S} : 2^{[d]} \rightarrow \mathbb{R}$
- | Prime Implicant Explanation² (mostly for $\mathcal{S} : 2^{[d]} \rightarrow \{0,1\}$)



Source: <https://clearcode.cc/blog/game-theory-attribution/>

$$S = \underset{S \subseteq [d]}{\operatorname{argmin}} |S^j| \quad \text{s.t.} \quad (S) = ([d])$$

1) "Explainable ai: A review of machine learning interpretability methods", Linardatos et al. [7] 2) "A symbolic approach to explaining bayesian network classifiers", Shih et al. [13]

How do we interpret black-box Functions?

- Core idea: Probe black-box function with different inputs
- Local interpretation (for specific input \mathbf{x}): vary features
- ML: Saliency $\&$ Coop. Game Theory: Surplus Attribution
- Characteristic function: $v: 2^{[d]} \rightarrow \mathbb{R}$
- Prime Implicant Explanation² (mostly for $v: 2^{[d]} \rightarrow \{0,1\}$)



Source: <https://clearcode.cc/blog/game-theory-attribution/>

$$S = \operatorname{argmin}_{S \subseteq [d]} |S| \text{ s.t. } v(S) = v([d])$$

- Shapley Values³ (linear, efficient, symmetric, null-player)

$$\phi_i = \frac{1}{d!} \sum_{\pi \in \Pi([d])} (v(P_i \cup \{i\} \cup \pi) - v(P_i \cup \pi))$$

1) "Explainable ai: A review of machine learning interpretability methods", Linardatos et al. [7] 2) "A symbolic approach to explaining bayesian network classifiers", Shih et al. [13] 3) "A value for n-person games" Shapley [12]

How do we interpret black-box Functions?

- | Core idea: Probe black-box function with different inputs
- | Local interpretation (for specific input \mathbf{x}): vary features
- | ML: Saliency $\&$ Coop. Game Theory: Surplus Attribution
- | Characteristic function: $v: 2^{[d]} \rightarrow \mathbb{R}$
- | Prime Implicant Explanation² (mostly for $v: 2^{[d]} \rightarrow [0; 1]$)



Source: <https://clearcode.cc/blog/game-theory-attribution/>

$$S = \operatorname{argmin}_{S \subseteq [d]} |S| \quad \text{s.t.} \quad v(S) = v([d])$$

- | Shapley Values³ (linear, efficient, symmetric, null-player)

$$\phi_i = \frac{1}{d!} \sum_{\pi \in \Pi([d])} (v(P_i \cup \{i\} \cup \pi) - v(P_i \cup \pi))$$

- | **Problem: We don't have characteristic functions!**

1) "Explainable ai: A review of machine learning interpretability methods", Linardatos et al. [7] 2) "A symbolic approach to explaining bayesian network classifiers", Shih et al. [13] 3) "A value for n-person games" Shapley [12]

| Approach from Lundberg et al¹:

$$\phi_{\mathbf{x}}(S) = E_{\mathbf{y}}[\mathbf{y} | \mathbf{y}_S = \mathbf{x}_S] = \int \mathbf{x} dP[\mathbf{x}_S | \mathbf{x}_S]:$$

1) "A unified approach to interpreting model predictions", Lundberg et al, [8]

- | Approach from Lundberg et al¹:

$$\phi_{\mathbf{x}}(S) = E_{\mathbf{y}}[\mathbf{y} | \mathbf{y}_S = \mathbf{x}_S] = \int \mathbf{x} dP[\mathbf{x}_{S^c} | \mathbf{x}_S]:$$

- | Needs good model of $P[\mathbf{x}_{S^c} | \mathbf{x}_S]$!

1) "A unified approach to interpreting model predictions", Lundberg et al, [8]

- | Approach from Lundberg et al¹:

$$\phi_{\mathbf{x}}(S) = E_{\mathbf{y}}[\mathbf{y} | \mathbf{y}_S = \mathbf{x}_S] = \int \mathbf{x} dP[\mathbf{x}_{S^c} | \mathbf{x}_S]:$$

- | Needs good model of $P[\mathbf{x}_{S^c} | \mathbf{x}_S]$!
- | Most methods simply approximate with baseline values (sometimes layer-wise)

1) "A unified approach to interpreting model predictions", Lundberg et al, [8]

- | Approach from Lundberg et al¹:

$$\phi_{\mathbf{x}}(S) = E_{\mathbf{y}}[\mathbf{y} | \mathbf{y}_S = \mathbf{x}_S] = \int \mathbf{y} dP[\mathbf{y}_{S^c} | \mathbf{x}_S]:$$

- | Needs good model of $P[\mathbf{x}_{S^c} | \mathbf{x}_S]$!
- | Most methods simply approximate with baseline values (sometimes layer-wise)
- | Change off-manifold behaviour to manipulate:
Gradient, Integrated gradients^{2;3}, LRP^{2;4;7},
LIME^{3;5}, DeepShap^{3;5}, Grad-Cam⁷,
Shapley-based⁶, Counterfactual explanations⁸,

1) "A unified approach to interpreting model predictions", Lundberg et al, [8] 2) "Fairwashing explanations with off-manifold detergent", Anders et al. [1] 3) "You Shouldn't Trust Me: Learning Models Which Conceal Unfairness From Multiple Explanation Methods", Dimanov et al. [3] 4) "Explanations can be manipulated and geometry is to blame" Dombrowski et al. [4] 5) "Fooling lime and shap: Adversarial attacks on post hoc explanation methods", Slack et al. [14] 6) "Shapley explainability on the data manifold", Frye et al. [5] 7) "Fooling neural network interpretations via adversarial model manipulation" Heo et al. [6] 8) "Counterfactual Explanations Can Be Manipulated" Slack et al. [15]

- Approach from Lundberg et al¹:

$$\phi_{\mathbf{x}}(S) = E_{\mathbf{y}}[\mathbf{y} | \mathbf{y}_S = \mathbf{x}_S] = \int \mathbf{y} dP[\mathbf{y}_{S^c} | \mathbf{x}_S]:$$

Image

FW

AFW

- Needs good model of $P[\mathbf{x}_{S^c} | \mathbf{x}_S]$!

LCG

LAFW

Sensitivity

- Most methods simply approximate with baseline values (sometimes layer-wise)

- Change manifold behaviour to manipulate:

Gradient, Integrated gradients^{2;3}, LRP^{2;4;7},

LIME^{3;5}, DeepShap^{3;5}, Grad-Cam⁷,

Shapley-based⁶, Counterfactual explanations⁸,

Best Performer RDE creates new features!⁸

1) "A unified approach to interpreting model predictions", Lundberg et al. [8] 2) "Fairwashing explanations with off-manifold detergent", Anders et al. [1] 3) "You Shouldn't Trust Me: Learning Models Which Conceal Unfairness From Multiple Explanation Methods", Dimanov et al. [3] 4) "Explanations can be manipulated and geometry is to blame" Dombrowski et al. [4] 5) "Fooling lime and shap: Adversarial attacks on post hoc explanation methods", Slack et al. [14] 6) "Shapley explainability on the data manifold", Frye et al. [5] 7) "Fooling neural network interpretations via adversarial model manipulation" Heo et al. [6] 8) "Counterfactual Explanations Can Be Manipulated" Slack et al. [15] 9) "Interpretable Neural Networks with Frank-Wolfe: Sparse Relevance Maps and Relevance Orderings", Macdonald et al. [9]

- Approach from Lundberg et al¹:

$$\phi_{\mathbf{x}}(S) = E_{\mathbf{y}}[\langle \mathbf{y} \rangle_{\mathbf{y}_S = \mathbf{x}_S}] = \int_{\mathcal{Z}} \langle \mathbf{x} \rangle dP[\mathbf{x}_{S^c} | \mathbf{x}_S]:$$

Image

FW

AFW

- Needs good model of $P[\mathbf{x}_{S^c} | \mathbf{x}_S]$!

LCG

LAFW

Sensitivity

- Most methods simply approximate with baseline values (sometimes layer-wise)

- Change ϕ -manifold behaviour to manipulate:

Gradient, Integrated gradients^{2;3}, LRP^{2;4;7},

LIME^{3;5}, DeepShap^{3;5}, Grad-Cam⁷,

Shapley-based⁶, Counterfactual explanations⁸,

Idea: Directly train a characteristic function!

1) "A unified approach to interpreting model predictions", Lundberg et al. [8] 2) "Fairwashing explanations with off-manifold detergent", Anders et al. [1] 3) "You Shouldn't Trust Me: Learning Models Which Conceal Unfairness From Multiple Explanation Methods", Dimanov et al. [3] 4) "Explanations can be manipulated and geometry is to blame" Dombrowski et al. [4] 5) "Fooling lime and shap: Adversarial attacks on post hoc explanation methods", Slack et al. [14] 6) "Shapley explainability on the data manifold", Frye et al. [5] 7) "Fooling neural network interpretations via adversarial model manipulation" Heo et al. [6] 8) "Counterfactual Explanations Can Be Manipulated" Slack et al. [15] 9) "Interpretable Neural Networks with Frank-Wolfe: Sparse Relevance Maps and Relevance Orderings", Macdonald et al. [9]

| Abstract Games: complex, yet low-dimensional

1) "Proximal policy optimization algorithms", Schulman et al [1] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Setup: Connect Four with hidden colour information

I Abstract Games: complex, yet low-dimensional

1) "Proximal policy optimization algorithms", Schulman et al [11] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Setup: Connect Four with hidden colour information

- I Abstract Games: complex, yet low-dimensional
- I Every turnt: $p_h \cup ([0; p_h^{\max}])$

1) "Proximal policy optimization algorithms", Schulman et al [11] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Setup: Connect Four with hidden colour information

- I Abstract Games: complex, yet low-dimensional
- I Every turn: $p_h \cup ([0; p_h^{\max}])$
- I Hide p_h 's colour features at random

1) "Proximal policy optimization algorithms", Schulman et al [11] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Setup: Connect Four with hidden colour information

- I Abstract Games: complex, yet low-dimensional
- I Every turn: $p_h \cup ([0; p_h^{\max}])$
- I Hide p_h 's colour features at random

1) "Proximal policy optimization algorithms", Schulman et al [11] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Setup: Connect Four with hidden colour information

- I Abstract Games: complex, yet low-dimensional
- I Every turn: $p_h \cup ([0; p_h^{\max}])$
- I Hide p_h 's colour features at random
- I Train agent with Proximal Policy Optimisation (PPO)^{1,2}

1) "Proximal policy optimization algorithms", Schulman et al [11] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Setup: Connect Four with hidden colour information

- I Abstract Games: complex, yet low-dimensional
- I Every turn: $p_h \cup ([0; p_h^{\max}])$
- I Hide p_h 's colour features at random
- I Train agent with Proximal Policy Optimisation (PPO)^{1,2}

FI: Full Information.

1) "Proximal policy optimization algorithms", Schulman et al [11] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Setup: Connect Four with hidden colour information

- I Abstract Games: complex, yet low-dimensional
- I Every turn: $p_h \sim \mathcal{U}([0; p_h^{\max}])$
- I Hide p_h 's colour features at random
- I Train agent with Proximal Policy Optimisation (PPO)^{1,2}

FI: Full Information.

PI-50: with $p_h \sim \mathcal{U}([0; 0.5])$

1) "Proximal policy optimization algorithms", Schulman et al [11] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Setup: Connect Four with hidden colour information

- I Abstract Games: complex, yet low-dimensional
- I Every turn: $p_h \sim \mathcal{U}([0; p_h^{\max}])$
- I Hide p_h 's colour features at random
- I Train agent with Proximal Policy Optimisation (PPO)^{1,2}

FI: Full Information.

PI-50: with $p_h \sim \mathcal{U}([0; 0.5])$

PI-100: with $p_h \sim \mathcal{U}([0; 1])$

1) "Proximal policy optimization algorithms", Schulman et al [11] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Setup: Connect Four with hidden colour information

- I Abstract Games: complex, yet low-dimensional
- I Every turn: $p_h \sim \mathcal{U}([0; p_h^{\max}])$
- I Hide p_h 's colour features at random
- I Train agent with Proximal Policy Optimisation (PPO)^{1,2}

FI: Full Information.

PI-50: with $p_h \sim \mathcal{U}([0; 0.5])$

PI-100: with $p_h \sim \mathcal{U}([0; 1])$

1) "Proximal policy optimization algorithms", Schulman et al [11] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Setup: Connect Four with hidden colour information

- | Abstract Games: complex, yet low-dimensional
- | Every turn: $p_h \sim \mathcal{U}([0; p_h^{\max}])$
- | Hide p_h 's colour features at random
- | Train agent with Proximal Policy Optimisation (PPO)^{1,2}

FI: Full Information.

PI-50: with $p_h \sim \mathcal{U}([0; 0.5])$

PI-100: with $p_h \sim \mathcal{U}([0; 1])$

1) "Proximal policy optimization algorithms", Schulman et al [11] 2) "Reinforcement Learning for Two-Player Zero-Sum Games", Crespo [2]

Interpretability with Characteristic Functions

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Interpretability with Characteristic Functions

I Let $t \geq 2$ [42]

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Interpretability with Characteristic Functions

1) Let $t \in [42]$, $x \in [0; 1]^3$ ^{6 7}

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Interpretability with Characteristic Functions

1) Let $t \in [42]$, $x \in [0; 1]^3$, $S \in [t]$

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Interpretability with Characteristic Functions

1) Let $t \in [0, 1]$, $x \in [0; 1]^3$, $S \in [t]$ and let $x^{(S)}$ be state with colour feature on S^c hidden

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Interpretability with Characteristic Functions

- | Let $t \in [0, 1]$, $x \in [0, 1]^3$, $S \in [t]$ and let $x^{(S)}$ be state with colour feature on S^c hidden
- | Let furthermore $a = \operatorname{argmax}_a P(a; x)$

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Interpretability with Characteristic Functions

- | Let $t \in [2, 7]$, $x \in [0; 1]^3$, $S \in [t]$ and let $x^{(S)}$ be state with colour feature on S^c hidden
- | Let furthermore $a = \operatorname{argmax}_a P(a; x)$
- | We define $\text{pol} : 2^{[t]} \rightarrow [0; 1]$ and $\text{val} : 2^{[t]} \rightarrow [0; 1]$ as

$$\text{pol}(S) = P(a; x^{(S)}) \quad \text{and} \quad \text{val}(S) = V(x^{(S)}):$$

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Interpretability with Characteristic Functions

- | Let $t \in [2, 7]$, $x \in [0; 1]^3$, $S \in \{t\}$ and let $x^{(S)}$ be state with colour feature on S^c hidden
- | Let furthermore $a = \operatorname{argmax}_a P(a; x)$
- | We define $\text{pol} : 2^{[t]} \rightarrow [0; 1]$ and $\text{val} : 2^{[t]} \rightarrow [0; 1]$ as

$$\text{pol}(S) = P(a; x^{(S)}) \quad \text{and} \quad \text{val}(S) = V(x^{(S)}):$$

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Interpretability with Characteristic Functions

- Let $t \in [2, 6, 7]$, $x \in [0; 1]^3$, $S \subseteq [t]$ and let $x^{(S)}$ be state with colour feature on S^c hidden
- Let furthermore $a = \operatorname{argmax}_a P(a; x)$
- We define $\text{pol} : 2^{[t]} \rightarrow [0; 1]$ and $\text{val} : 2^{[t]} \rightarrow [0; 1]$ as

$$\text{pol}(S) = P(a; x^{(S)}) \quad \text{and} \quad \text{val}(S) = V(x^{(S)}):$$

- We can approximate the Shapley sum by sampling from $\mathcal{B}([t])$:

$$\phi_i = \frac{1}{t!} \sum_{S \in \mathcal{B}([t])} (P_i(f|_S) - P_i):$$

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Interpretability with Characteristic Functions

- Let $t \in [2, 6, 7]$, $x \in [0; 1]^3$, $S \subseteq [t]$ and let $x^{(S)}$ be state with colour feature on S^c hidden
- Let furthermore $a = \operatorname{argmax}_a P(a; x)$
- We define $\text{pol} : 2^{[t]} \rightarrow [0; 1]$ and $\text{val} : 2^{[t]} \rightarrow [0; 1]$ as

$$\text{pol}(S) = P(a; x^{(S)}) \quad \text{and} \quad \text{val}(S) = V(x^{(S)}):$$

- We can approximate the Shapley sum by sampling from $\mathcal{M}([t])$:

$$\phi_i = \frac{1}{t!} \sum_{S \in \mathcal{M}([t])} (P_i(f|_S) - P_i):$$

- P_i is a $(0;01; 0;01)$ -approximation 26 500 samples (Hoeffding)

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Interpretability with Characteristic Functions

- Let $t \in [2, 6, 7]$, $x \in [0; 1]^3$, $S \subseteq [t]$ and let $x^{(S)}$ be state with colour feature on S^c hidden
- Let furthermore $a = \operatorname{argmax}_a P(a; x)$
- We define $\text{pol} : 2^{[t]} \rightarrow [0; 1]$ and $\text{val} : 2^{[t]} \rightarrow [0; 1]$ as

$$\text{pol}(S) = P(a; x^{(S)}) \quad \text{and} \quad \text{val}(S) = V(x^{(S)}):$$

- We can approximate the Shapley sum by sampling from $\binom{[t]}{t}$:

$$\phi_i = \frac{1}{t!} \sum_{S \subseteq [t], i \in S} (P_i([f]_S) - P_i(S \setminus \{i\}))$$

- P_i is $(0; 0.01; 0; 0.01)$ -approximation 26 500 samples (Hoeffding)
- Calculate PIE with Frank-Wolfe optimiser solving convex relaxation of

$$S = \operatorname{argmin}_{|S| \leq \rho n t c} \sum_{i \in S} \phi_i^2:$$

1) "Deep Neural Network Training with Frank-Wolfe", Pokutta et al. [10]

Example Saliencies for different Methods

Information-Performance Comparison

Round-Robin Tournament

Round-Robin Tournament

Round-Robin Tournament

Limitations and Outlook

Limitations and Outlook

- I So far works only for certain abstract games (Connect Four, Hex, Go, ...)

Limitations and Outlook

- | So far works only for certain abstract games (Connect Four, Hex, Go, ...)
-) Filter out illegal moves with model-based approaches (see e.g. AlphaGo)

Limitations and Outlook

- | So far works only for certain abstract games (Connect Four, Hex, Go, ...)
-) Filter out illegal moves with model-based approaches (see e.g. AlphaGo)
- | Further extension to real-world tasks (high-dimensional, high redundancy) is challenging

Limitations and Outlook

- | So far works only for certain abstract games (Connect Four, Hex, Go, ...)
-) Filter out illegal moves with model-based approaches (see e.g. AlphaGo)
- | Further extension to real-world tasks (high-dimensional, high redundancy) is challenging
-) Our approach is should be used primarily for evaluation of saliency methods

Limitations and Outlook

- | So far works only for certain abstract games (Connect Four, Hex, Go, ...)
-) Filter out illegal moves with model-based approaches (see e.g. AlphaGo)
- | Further extension to real-world tasks (high-dimensional, high redundancy) is challenging
-) Our approach is should be used primarily for evaluation of saliency methods
- | Shapley sampling suffered from unstable policy layer for large hidden information

Limitations and Outlook

- | So far works only for certain abstract games (Connect Four, Hex, Go, ...)
-) Filter out illegal moves with model-based approaches (see e.g. AlphaGo)
- | Further extension to real-world tasks (high-dimensional, high redundancy) is challenging
-) Our approach is should be used primarily for evaluation of saliency methods
- | Shapley sampling suffered from unstable policy layer for large hidden information
-) Train value function instead

Limitations and Outlook

- | So far works only for certain abstract games (Connect Four, Hex, Go, ...)
-) Filter out illegal moves with model-based approaches (see e.g. AlphaGo)
- | Further extension to real-world tasks (high-dimensional, high redundancy) is challenging
-) Our approach is should be used primarily for evaluation of saliency methods
- | Shapley sampling suffered from unstable policy layer for large hidden information
-) Train value function instead
-) Q-Learning could be a more stable approach

Conclusion:

Conclusion:

- | Interpretability relies on a good model of the data distribution

Conclusion:

- | Interpretability relies on a good model of the data distribution
- | We can design proxy-task where we know the distribution via abstract games with missing information

Conclusion:

- | Interpretability relies on a good model of the data distribution
- | We can design proxy-task where we know the distribution via abstract games with missing information
- | Use these tasks to evaluate saliency methods without going o -manifold

Conclusion:

- | Interpretability relies on a good model of the data distribution
- | We can design proxy-task where we know the distribution via abstract games with missing information
- | Use these tasks to evaluate saliency methods without going o -manifold

Thank You!

Conclusion:

- | Interpretability relies on a good model of the data distribution
- | We can design proxy-task where we know the distribution via abstract games with missing information
- | Use these tasks to evaluate saliency methods without going o -manifold

Thank You!

Contact: waeldchen@zib.de

Paper: Training Characteristic Functions with Reinforcement Learning:
XAI-methods play Connect Four
Stephan Waldchen, F Huber, S Pokutta
arXiv preprint [arXiv:2202.11797](https://arxiv.org/abs/2202.11797)

C. Anders, P. Pasliev, A.-K. Dombrowski, K.-R. Müller, and P. Kessel.

Fairwashing explanations with o-manifold detergent.

In International Conference on Machine Learning, pages 314–323. PMLR, 2020.

J. Crespo.

Reinforcement learning for two-player zero-sum games.

Master's thesis, Técnico Lisboa, https://fenix.tecnico.ulisboa.pt/downloadFile/1689244997260153/81811-joao-crespo_dissertacao.pdf, 2019.

B. Dimanov, U. Bhatt, M. Jamnik, and A. Weller.

You shouldn't trust me: Learning models which conceal unfairness from multiple explanation methods.

In SafeAI@ AAAI2020.

A.-K. Dombrowski, M. Alber, C. J. Anders, M. Ackermann, K.-R. Müller, and P. Kessel.
Explanations can be manipulated and geometry is to blame.

arXiv preprint arXiv:1906.07983, 2019.

C. Frye, D. de Mijolla, T. Begley, L. Cowton, M. Stanley, and I. Feige.
Shapley explainability on the data manifold.

arXiv preprint arXiv:2006.01272, 2020.

J. Heo, S. Joo, and T. Moon.

Fooling neural network interpretations via adversarial model manipulation.

Advances in Neural Information Processing Systems, 32:2925–2936, 2019.

P. Linardatos, V. Papastefanopoulos, and S. Kotsiantis.

Explainable ai: A review of machine learning interpretability methods.

Entropy, 23(1):18, 2020.

S. M. Lundberg and S.-I. Lee.

A uni ed approach to interpreting model predictions.

In Proceedings of the 31st international conference on neural information processing systems pages 4768-4777, 2017.

J. Macdonald, M. Besarcon, and S. Pokutta.

Interpretable neural networks with frank-wolfe: Sparse relevance maps and relevance orderings.

arXiv preprint arXiv:2110.08105, 2021.

S. Pokutta, C. Spiegel, and M. Zimmer.

Deep neural network training with frank-wolfe.

arXiv preprint arXiv:2010.07243, 2020.

J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov.

Proximal policy optimization algorithms.

arXiv preprint arXiv:1707.06347, 2017.

L. S. Shapley.

17. A value for n-person games.

Princeton University Press, 2016.

A. Shih, A. Choi, and A. Darwiche.

A symbolic approach to explaining bayesian network classifiers.

arXiv preprint arXiv:1805.03364, 2018.

D. Slack, S. Hilgard, E. Jia, S. Singh, and H. Lakkaraju.

Fooling lime and shap: Adversarial attacks on post hoc explanation methods.

In Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society, pages 180–186, 2020.

D. Slack, S. Hilgard, H. Lakkaraju, and S. Singh.

Counterfactual explanations can be manipulated.

arXiv preprint arXiv:2106.02666, 2021.

Ground Truth Comparison: Winning Move

Tournament: Standard Deviation and Illegal Move Rate

