# Training Characteristic Functions with Reinforcement Learning

**Stephan Wäldchen, Felix Huber, Sebastian Pokutta**

Zuse Institut Berlin

▶ Core idea: Probe black-box function with different inputs         $\mathbf{x} \longrightarrow \boxed{\Phi} \longrightarrow y$

## How do we interpret black-box Functions?

- ▶ Core idea: Probe black-box function with different inputs
- ▶ Local interpretation (for specific input **x**): vary features

$$\mathbf{x} \longrightarrow \boxed{\Phi} \longrightarrow y$$

# 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



LIME saliency[1]

---

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 $\leftrightarrow$ Coop. Game Theory: Surplus Attribution



LIME saliency[1]

---

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

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



LIME saliency[1]

---

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: $\nu \colon 2^{[d]} \to \mathbb{R}$



LIME saliency[1]

- ▶ 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: $\nu\colon 2^{[d]} \to \mathbb{R}$
- ▶ Prime Implicant Explanation[2] (mostly for $\nu : 2^{[d]} \to \{0,1\}$)



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

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

---

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]

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

▶ Local interpretation (for specific input **x**): vary features

▶ ML: Saliency $\leftrightarrow$ Coop. Game Theory: Surplus Attribution

▶ Characteristic function: $\nu \colon 2^{[d]} \to \mathbb{R}$

▶ Prime Implicant Explanation[2] (mostly for $\nu : 2^{[d]} \to \{0, 1\}$)

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

▶ Shapley Values[3] (linear, efficient, symmetric, null-player)

$$\phi_i = \frac{1}{d!} \sum_{\pi \in \Pi([d])} \left( \nu(P_i^\pi \cup \{i\}) - \nu(P_i^\pi) \right).$$



$0   $7   $4   $6
$7   $15   $9   $19

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

---

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 **x**): vary features

▶ ML: Saliency ↔ Coop. Game Theory: Surplus Attribution

▶ Characteristic function: $\nu\colon 2^{[d]} \to \mathbb{R}$

▶ Prime Implicant Explanation[2] (mostly for $\nu : 2^{[d]} \to \{0,1\}$)


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

$$S^* = \operatorname*{argmin}_{S \subset [d]} |S| \quad \text{s.t.} \quad \nu(S) = \nu([d])$$

▶ Shapley Values[3] (linear, efficient, symmetric, null-player)

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

▶ 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 al[1]:

$$\nu_{\Phi,\mathbf{x}}(S) = \mathbb{E}_{\mathbf{y}}[\Phi(\mathbf{y}) \,|\, \mathbf{y}_S = \mathbf{x}_S] = \int \Phi(\mathbf{x}) d\mathbb{P}[\mathbf{x}_{S^c} | \mathbf{x}_S].$$

---

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

- Approach from Lundberg et al[1]:

$$\nu_{\Phi,\mathbf{x}}(S) = \mathbb{E}_{\mathbf{y}}[\Phi(\mathbf{y}) \,|\, \mathbf{y}_S = \mathbf{x}_S] = \int \Phi(\mathbf{x}) d\mathbb{P}[\mathbf{x}_{S^c} | \mathbf{x}_S].$$

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

---

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

## Interpretations Rely on a Model of the Data Manifold

▶ Approach from Lundberg et al[1]:

$$\nu_{\Phi,\mathbf{x}}(S) = \mathbb{E}_{\mathbf{y}}[\Phi(\mathbf{y}) \,|\, \mathbf{y}_S = \mathbf{x}_S] = \int \Phi(\mathbf{x}) d\mathbb{P}[\mathbf{x}_{S^c}|\mathbf{x}_S].$$

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

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

---

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

## Interpretations Rely on a Model of the Data Manifold

▶ Approach from Lundberg et al[1]:

$$\nu_{\Phi,\mathbf{x}}(S) = \mathbb{E}_{\mathbf{y}}[\Phi(\mathbf{y}) \,|\, \mathbf{y}_S = \mathbf{x}_S] = \int \Phi(\mathbf{x}) \mathrm{d}\mathbb{P}[\mathbf{x}_{S^c} | \mathbf{x}_S].$$

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

▶ 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-Cam[7],
Shapley-based[6], Counterfactual explanations[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]

# Interpretations Rely on a Model of the Data Manifold

► Approach from Lundberg et al[1]:

$$\nu_{\Phi,\mathbf{x}}(S) = \mathbb{E}_{\mathbf{y}}[\Phi(\mathbf{y}) \,|\, \mathbf{y}_S = \mathbf{x}_S] = \int \Phi(\mathbf{x}) \mathrm{d}\mathbb{P}[\mathbf{x}_{S^c} | \mathbf{x}_S].$$

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

► 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-Cam[7], Shapley-based[6], Counterfactual explanations[8],



Image    FW    AFW

LCG    LAFW    Sensitivity

► 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]

# Interpretations Rely on a Model of the Data Manifold

▶ Approach from Lundberg et al[1]:

$$\nu_{\Phi, \mathbf{x}}(S) = \mathbb{E}_{\mathbf{y}}[\Phi(\mathbf{y}) \mid \mathbf{y}_S = \mathbf{x}_S] = \int \Phi(\mathbf{x}) d\mathbb{P}[\mathbf{x}_{S^c} | \mathbf{x}_S].$$

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

▶ 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-Cam[7], Shapley-based[6], Counterfactual explanations[8],



Image    FW    AFW

LCG    LAFW    Sensitivity

▶ 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

---

▶ Abstract Games: complex, yet low-dimensional



Game Board     Channel1     Channel2     Channel3

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

- Abstract Games: complex, yet low-dimensional
- Every turn $t$: $p_h \sim \mathcal{U}([0, p_h^{\max}])$



Game Board    Channel1    Channel2    Channel3

---

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

- Abstract Games: complex, yet low-dimensional
- Every turn $t$: $p_h \sim \mathcal{U}([0, p_h^{\max}])$
- Hide $\lfloor p_h t \rfloor$ colour features at random



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

ICML 2022 (Stephan Wäldchen)

3/15

- Abstract Games: complex, yet low-dimensional
- Every turn $t$: $p_h \sim \mathcal{U}([0, p_h^{\max}])$
- Hide $\lfloor p_h t \rfloor$ colour features at random

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

ICML 2022 (Stephan Wäldchen)                                                                 3/15

- Abstract Games: complex, yet low-dimensional
- Every turn $t$: $p_h \sim \mathcal{U}([0, p_h^{\max}])$
- Hide $\lfloor p_h t \rfloor$ colour features at random
- 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]

- ▶ Abstract Games: complex, yet low-dimensional
- ▶ Every turn $t$: $p_h \sim \mathcal{U}([0, p_h^{\max}])$
- ▶ Hide $\lfloor p_h t \rfloor$ colour features at random
- ▶ 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

▶ Abstract Games: complex, yet low-dimensional

▶ Every turn $t$: $p_h \sim \mathcal{U}([0, p_h^{\max}])$

▶ Hide $\lfloor p_h t \rfloor$ colour features at random

▶ Train agent with Proximal Policy Optimisation (PPO)[1,2]

◆ FI: Full Information.

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

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 $t$: $p_h \sim \mathcal{U}([0, p_h^{\max}])$
- Hide $\lfloor p_h t \rfloor$ 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]

# Setup: Connect Four with hidden colour information

- ▶ Abstract Games: complex, yet low-dimensional
- ▶ Every turn $t$: $p_h \sim \mathcal{U}([0, p_h^{\max}])$
- ▶ Hide $\lfloor p_h t \rfloor$ 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]

# Setup: Connect Four with hidden colour information

- ▶ Abstract Games: complex, yet low-dimensional
- ▶ Every turn $t$: $p_h \sim \mathcal{U}([0, p_h^{\max}])$
- ▶ Hide $\lfloor p_h t \rfloor$ 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]

---

► Let $t \in [42]$

---

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

▶ Let $t \in [42]$, $\mathbf{x} \in [0, 1]^{3 \times 6 \times 7}$

---

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

▶ Let $t \in [42]$, $\mathbf{x} \in [0,1]^{3 \times 6 \times 7}$, $S \in [t]$

---

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

▶ Let $t \in [42]$, $\mathbf{x} \in [0,1]^{3 \times 6 \times 7}$, $S \in [t]$ and let $\mathbf{x}^{(S)}$ be state with colour feature on $S^c$ hidden

---

► Let $t \in [42]$, $\mathbf{x} \in [0,1]^{3 \times 6 \times 7}$, $S \in [t]$ and let $\mathbf{x}^{(S)}$ be state with colour feature on $S^c$ hidden

► Let furthermore $a^* = \text{argmax}_a P(a\,;\mathbf{x})$

---

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

## Interpretability with Characteristic Functions

- Let $t \in [42]$, $\mathbf{x} \in [0,1]^{3 \times 6 \times 7}$, $S \in [t]$ and let $\mathbf{x}^{(S)}$ be state with colour feature on $S^c$ hidden
- Let furthermore $a^* = \text{argmax}_a P(a\,;\mathbf{x})$
- We define $\nu^{\text{pol}} : 2^{[t]} \to [0,1]$ and $\nu^{\text{val}} : 2^{[t]} \to [-1,1]$ as

$$\nu^{\text{pol}}(S) = P(a^*;\mathbf{x}^{(S)}) \quad \text{and} \quad \nu^{\text{val}}(S) = V(\mathbf{x}^{(S)}).$$

---

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

- Let $t \in [42]$, $\mathbf{x} \in [0,1]^{3 \times 6 \times 7}$, $S \in [t]$ and let $\mathbf{x}^{(S)}$ be state with colour feature on $S^c$ hidden
- Let furthermore $a^* = \text{argmax}_a P(a\,;\mathbf{x})$
- We define $\nu^{\text{pol}} : 2^{[t]} \to [0,1]$ and $\nu^{\text{val}} : 2^{[t]} \to [-1,1]$ as

$$\nu^{\text{pol}}(S) = P(a^*;\mathbf{x}^{(S)}) \quad \text{and} \quad \nu^{\text{val}}(S) = V(\mathbf{x}^{(S)}).$$

---

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

# Interpretability with Characteristic Functions

- Let $t \in [42]$, $\mathbf{x} \in [0,1]^{3 \times 6 \times 7}$, $S \in [t]$ and let $\mathbf{x}^{(S)}$ be state with colour feature on $S^c$ hidden
- Let furthermore $a^* = \operatorname{argmax}_a P(a \, ; \mathbf{x})$
- We define $\nu^{\text{pol}} : 2^{[t]} \to [0,1]$ and $\nu^{\text{val}} : 2^{[t]} \to [-1,1]$ as

$$\nu^{\text{pol}}(S) = P(a^*; \mathbf{x}^{(S)}) \quad \text{and} \quad \nu^{\text{val}}(S) = V(\mathbf{x}^{(S)}).$$

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

$$\phi_i = \frac{1}{t!} \sum_{\pi \in \Pi([t])^{\text{sample}}} (\nu(P_i^\pi \cup \{i\}) - \nu(P_i^\pi)).$$

---

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

► Let $t \in [42]$, $\mathbf{x} \in [0,1]^{3 \times 6 \times 7}$, $S \in [t]$ and let $\mathbf{x}^{(S)}$ be state with colour feature on $S^c$ hidden

► Let furthermore $a^* = \text{argmax}_a P(a \, ; \mathbf{x})$

► We define $\nu^{\text{pol}} : 2^{[t]} \to [0,1]$ and $\nu^{\text{val}} : 2^{[t]} \to [-1,1]$ as

$$\nu^{\text{pol}}(S) = P(a^* ; \mathbf{x}^{(S)}) \quad \text{and} \quad \nu^{\text{val}}(S) = V(\mathbf{x}^{(S)}).$$

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

$$\phi_i = \frac{1}{t!} \sum_{\pi \in \Pi([t])^{\text{sample}}} (\nu(P_i^\pi \cup \{i\}) - \nu(P_i^\pi)).$$

► $\mathbb{P}\big[\big|\phi_i - \bar{\phi}_i\big| \leq \epsilon\big] \geq 1 - \delta \to (0.01, 0.01)$-approximation $\approx$26 500 samples (Hoeffding)

---

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

# Interpretability with Characteristic Functions

▶ Let $t \in [42]$, $\mathbf{x} \in [0,1]^{3 \times 6 \times 7}$, $S \in [t]$ and let $\mathbf{x}^{(S)}$ be state with colour feature on $S^c$ hidden

▶ Let furthermore $a^* = \text{argmax}_a P(a ; \mathbf{x})$

▶ We define $\nu^{\text{pol}} : 2^{[t]} \to [0,1]$ and $\nu^{\text{val}} : 2^{[t]} \to [-1,1]$ as

$$\nu^{\text{pol}}(S) = P(a^* ; \mathbf{x}^{(S)}) \quad \text{and} \quad \nu^{\text{val}}(S) = V(\mathbf{x}^{(S)}).$$

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

$$\phi_i = \frac{1}{t!} \sum_{\pi \in \Pi([t])^{\text{sample}}} (\nu(P_i^\pi \cup \{i\}) - \nu(P_i^\pi)).$$

▶ $\mathbb{P}\big[\big|\phi_i - \bar{\phi}_i\big| \leq \epsilon\big] \geq 1 - \delta \to (0.01, 0.01)$-approximation $\approx 26\ 500$ samples (Hoeffding)

▶ Calculate PIE with Frank-Wolfe optimiser solving[1] convex relaxation of

$$S^* = \underset{|S| \leq \lfloor p_h t \rfloor}{\text{argmin}} \, (\nu([t]) - \nu(S))^2.$$

---

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

Gradient    DeepShap    GuidedBP    SmoothGrad    LRP

DeepTaylor    Random    **Shapley Sampling**    **FW**

Masker1

Player1

Player2

Masker2

▶ 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

ZIB

- ▶ 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

ZIB

- ▶ 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 off-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 off-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 off-manifold

## Thank You!

🌐 **Contact:** waeldchen@zib.de

📄 **Paper:** *Training Characteristic Functions with Reinforcement Learning:
  XAI-methods play Connect Four*, S Wäldchen, F Huber, S Pokutta
  *arXiv preprint* arXiv:2202.11797

📄 C. Anders, P. Pasliev, A.-K. Dombrowski, K.-R. Müller, and P. Kessel.
**Fairwashing explanations with off-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, Tecnico 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@ AAAI*, 2020.

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 unified approach to interpreting model predictions.**
In *Proceedings of the 31st international conference on neural information processing systems*, pages 4768–4777, 2017.

📄 J. Macdonald, M. Besançon, 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.

# Appendix