

Practical Approaches to Explainable Machine Learning

Umang Bhatt

PhD Candidate, University of Cambridge

Fellow, Mozilla Foundation

Student Fellow, Leverhulme Center for the Future of Intelligence

@umangsbhatt

usb20@cam.ac.uk



Our Starting Point

How are existing approaches to explainability used in practice?

Explainable Machine Learning in Deployment

Umang Bhatt, Alice Xiang, Shubham Sharma, Adrian Weller, Ankur Taly,
Yuhuan Jia, Joydeep Ghosh, Ruchir Puri, José Moura, and Peter Eckersley

Appeared at the ACM Conference on Fairness, Accountability, and Transparency 2020

<https://arxiv.org/abs/1909.06342>



Growth of Transparency Literature

Many algorithms proposed to “explain” machine learning model output

We study how organizations use these algorithms, if at all

Our Approach

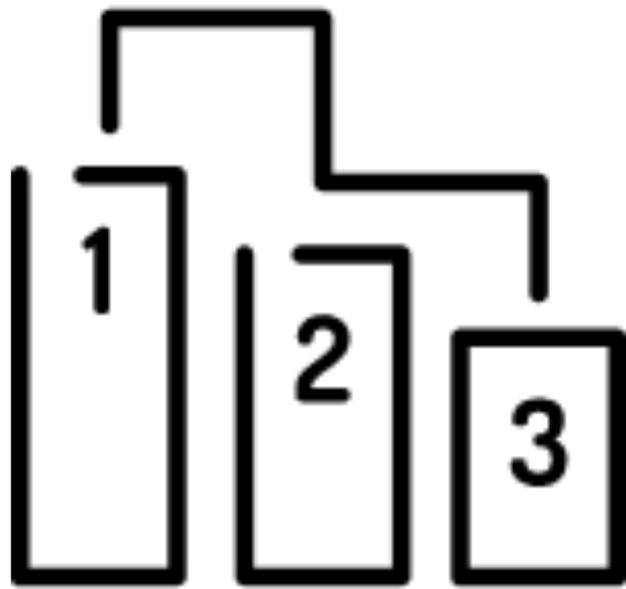
30 minute to 2 hour semi-structured interviews

50 individuals from 30 organizations interviewed

Shared Language

- **Transparency:** Providing stakeholders with relevant information about how the model works: this includes documentation of the training procedure, analysis of training data distribution, code releases, feature-level explanations, etc.
- **Explainability:** Providing insights into a model's behavior for specific datapoint(s)

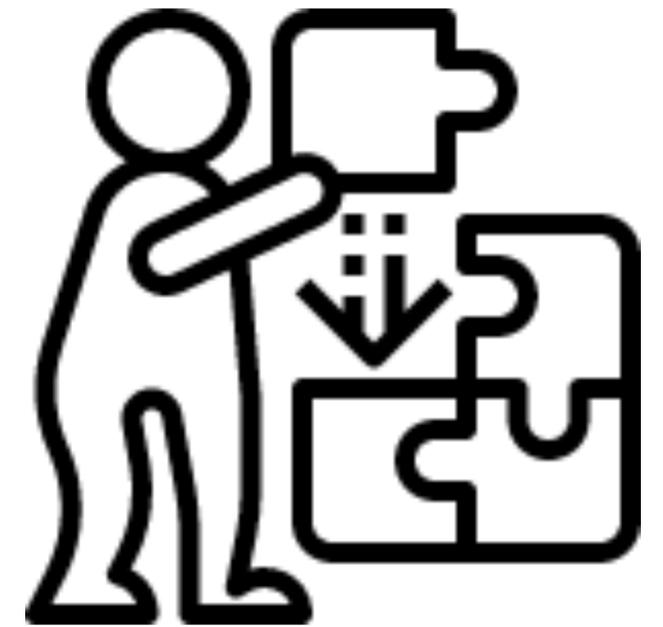
Types of Explanations



Feature Importance



Sample Importance

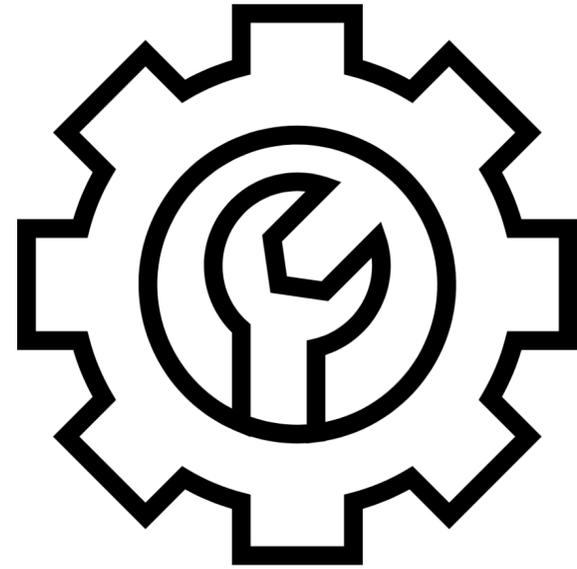


Counterfactuals

Stakeholders



Executives



Engineers



End Users



Regulators

Findings

1. Explainability is used for **debugging** internally
2. **Goals** of explainability are not clearly defined within organizations
3. Technical **limitations** make explainability hard to deploy in real-time

Machine Learning Explainability for External Stakeholders

Umang Bhatt, McKane Andrus, Adrian Weller, and Alice Xiang

Convening hosted by CFI, PAI, and IBM

Appeared at the ICML 2020 Workshop on Extending Explainable AI: Beyond Deep Models and Classifiers

<https://arxiv.org/abs/2007.05408>



Overview

- 33 participants from 5 countries
 - 15 ML experts, 3 designers, 6 legal experts, 9 policymakers
 - Domain expertise: Finance, Healthcare, Media, and Social Services
- Goal: ***facilitate an inter-stakeholder conversation around explainable machine learning***

Community Engagement

1. *In which **context** will this explanation be used? Does the context change the properties of the explanations we expose?*
2. *How should the explanation be **evaluated**? Both quantitatively and qualitatively...*
3. *Can we prevent data misuse and preferential treatment by involving **affected groups** in the development process?*
4. *Can we **educate** external stakeholders (and data scientists) regarding the functionalities and limitations of explainable machine learning?*

Deploying Explainability

1. How does **uncertainty** in the model and introduced by the (potentially approximate) explanation technique affect the resulting explanations?
2. How can stakeholders **interact** with the resulting explanations? Can explanations be a conduit for interacting with the model?
3. How, if at all, will stakeholder **behavior** change as a result of the explanation shown?
4. Over **time**, how will the explanation technique adapt to changes in stakeholder behavior?

Practitioner-Driven Questions

1. Can we provide methodology for practitioners to **evaluate** explanations?
2. Can existing explainability tools be used to identify model **unfairness**?
3. Can we quantify how much **uncertainty** is associated with given explanations?

Practitioner-Driven Questions

1. Can we provide methodology for practitioners to **evaluate** explanations?
2. Can existing explainability tools be used to identify model unfairness?
3. Can we quantify how much uncertainty is associated with given explanations?

Evaluating and Aggregating Feature-based Model Explanations

Umang Bhatt, Adrian Weller, and José Moura

Appeared at the International Joint Conference on Artificial Intelligence 2020

<https://arxiv.org/abs/2005.00631>



Carnegie
Mellon
University

The
Alan Turing
Institute

Feature Importance Desiderata

Explain $\mathbf{f}(\mathbf{x}) : \mathbb{R}^d \rightarrow \{C_1 \dots C_j\}$

Feature Importance Desiderata

Explain $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \{C_1 \dots C_j\}$

- Restrict ourselves to feature-based explanations in the classification setting (i.e. attributions, saliency maps, etc.)

Feature Importance Desiderata

Explain $\mathbf{f}(\mathbf{x}) : \mathbb{R}^d \rightarrow \{C_1 \dots C_j\}$

- Restrict ourselves to feature-based explanations in the classification setting (i.e. attributions, saliency maps, etc.)
 - Which feature(s) x_i contributes to classifying \mathbf{x} with \mathbf{f}

Feature Importance Desiderata

Explain $\mathbf{f}(\mathbf{x}) : \mathbb{R}^d \rightarrow \{C_1 \dots C_j\}$

- Restrict ourselves to feature-based explanations in the classification setting (i.e. attributions, saliency maps, etc.)
 - Which feature(s) x_i contributes to classifying \mathbf{x} with \mathbf{f}
 - Let \mathbf{g} be our explanation technique and let $\mathbf{g}(\mathbf{f}, \mathbf{x})$ be the explanation for $\mathbf{f}(\mathbf{x})$

Feature Importance Desiderata

Explain $\mathbf{f}(\mathbf{x}) : \mathbb{R}^d \rightarrow \{C_1 \dots C_j\}$

- Restrict ourselves to feature-based explanations in the classification setting (i.e. attributions, saliency maps, etc.)
 - Which feature(s) x_i contributes to classifying \mathbf{x} with \mathbf{f}
 - Let \mathbf{g} be our explanation technique and let $\mathbf{g}(\mathbf{f}, \mathbf{x})$ be the explanation for $\mathbf{f}(\mathbf{x})$
- Feature Importance: $\mathbf{g}_m = |\mathbf{g}(\mathbf{f}, \mathbf{x})| \in \mathbb{R}_{\geq 0}^d$

Feature Importance Desiderata

Explain $\mathbf{f}(\mathbf{x}) : \mathbb{R}^d \rightarrow \{C_1 \dots C_j\}$

- Restrict ourselves to feature-based explanations in the classification setting (i.e. attributions, saliency maps, etc.)
 - Which feature(s) x_i contributes to classifying \mathbf{x} with \mathbf{f}
 - Let \mathbf{g} be our explanation technique and let $\mathbf{g}(\mathbf{f}, \mathbf{x})$ be the explanation for $\mathbf{f}(\mathbf{x})$
- Feature Importance: $\mathbf{g}_m = |\mathbf{g}(\mathbf{f}, \mathbf{x})| \in \mathbb{R}_{\geq 0}^d$
 - Normalize \mathbf{g}_m to induce a probability distribution over the features in the input space

Feature Importance Desiderata

Explain $\mathbf{f}(\mathbf{x}) : \mathbb{R}^d \rightarrow \{C_1 \dots C_j\}$

- Restrict ourselves to feature-based explanations in the classification setting (i.e. attributions, saliency maps, etc.)
 - Which feature(s) x_i contributes to classifying \mathbf{x} with \mathbf{f}
 - Let \mathbf{g} be our explanation technique and let $\mathbf{g}(\mathbf{f}, \mathbf{x})$ be the explanation for $\mathbf{f}(\mathbf{x})$
- Feature Importance: $\mathbf{g}_m = |\mathbf{g}(\mathbf{f}, \mathbf{x})| \in \mathbb{R}_{\geq 0}^d$
 - Normalize \mathbf{g}_m to induce a probability distribution over the features in the input space
- Feature Contribution: $\mathbf{g}_c = \mathbf{g}(\mathbf{f}, \mathbf{x}) \in \mathbb{R}^d$ and $|\mathbf{g}_c| = \mathbf{g}_m$

Feature Importance Desiderata

Explain $\mathbf{f}(\mathbf{x}) : \mathbb{R}^d \rightarrow \{C_1 \dots C_j\}$

- Restrict ourselves to feature-based explanations in the classification setting (i.e. attributions, saliency maps, etc.)
 - Which feature(s) x_i contributes to classifying \mathbf{x} with \mathbf{f}
 - Let \mathbf{g} be our explanation technique and let $\mathbf{g}(\mathbf{f}, \mathbf{x})$ be the explanation for $\mathbf{f}(\mathbf{x})$
- Feature Importance: $\mathbf{g}_m = |\mathbf{g}(\mathbf{f}, \mathbf{x})| \in \mathbb{R}_{\geq 0}^d$
 - Normalize \mathbf{g}_m to induce a probability distribution over the features in the input space
- Feature Contribution: $\mathbf{g}_c = \mathbf{g}(\mathbf{f}, \mathbf{x}) \in \mathbb{R}^d$ and $|\mathbf{g}_c| = \mathbf{g}_m$
 - Presence of rain increases the chance of precipitation (+)

Feature Importance Desiderata

Explain $\mathbf{f}(\mathbf{x}) : \mathbb{R}^d \rightarrow \{C_1 \dots C_j\}$

- Restrict ourselves to feature-based explanations in the classification setting (i.e. attributions, saliency maps, etc.)
 - Which feature(s) x_i contributes to classifying \mathbf{x} with \mathbf{f}
 - Let \mathbf{g} be our explanation technique and let $\mathbf{g}(\mathbf{f}, \mathbf{x})$ be the explanation for $\mathbf{f}(\mathbf{x})$
- Feature Importance: $\mathbf{g}_m = |\mathbf{g}(\mathbf{f}, \mathbf{x})| \in \mathbb{R}_{\geq 0}^d$
 - Normalize \mathbf{g}_m to induce a probability distribution over the features in the input space
- Feature Contribution: $\mathbf{g}_c = \mathbf{g}(\mathbf{f}, \mathbf{x}) \in \mathbb{R}^d$ and $|\mathbf{g}_c| = \mathbf{g}_m$
 - Presence of rain increases the chance of precipitation (+)
 - Presence of clouds detracts from the chance of sunny weather (-)

Common feature-based explanations

LIME (Ribeiro et al. KDD 2016)

$$\mathbf{g}(\mathbf{f}, \mathbf{x})_i = \arg \min_{\mathbf{g} \in \mathcal{G}} \mathcal{L}(\mathbf{f}, \mathbf{g}, \pi_{\mathbf{x}}) + \Omega(\mathbf{g})$$

Common feature-based explanations

LIME (Ribeiro et al. KDD 2016)

$$\mathbf{g}(\mathbf{f}, \mathbf{x})_i = \arg \min_{\mathbf{g} \in \mathcal{G}} \mathcal{L}(\mathbf{f}, \mathbf{g}, \pi_{\mathbf{x}}) + \Omega(\mathbf{g})$$

Local surrogate model, \mathbf{g} , to approximate f in some kernelized region $\pi_{\mathbf{x}}$, and encourages **sparsity** by keeping model complexity, $\Omega(\mathbf{g})$, low

SHAP (Lundberg and Lee. NeurIPS 2017)

$$\mathbf{g}_{ci} = \mathbf{g}(\mathbf{f}, \mathbf{x})_i = \phi_i = \frac{1}{|F|} \sum_{S \subseteq F \setminus \{i\}} \binom{F-1}{S}^{-1} (\mathbf{f}(\mathbf{x}_{S \cup \{i\}}) - \mathbf{f}(\mathbf{x}_S))$$

Considers contribution over all features, as Shapley values satisfy **efficiency**, symmetry, additivity, and dummy (zero).

Common feature-based explanations

LIME (Ribeiro et al. KDD 2016)

$$\mathbf{g}(\mathbf{f}, \mathbf{x})_i = \arg \min_{\mathbf{g} \in \mathcal{G}} \mathcal{L}(\mathbf{f}, \mathbf{g}, \pi_{\mathbf{x}}) + \Omega(\mathbf{g})$$

Local surrogate model, \mathbf{g} , to approximate f in some kernelized region $\pi_{\mathbf{x}}$, and encourages **sparsity** by keeping model complexity, $\Omega(\mathbf{g})$, low

SHAP (Lundberg and Lee. NeurIPS 2017)

$$\mathbf{g}_{ci} = \mathbf{g}(\mathbf{f}, \mathbf{x})_i = \phi_i = \frac{1}{|F|} \sum_{S \subseteq F \setminus \{i\}} \binom{F-1}{S}^{-1} (\mathbf{f}(\mathbf{x}_{S \cup \{i\}}) - \mathbf{f}(\mathbf{x}_S))$$

Considers contribution over all features, as Shapley values satisfy **efficiency**, symmetry, additivity, and dummy (zero).

Integrated Gradients (Sundarajan et al. ICML 2017)

Accumulates the gradients along a straight line path between \mathbf{x} and $\bar{\mathbf{x}}$, where $\mathbf{f}(\bar{\mathbf{x}}) \approx 0$, and satisfies **completeness**,

$$\sum_{i=1}^d \mathbf{g}(\mathbf{f}, \mathbf{x})_i = \mathbf{f}(\mathbf{x}) - \mathbf{f}(\bar{\mathbf{x}}).$$

Evaluating explanations

Evaluating explanations

Sensitivity

Do similar inputs have similar explanations?

Evaluating explanations

Sensitivity

Do similar inputs have similar explanations?

$$\mu_{\text{AVG}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \int_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z})) \mathbb{P}_{\mathbf{x}}(\mathbf{z}) d\mathbf{z}$$

Evaluating explanations

Sensitivity

Do similar inputs have similar explanations?

$$\mu_{\text{AVG}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \int_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z})) \mathbb{P}_{\mathbf{x}}(\mathbf{z}) d\mathbf{z}$$

$$\mu_{\text{MAX}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \max_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z}))$$

Let D be the distance between explanations and ρ be the distance between inputs

Evaluating explanations

Sensitivity

Do similar inputs have similar explanations?

$$\mu_{\text{AVG}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \int_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z})) \mathbb{P}_{\mathbf{x}}(\mathbf{z}) d\mathbf{z}$$

$$\mu_{\text{MAX}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \max_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z}))$$

Let D be the distance between explanations and ρ be the distance between inputs

Faithfulness

Does the explanation capture features important to the prediction?

Evaluating explanations

Sensitivity

Do similar inputs have similar explanations?

$$\mu_{\text{AVG}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \int_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z})) \mathbb{P}_{\mathbf{x}}(\mathbf{z}) d\mathbf{z}$$

$$\mu_{\text{MAX}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \max_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z}))$$

Let D be the distance between explanations and ρ be the distance between inputs

Faithfulness

Does the explanation capture features important to the prediction?

$$\mu_{\text{F}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, S) = \text{corr}\left(\frac{1}{|S|} \sum_{i \in S} \mathbf{g}(\mathbf{f}, \mathbf{x})_i, \mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_{[x_s = \bar{x}_s]})\right)$$

Fix a subset size and randomly sample subsets of that size from \mathbf{x} to estimate the Pearson Correlation Coefficient

Evaluating explanations (cont.)

Evaluating explanations (cont.)

Complexity

Is the explanation digestible?

Evaluating explanations (cont.)

Complexity

Is the explanation digestible? We define an attribution contribution distribution:

$$\mathbb{P}_A = \left\{ \frac{|\mathbf{g}(\mathbf{x})_1|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|}, \frac{|\mathbf{g}(\mathbf{x})_2|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|}, \dots, \frac{|\mathbf{g}(\mathbf{x})_d|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|} \right\}$$
$$\mu_C(\mathbf{f}, \mathbf{g}, \mathbf{x}) = H(\mathbf{x}) = \mathbb{E}_i \left[-\ln(\mathbb{P}_A(i)) \right] = -\sum_{i=1}^d \mathbb{P}_A(i) \ln(\mathbb{P}_A(i))$$

Evaluating explanations (cont.)

Complexity

Is the explanation digestible? We define an attribution contribution distribution:

$$\mathbb{P}_A = \left\{ \frac{|\mathbf{g}(\mathbf{x})_1|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|}, \frac{|\mathbf{g}(\mathbf{x})_2|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|}, \dots, \frac{|\mathbf{g}(\mathbf{x})_d|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|} \right\}$$

$$\mu_C(\mathbf{f}, \mathbf{g}, \mathbf{x}) = H(\mathbf{x}) = \mathbb{E}_i \left[-\ln(\mathbb{P}_A(i)) \right] = -\sum_{i=1}^d \mathbb{P}_A(i) \ln(\mathbb{P}_A(i))$$

The least complex explanation is one where $\mathbf{g}(\mathbf{x})_i = 1$ and the most complex explanation is one where $\mathbf{g}(\mathbf{x})_i = \frac{1}{d}$.

Aggregating Existing Techniques

Can we learn an aggregate explanation of existing techniques that does better with respect to a criterion of interest? An approach to study \mathbf{g}_{agg} can be to set the problem up as follows:

$$\mathbf{g}_{\text{agg}} = \arg \max_{\mathbf{g} \in \mathcal{G}} \mu(\mathbf{f}, \mathbf{g}), \text{ s.t. } \mathbf{g} = h(\mathcal{G}_m)$$

Aggregating Existing Techniques

Can we learn an aggregate explanation of existing techniques that does better with respect to a criterion of interest? An approach to study \mathbf{g}_{agg} can be to set the problem up as follows:

$$\mathbf{g}_{\text{agg}} = \arg \max_{\mathbf{g} \in \mathcal{G}} \mu(\mathbf{f}, \mathbf{g}), \text{ s.t. } \mathbf{g} = h(\mathcal{G}_m)$$

Three candidate methods for $h(\cdot)$.

- Convex Combination: $\mathbf{g}_{\text{agg}} = w\mathbf{g}_1 + (1 - w)\mathbf{g}_2$

Aggregating Existing Techniques

Can we learn an aggregate explanation of existing techniques that does better with respect to a criterion of interest? An approach to study \mathbf{g}_{agg} can be to set the problem up as follows:

$$\mathbf{g}_{\text{agg}} = \arg \max_{\mathbf{g} \in \mathcal{G}} \mu(\mathbf{f}, \mathbf{g}), \text{ s.t. } \mathbf{g} = h(\mathcal{G}_m)$$

Three candidate methods for $h(\cdot)$.

- Convex Combination: $\mathbf{g}_{\text{agg}} = w\mathbf{g}_1 + (1 - w)\mathbf{g}_2$
- Centroid Aggregation: $\mathbf{g}_{\text{agg}} \in \arg \min_{\mathbf{g} \in \mathcal{G}} \sum_{i=1}^m d(\mathbf{g}, \mathbf{g}_i)$

Aggregating Existing Techniques

Can we learn an aggregate explanation of existing techniques that does better with respect to a criterion of interest? An approach to study \mathbf{g}_{agg} can be to set the problem up as follows:

$$\mathbf{g}_{\text{agg}} = \arg \max_{\mathbf{g} \in \mathcal{G}} \mu(\mathbf{f}, \mathbf{g}), \text{ s.t. } \mathbf{g} = h(\mathcal{G}_m)$$

Three candidate methods for $h(\cdot)$.

- Convex Combination: $\mathbf{g}_{\text{agg}} = w\mathbf{g}_1 + (1 - w)\mathbf{g}_2$
- Centroid Aggregation: $\mathbf{g}_{\text{agg}} \in \arg \min_{\mathbf{g} \in \mathcal{G}} \sum_{i=1}^m d(\mathbf{g}, \mathbf{g}_i)$
- Bayesian Optimization: $\max_{\mathbf{g}_{\text{agg}} \in \mathcal{G}} \mu(\mathbf{g}_{\text{agg}})$ where

$$k(\mathbf{g}_i, \mathbf{g}_j) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_x} \left[k(\mathbf{g}_i(\mathbf{x}), \mathbf{g}_j(\mathbf{x})) \right] = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_x} \left[e^{-\frac{1}{2} \|\mathbf{g}_i(\mathbf{x}) - \mathbf{g}_j(\mathbf{x})\|^2} \right]$$

Generalized Aggregation

Convex Combination

$$\mathbf{g}_{agg} = \mathbf{w}^T \mathbf{G}$$

$$\mathbf{G}^T = \left(\begin{array}{c|c|c|c} & | & | & | \\ \text{SHAP}_1 & \text{LIME}_2 & \dots & \text{IG}_m \\ & | & | & | \end{array} \right)$$

Generalized Aggregation

Convex Combination

$$\mathbf{g}_{agg} = w^T G$$

$$G^T = \left(\begin{array}{c|c|c|c} & | & | & | \\ \text{SHAP}_1 & \text{LIME}_2 & \dots & \text{IG}_m \\ & | & | & | \end{array} \right)$$

$$w_{agg} \in \arg \max_w \sum_{i=1}^m \mu(w^T G_i)$$

Generalized Aggregation

Classical Rank Aggregation

$$\mathcal{G}_m = \{SHAP_1, LIME_2, \dots, IG_m\}$$

Generalized Aggregation

Classical Rank Aggregation

$$\mathcal{G}_m = \{SHAP_1, LIME_2, \dots, IG_m\}$$

$$\mathbf{g}_c^S = [1 \quad -2 \quad 7] \rightarrow \mathbf{g}_m^S = [.1 \quad .2 \quad .7] \rightarrow \text{rank}^S = [C \quad B \quad A]$$

Generalized Aggregation

Classical Rank Aggregation

$$\mathcal{G}_m = \{SHAP_1, LIME_2, \dots, IG_m\}$$

$$\mathbf{g}_c^S = [1 \quad -2 \quad 7] \rightarrow \mathbf{g}_m^S = [.1 \quad .2 \quad .7] \rightarrow \text{rank}^S = [C \quad B \quad A]$$

$$\text{rank}^{S_1} = [C \quad B \quad A] \quad \text{rank}^{S_2} = [C \quad A \quad B] \quad \text{rank}^{S_3} = [A \quad B \quad C]$$

Generalized Aggregation

Classical Rank Aggregation

$$\mathcal{G}_m = \{SHAP_1, LIME_2, \dots, IG_m\}$$

$$\mathbf{g}_c^S = [1 \quad -2 \quad 7] \rightarrow \mathbf{g}_m^S = [.1 \quad .2 \quad .7] \rightarrow \text{rank}^S = [C \quad B \quad A]$$

$$\text{rank}^{S_1} = [C \quad B \quad A] \quad \text{rank}^{S_2} = [C \quad A \quad B] \quad \text{rank}^{S_3} = [A \quad B \quad C]$$

$$\text{Borda Count: } \mathbf{g}_{agg} = \text{rank}^{agg} = [C \quad A \quad B]$$

Generalized Aggregation

Classical Rank Aggregation

$$\mathcal{G}_m = \{SHAP_1, LIME_2, \dots, IG_m\}$$

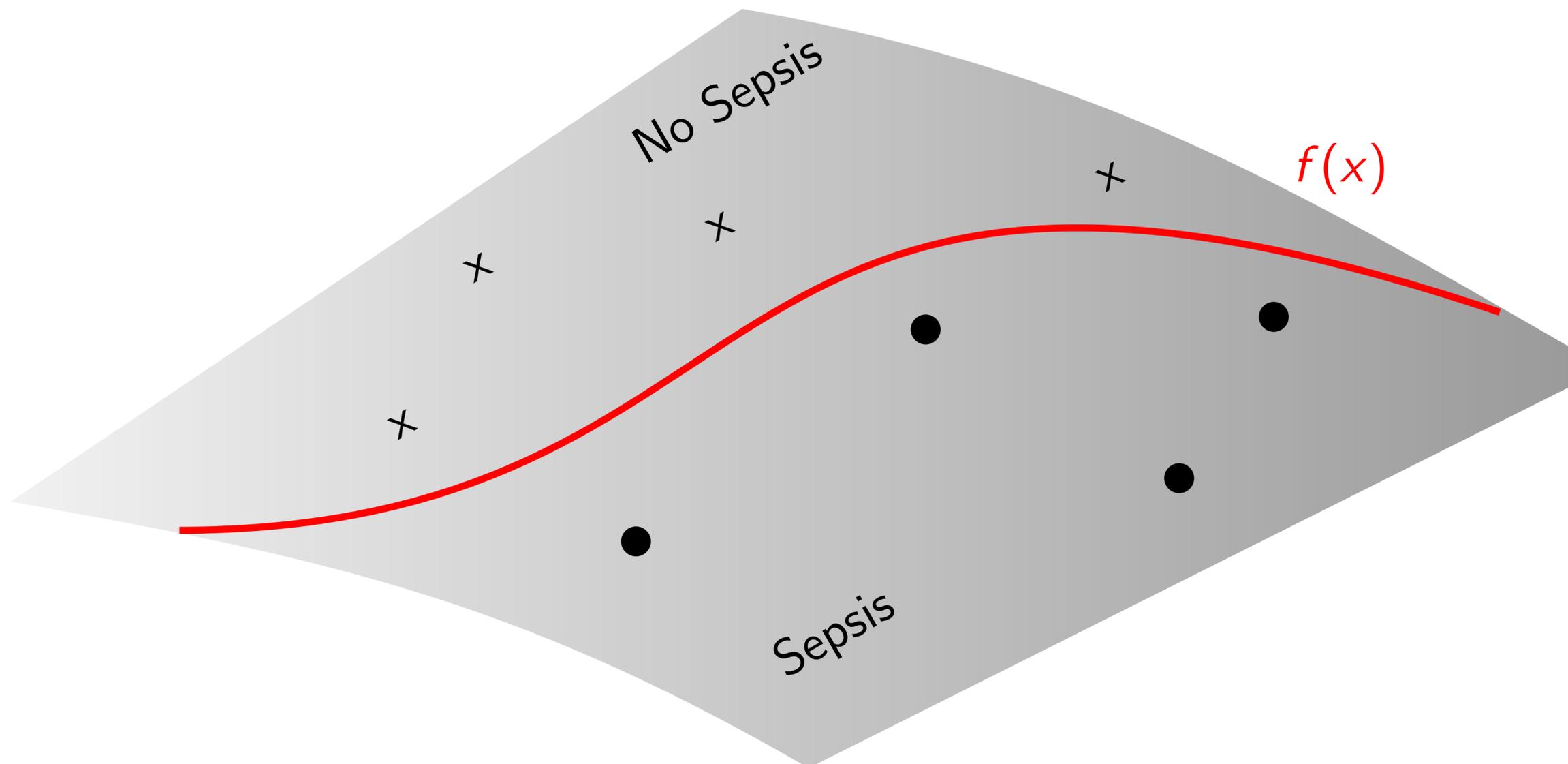
$$\mathbf{g}_c^S = [1 \quad -2 \quad 7] \rightarrow \mathbf{g}_m^S = [.1 \quad .2 \quad .7] \rightarrow \text{rank}^S = [C \quad B \quad A]$$

$$\text{rank}^{S_1} = [C \quad B \quad A] \quad \text{rank}^{S_2} = [C \quad A \quad B] \quad \text{rank}^{S_3} = [A \quad B \quad C]$$

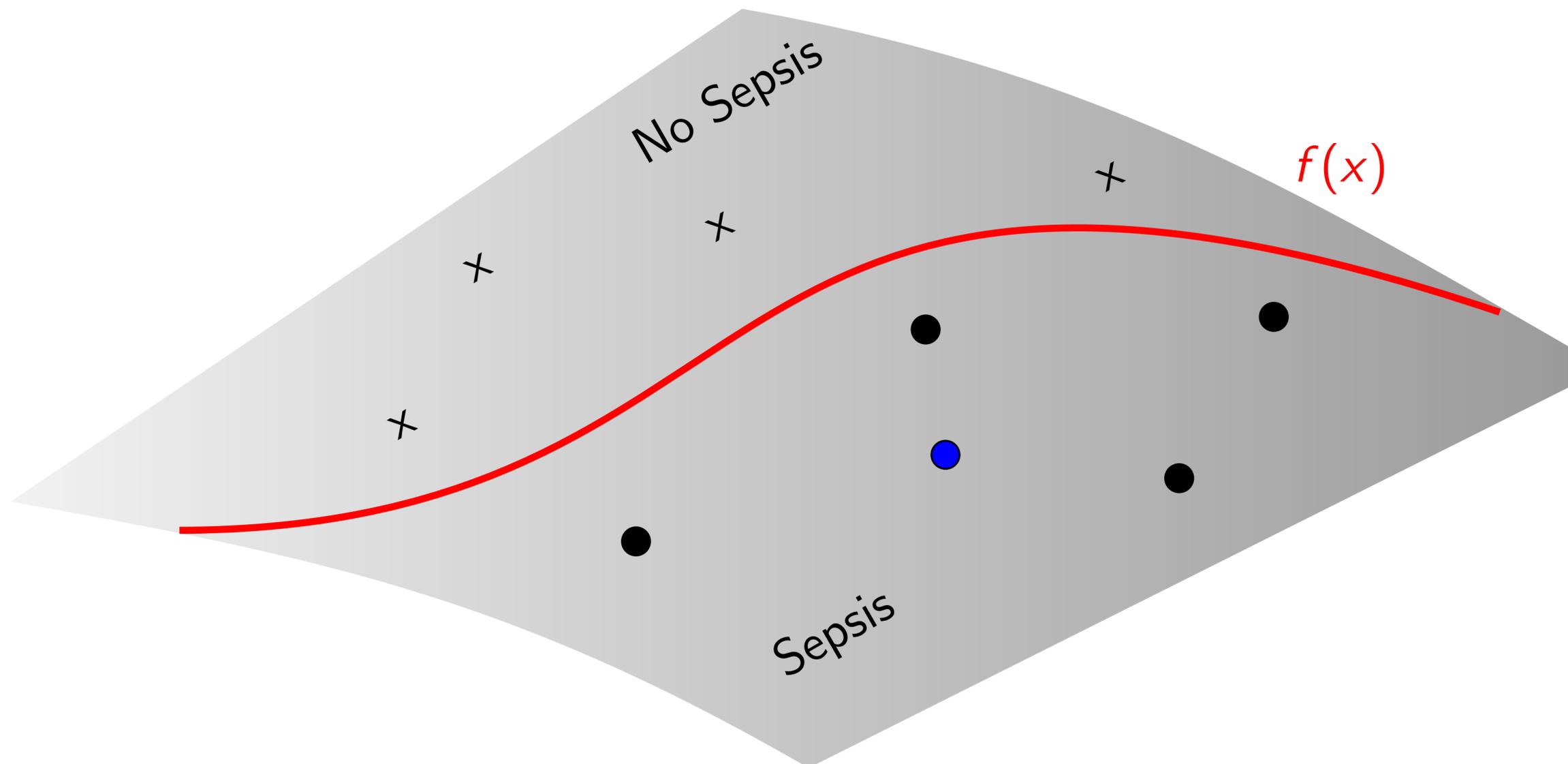
$$\text{Borda Count: } \mathbf{g}_{agg} = \text{rank}^{agg} = [C \quad A \quad B]$$

$$\mathbf{g}_{agg} \in \arg \min_{\mathbf{g}} \sum_{\mathbf{g}_i \in \mathcal{G}_m} d(\mathbf{g}, \mathbf{g}_i)$$

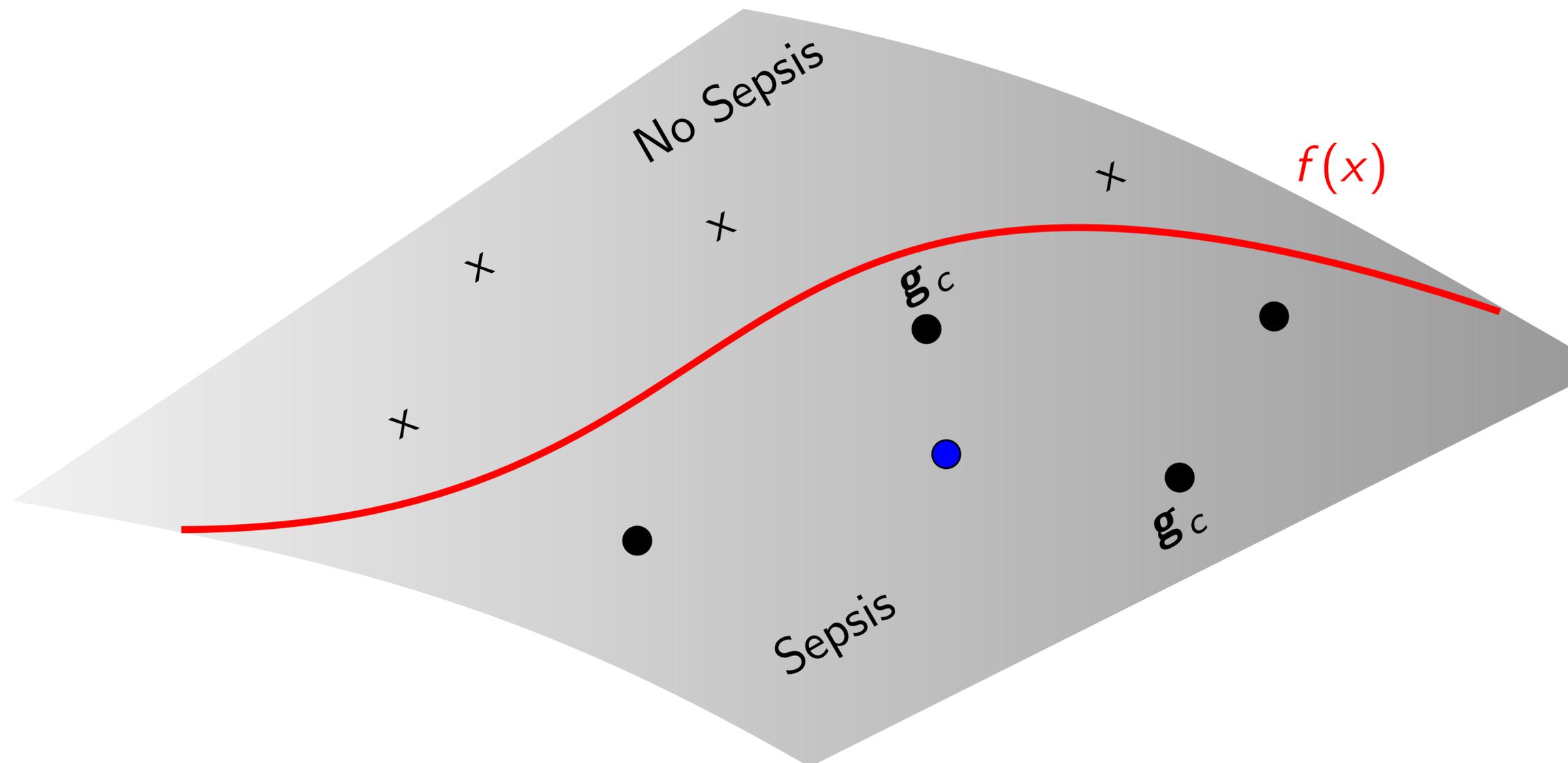
Aggregating Local Explanations



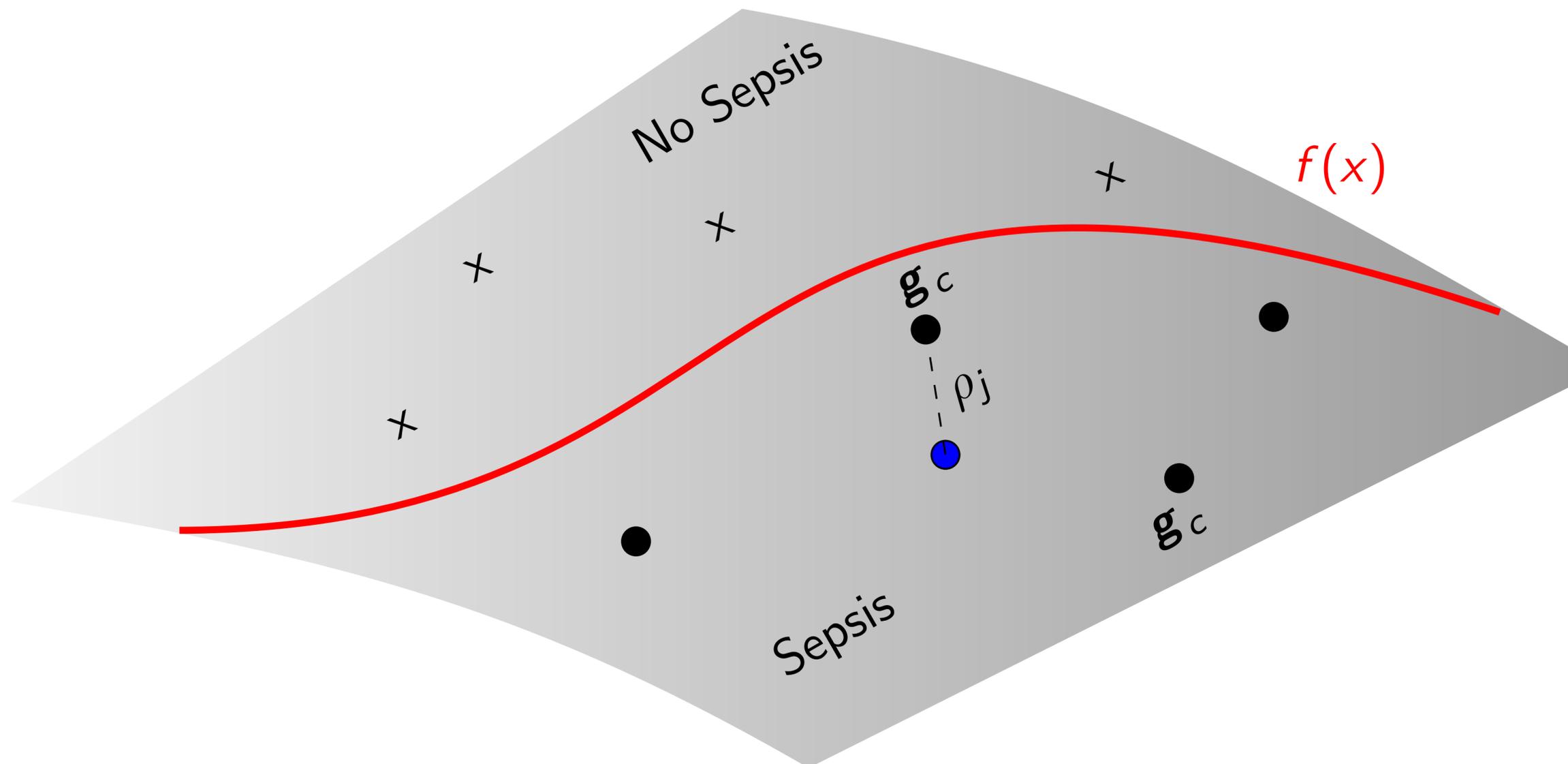
Aggregating Local Explanations



Aggregating Local Explanations



Aggregating Local Explanations



AVA: Aggregate Valuation of Antecedents

Can we use weighted Shapley values (Kalai et al. Journal of Game Theory 1987) to aggregate feature-based explanations with lower sensitivity?

AVA: Aggregate Valuation of Antecedents

Can we use weighted Shapley values (Kalai et al. Journal of Game Theory 1987) to aggregate feature-based explanations with lower sensitivity?

- 1 Find k nearest neighbors, \mathcal{N}_k , of x_{test} and their weights, ρ_j

$$\rho_j = \frac{d}{d\epsilon} \mathcal{L}(f_{\epsilon, x^{(j)}}, x_{\text{test}}) \Big|_{\epsilon=0}$$

$$\mathcal{N}_k(x_{\text{test}}, \mathcal{D}) = \arg \max_{\mathcal{N} \subset \mathcal{D}, |\mathcal{N}|=k} \sum_{x^{(j)} \in \mathcal{N}} \rho_j$$

AVA: Aggregate Valuation of Antecedents

Can we use weighted Shapley values (Kalai et al. Journal of Game Theory 1987) to aggregate feature-based explanations with lower sensitivity?

- 1 Find k nearest neighbors, \mathcal{N}_k , of x_{test} and their weights, ρ_j

$$\rho_j = \frac{d}{d\epsilon} \mathcal{L}(f_{\epsilon, x^{(j)}}, x_{\text{test}}) \Big|_{\epsilon=0}$$

$$\mathcal{N}_k(x_{\text{test}}, \mathcal{D}) = \arg \max_{\mathcal{N} \subset \mathcal{D}, |\mathcal{N}|=k} \sum_{x^{(j)} \in \mathcal{N}} \rho_j$$

- 2 Calculate the attributions, \mathbf{g}_c , for all points in \mathcal{N}_k

$$\mathbf{g}_{ci} = \phi_i = \frac{1}{|F|} \sum_{S \subseteq F \setminus \{i\}} \binom{F-1}{S}^{-1} (f(x_{S \cup \{i\}}) - f(x_S))$$

AVA: Aggregate Valuation of Antecedents

Can we use weighted Shapley values (Kalai et al. Journal of Game Theory 1987) to aggregate feature-based explanations with lower sensitivity?

- 1 Find k nearest neighbors, \mathcal{N}_k , of x_{test} and their weights, ρ_j

$$\rho_j = \frac{d}{d\epsilon} \mathcal{L}(f_{\epsilon, x^{(j)}}, x_{\text{test}}) \Big|_{\epsilon=0}$$

$$\mathcal{N}_k(x_{\text{test}}, \mathcal{D}) = \arg \max_{\mathcal{N} \subset \mathcal{D}, |\mathcal{N}|=k} \sum_{x^{(j)} \in \mathcal{N}} \rho_j$$

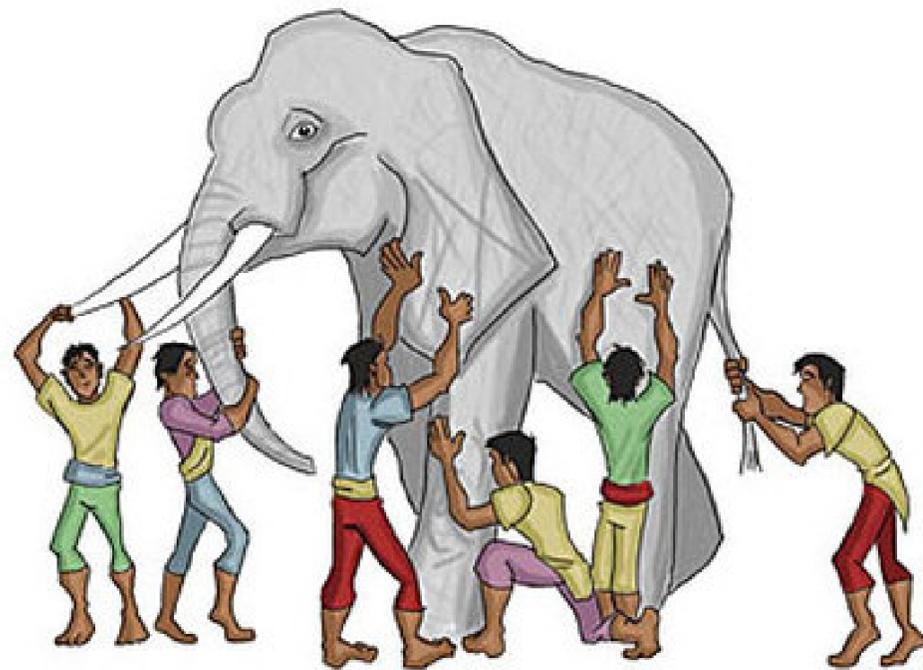
- 2 Calculate the attributions, \mathbf{g}_c , for all points in \mathcal{N}_k

$$\mathbf{g}_{ci} = \phi_i = \frac{1}{|F|} \sum_{S \subseteq F \setminus \{i\}} \binom{F-1}{S}^{-1} (f(x_{S \cup \{i\}}) - f(x_S))$$

- 3 Aggregate the k explanations into a consensus, \mathbf{g}_{agg}

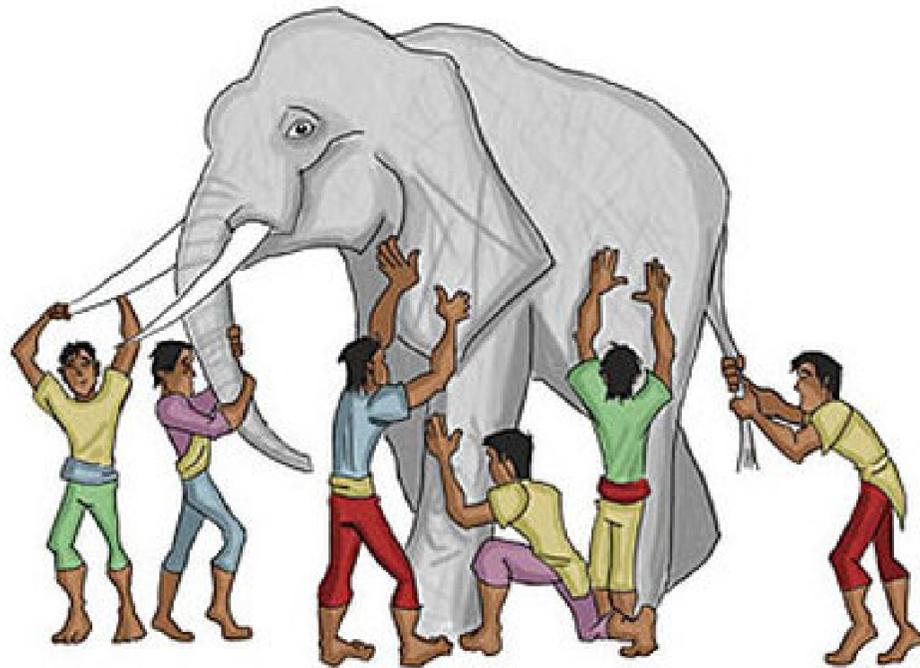
$$\mathbf{g}_{\text{agg}} = \sum_{x^{(j)} \in \mathcal{N}_k} \frac{\rho_j}{\rho} \mathbf{g}_c^j$$

Why aggregate?

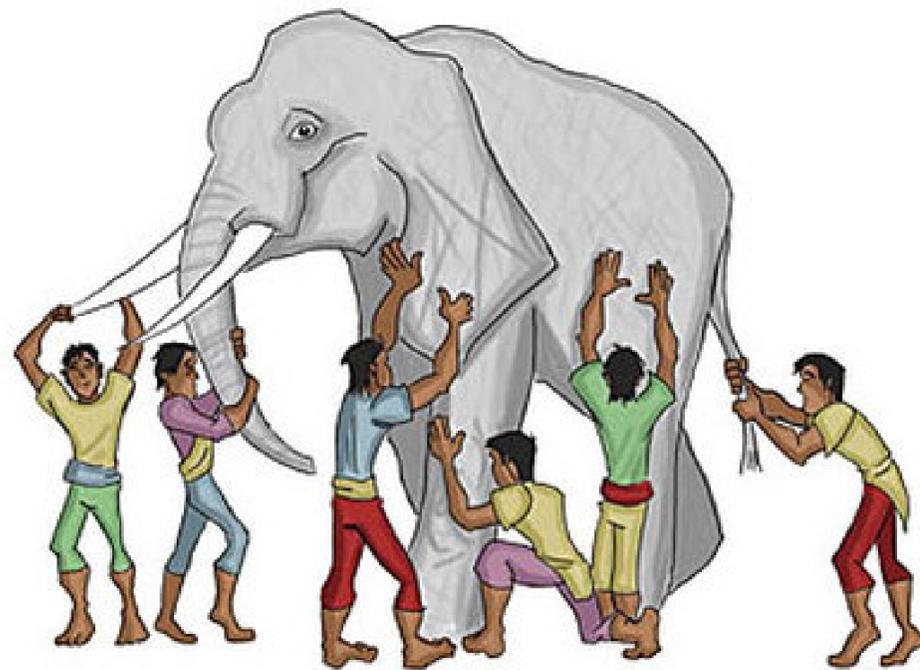


Why aggregate?

- Each training point has its own “learned” attribution

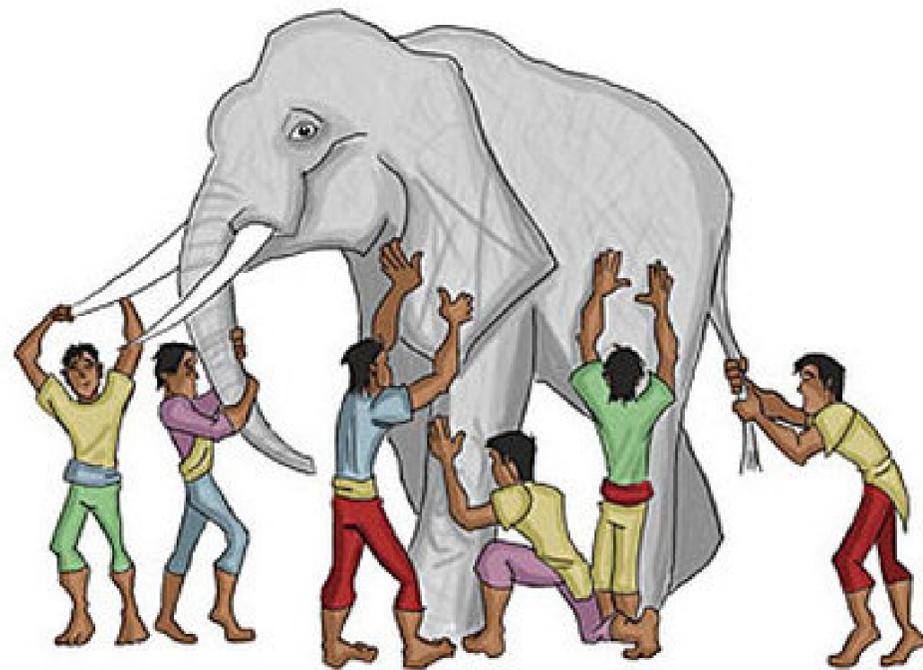


Why aggregate?



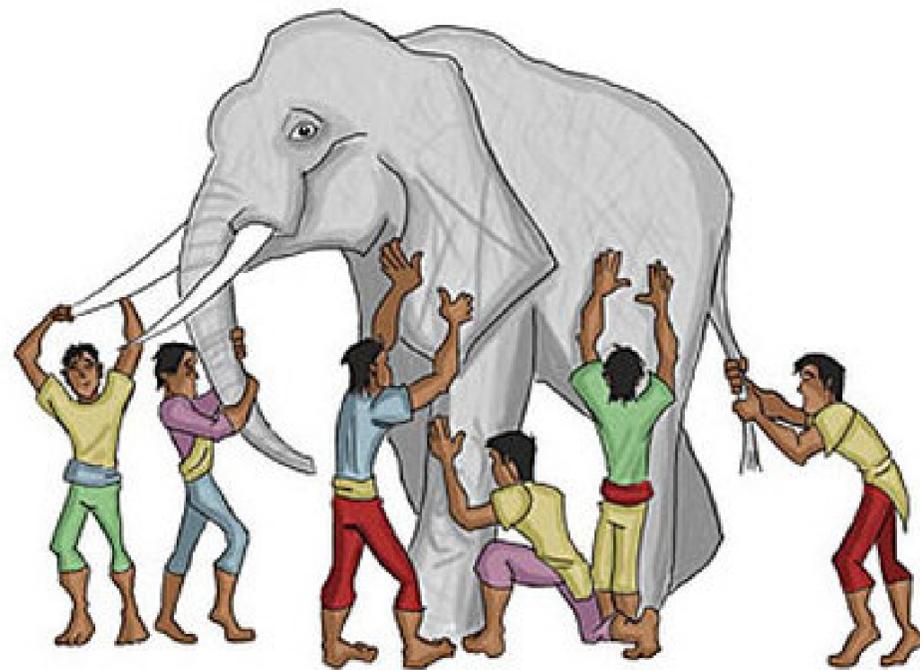
- Each training point has its own “learned” attribution
- Aggregate explanation now has lower sensitivity

Why aggregate?



- Each training point has its own “learned” attribution
- Aggregate explanation now has lower sensitivity
- Resulting attribution uses motivating reasoning of a doctor

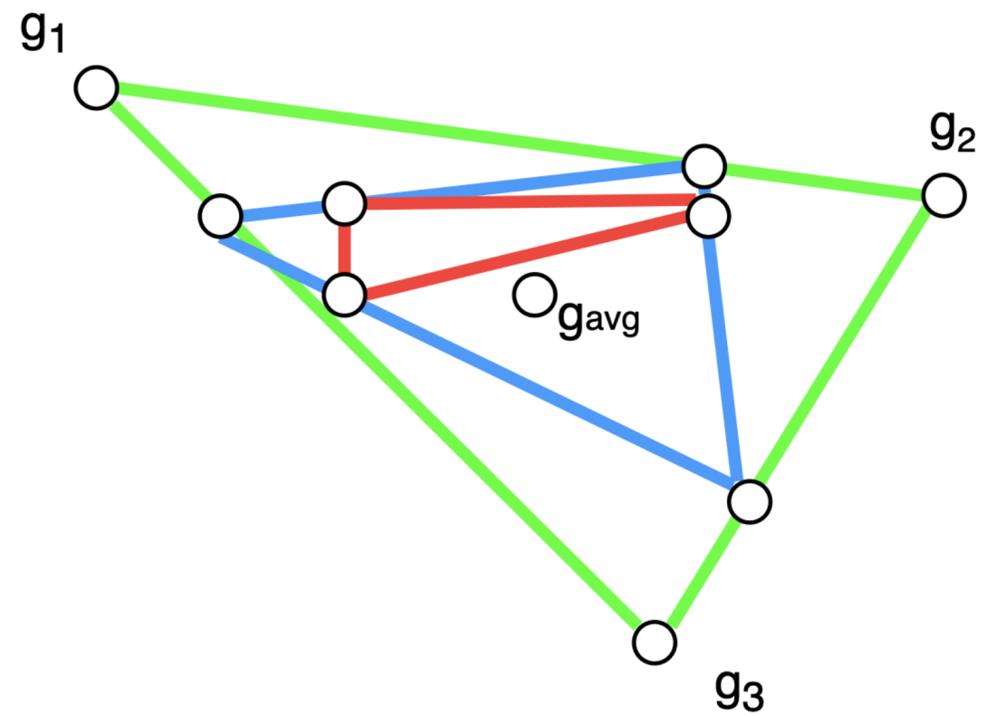
Why aggregate?



- Each training point has its own “learned” attribution
- Aggregate explanation now has lower sensitivity
- Resulting attribution uses motivating reasoning of a doctor
- **SUPER SUPER** cheap to compute

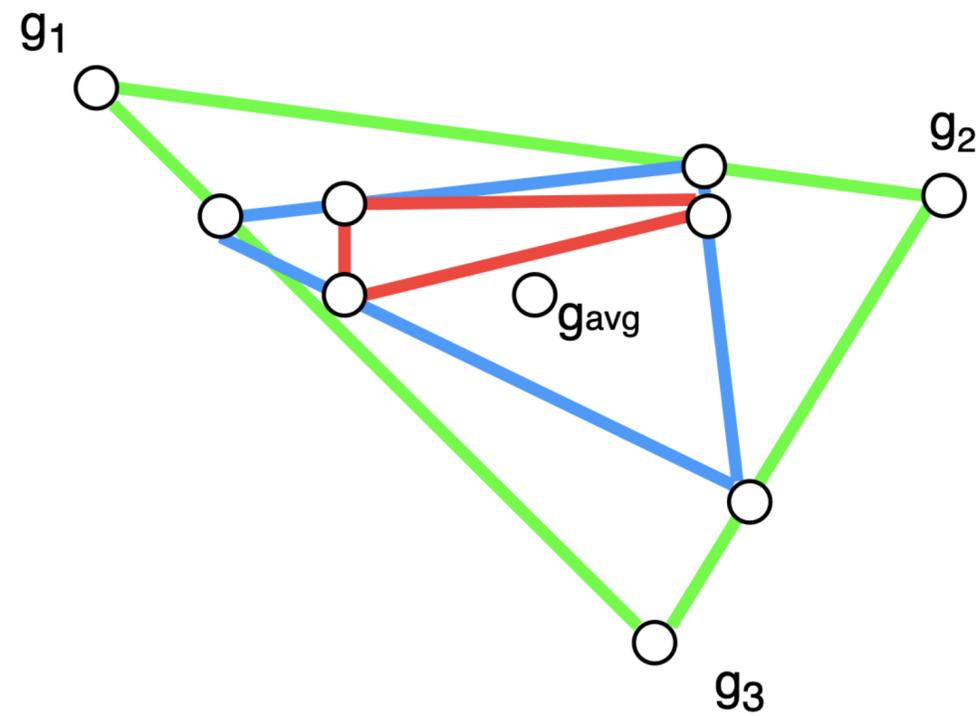
Minimizing Complexity

Region Shrinking Method

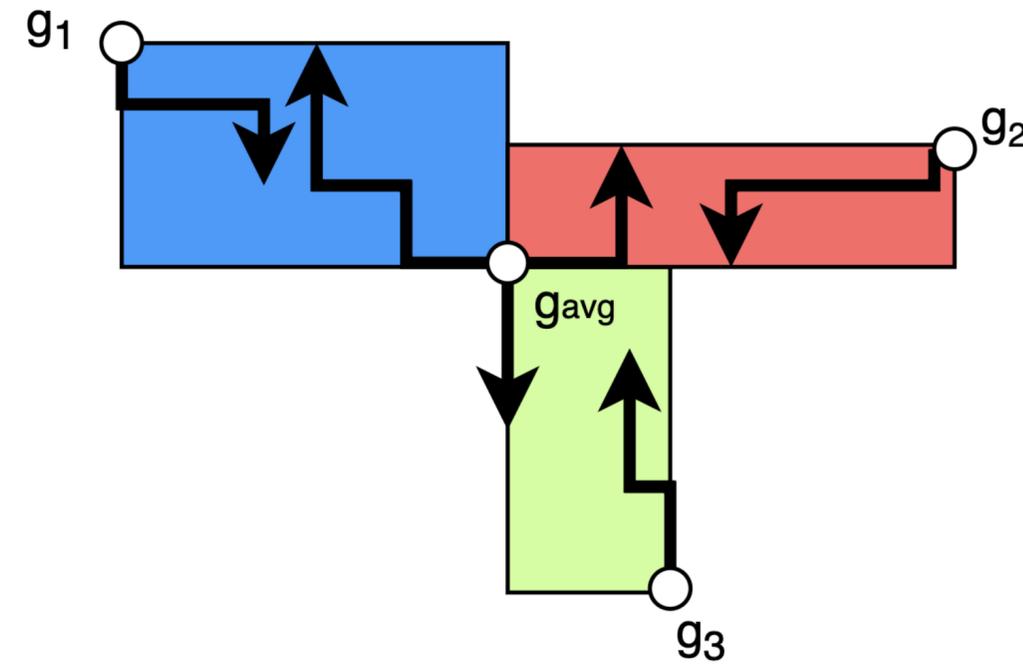


Minimizing Complexity

Region Shrinking Method



Gradient-Descent Style Method



Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful

Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value

Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation

Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.

Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.
- 2 Are feature-based explanations even useful?

Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.
- 2 Are feature-based explanations even useful?
 - Consider counterfactuals, natural language explanations, etc.

Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.
- 2 Are feature-based explanations even useful?
 - Consider counterfactuals, natural language explanations, etc.
- 3 Working with experts to find a \mathbf{g}^*

Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.
- 2 Are feature-based explanations even useful?
 - Consider counterfactuals, natural language explanations, etc.
- 3 Working with experts to find a \mathbf{g}^*
- 4 Multi-Objective optimization

Summary

This Work

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.
- 2 Are feature-based explanations even useful?
 - Consider counterfactuals, natural language explanations, etc.
- 3 Working with experts to find a \mathbf{g}^*
- 4 Multi-Objective optimization
 - Resulting Setup

$$\max \text{faithfulness}(\mathbf{g}_{agg}) + \text{sensitivity}(\mathbf{g}_{agg})$$

Practitioner-Driven Questions

1. Can we provide methodology for practitioners to **evaluate** explanations?
2. Can existing explainability tools be used to identify model unfairness?
3. Can we quantify how much uncertainty is associated with given explanations?

Practitioner-Driven Questions

1. Can we provide methodology for practitioners to evaluate explanations?
2. Can existing explainability tools be used to identify model **unfairness**?
3. Can we quantify how much uncertainty is associated with given explanations?

You Shouldn't Trust Me: Learning Models Which Conceal Unfairness From Multiple Explanation Methods

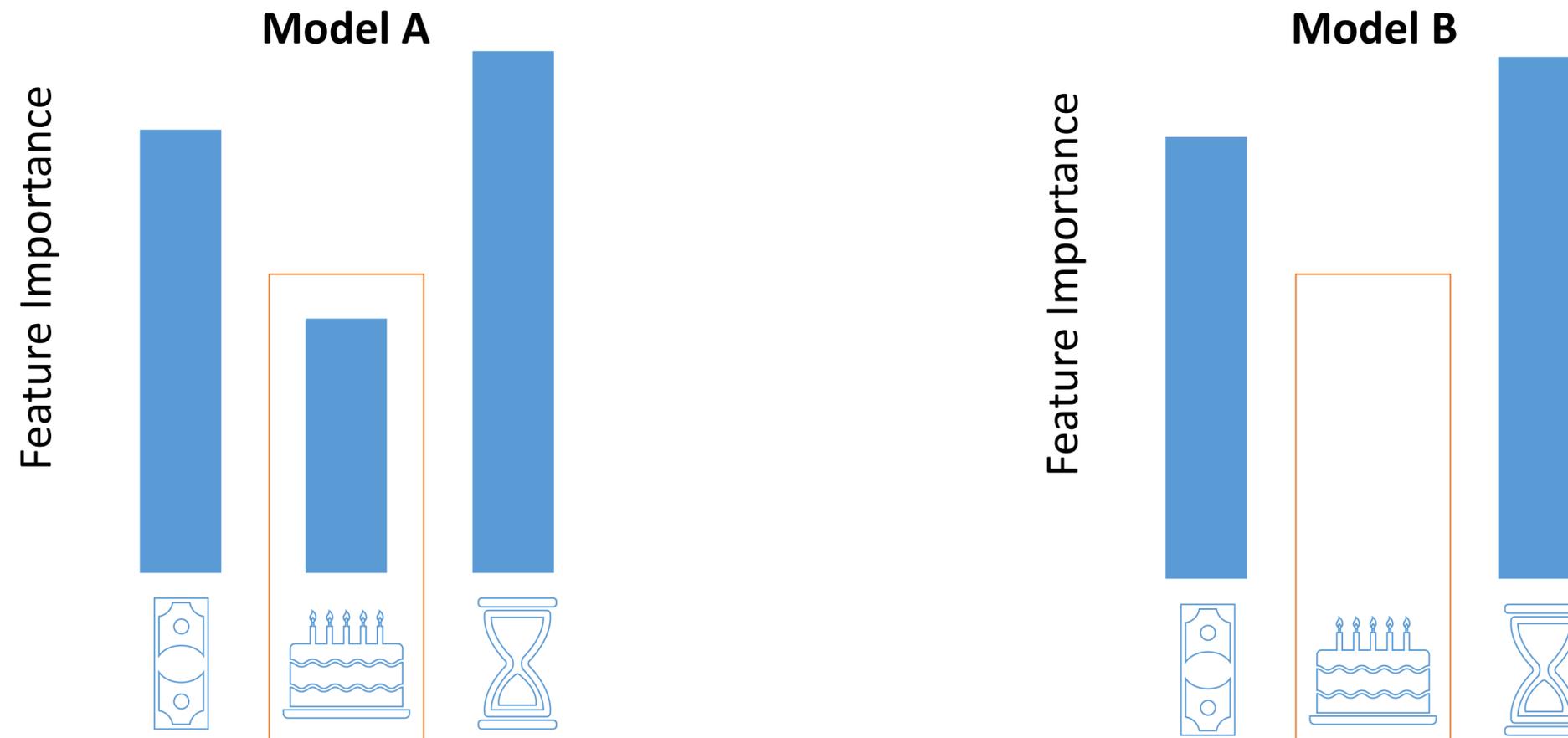
Botty Dimanov, **Umang Bhatt**, Mateja Jamnik, and Adrian Weller

Appeared at the European Conference on Artificial Intelligence 2020
http://ecai2020.eu/papers/72_paper.pdf

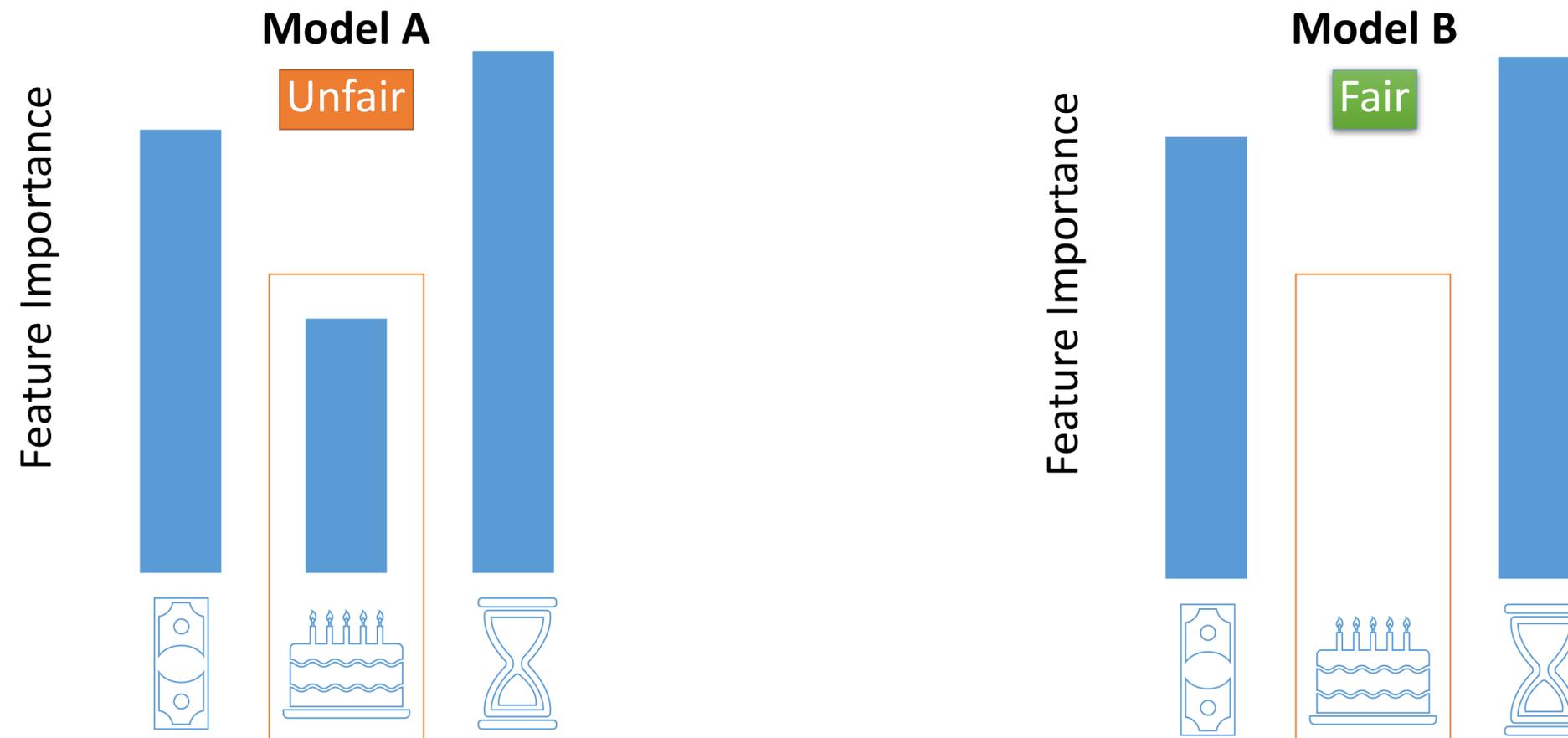


The
Alan Turing
Institute

Why do we care?



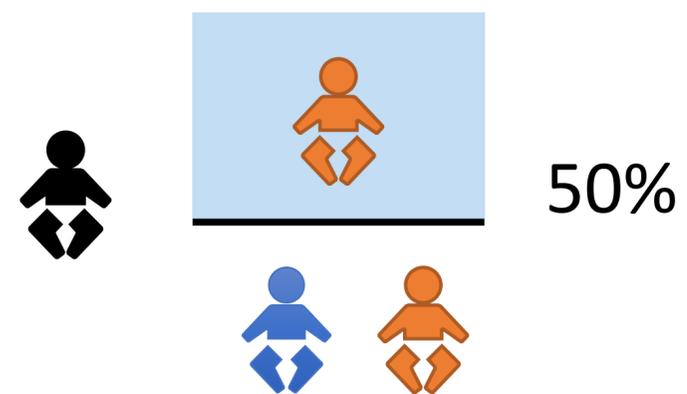
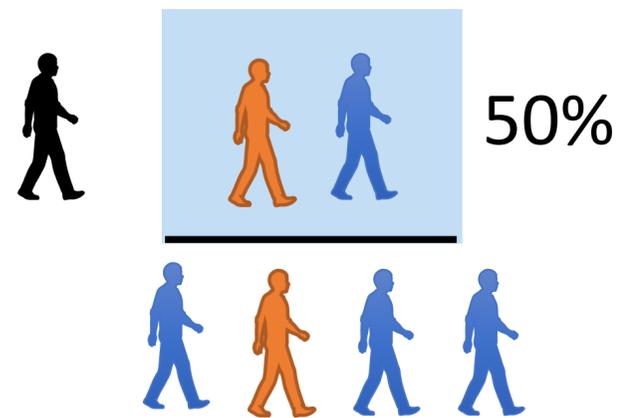
Why do we care?



What is fairness?

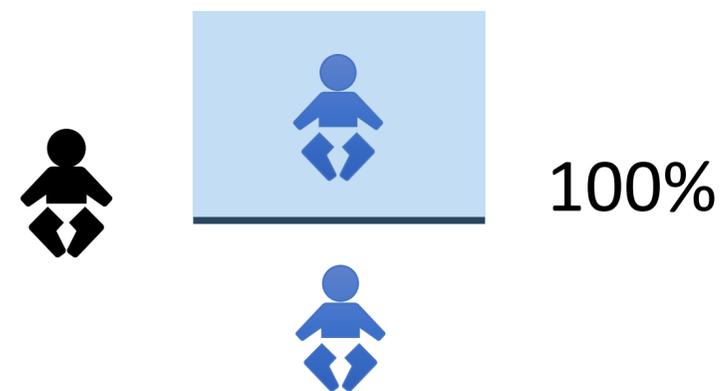
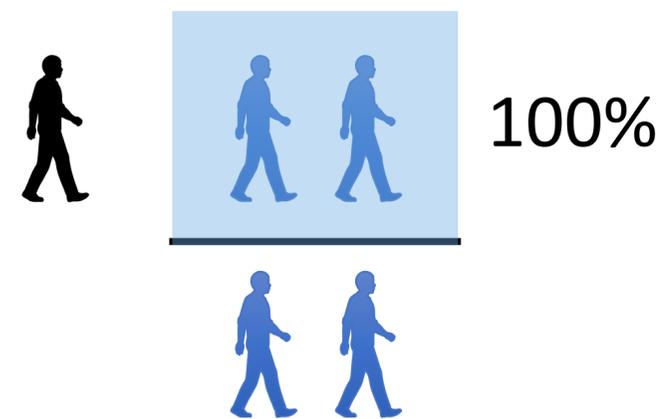
Demographic Parity

positive rate



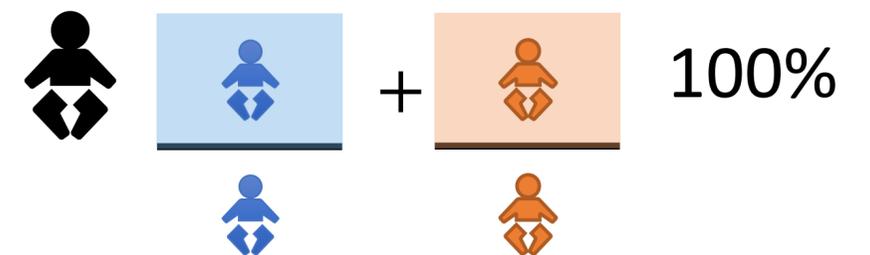
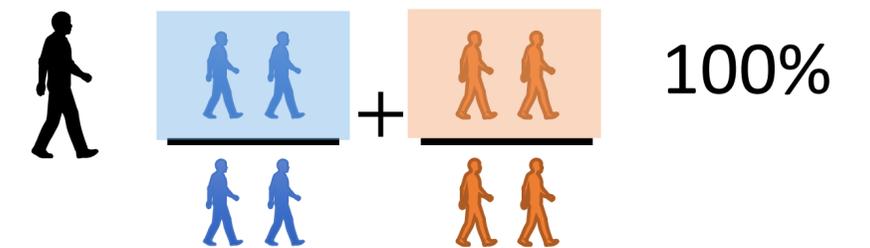
Equal Opportunity

true positive rate



Equal Accuracy

true positive +
true negative



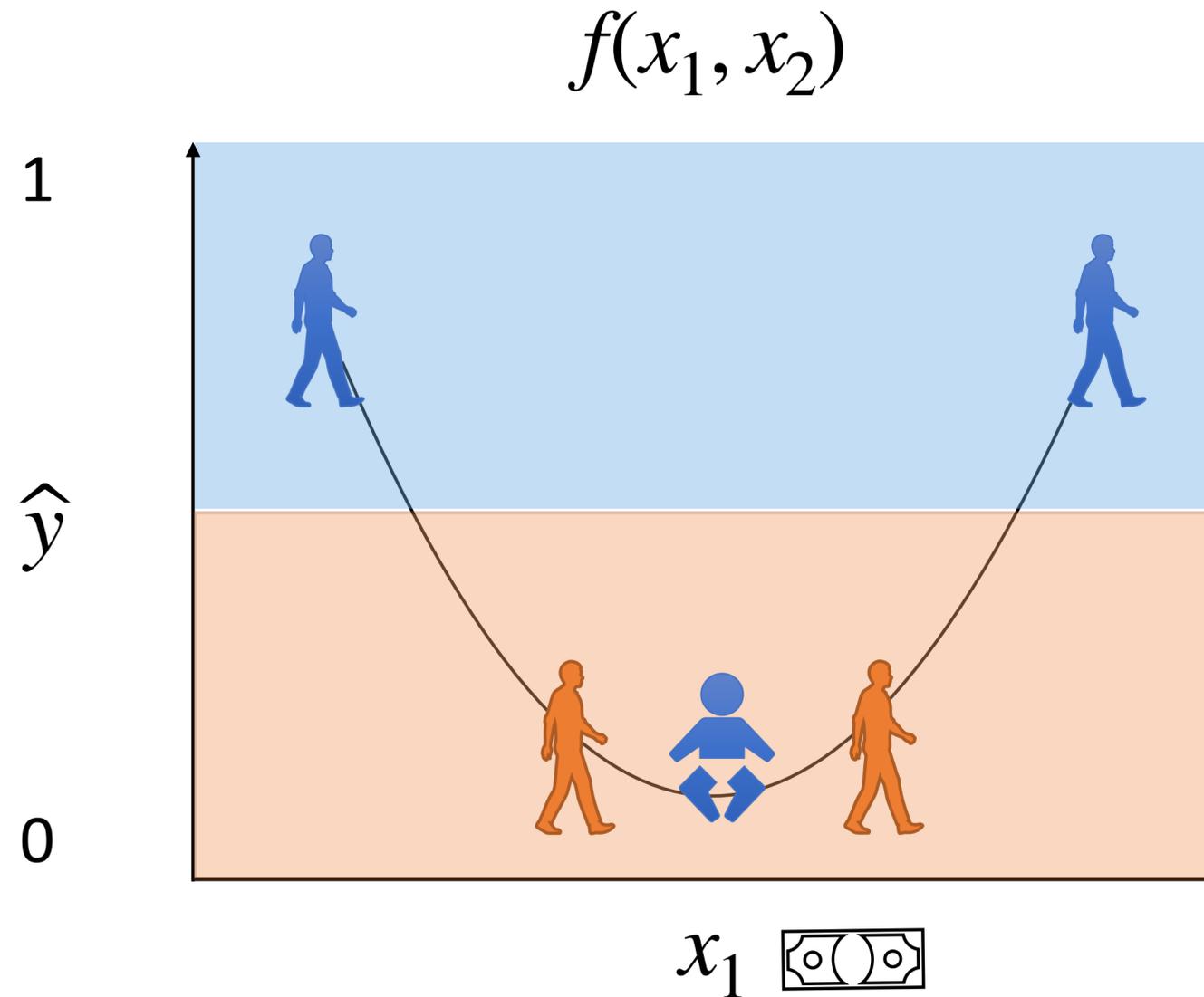
Toy Example

$x_1 \in R$ 

$x_2 \in \{0,1\}$  , 

$y \in \{0,1\}$  , 

$$\frac{\partial f(x_1, x_2)}{\partial x_2} = 0$$



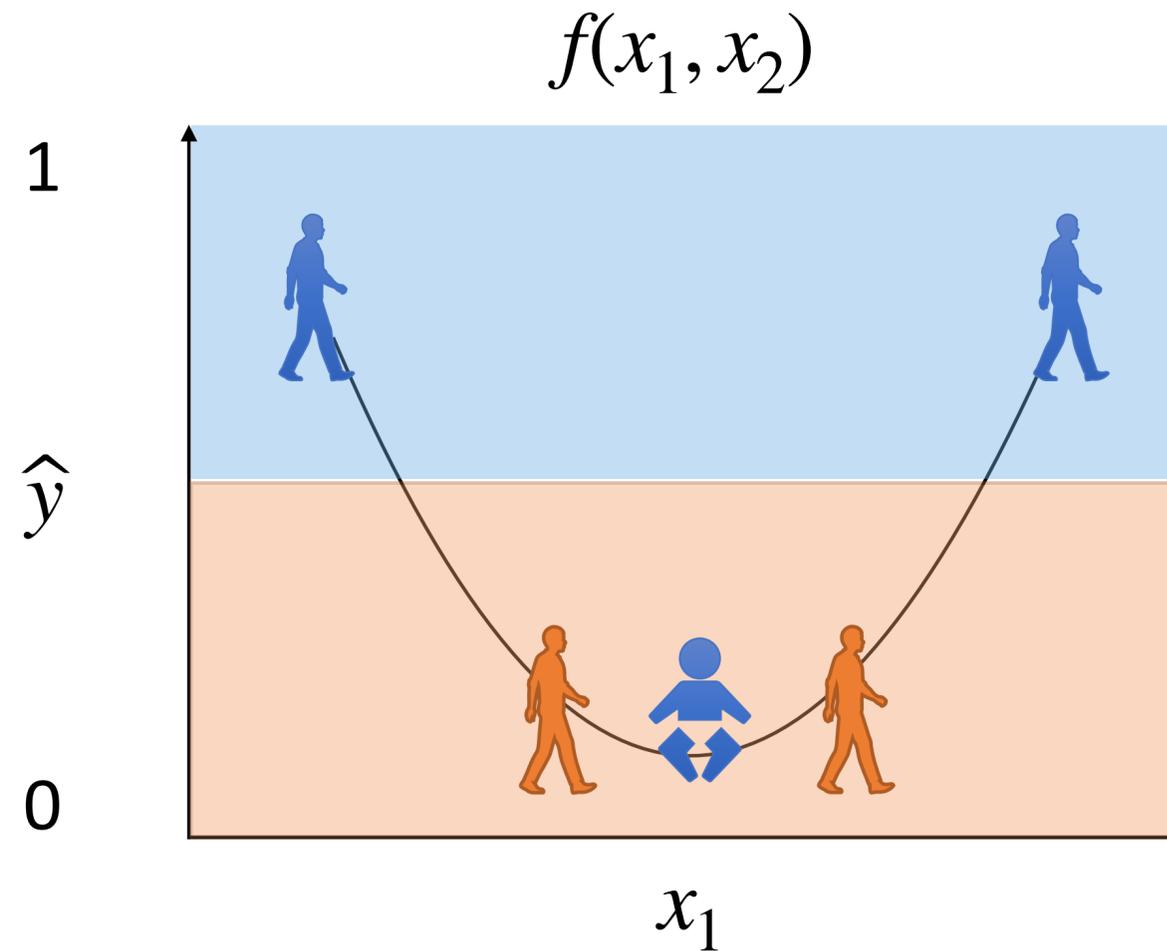
Toy Example

$x_1 \in R$ 

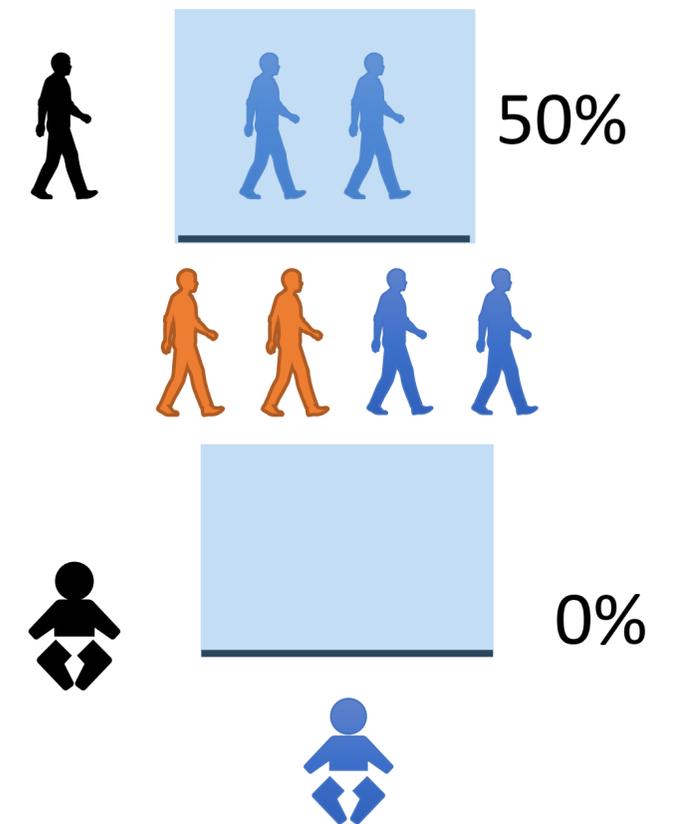
$x_2 \in \{0,1\}$ , 

$y \in \{0,1\}$ , 

$\frac{\partial f(x_1, x_2)}{\partial x_2} = 0$



Demographic Parity



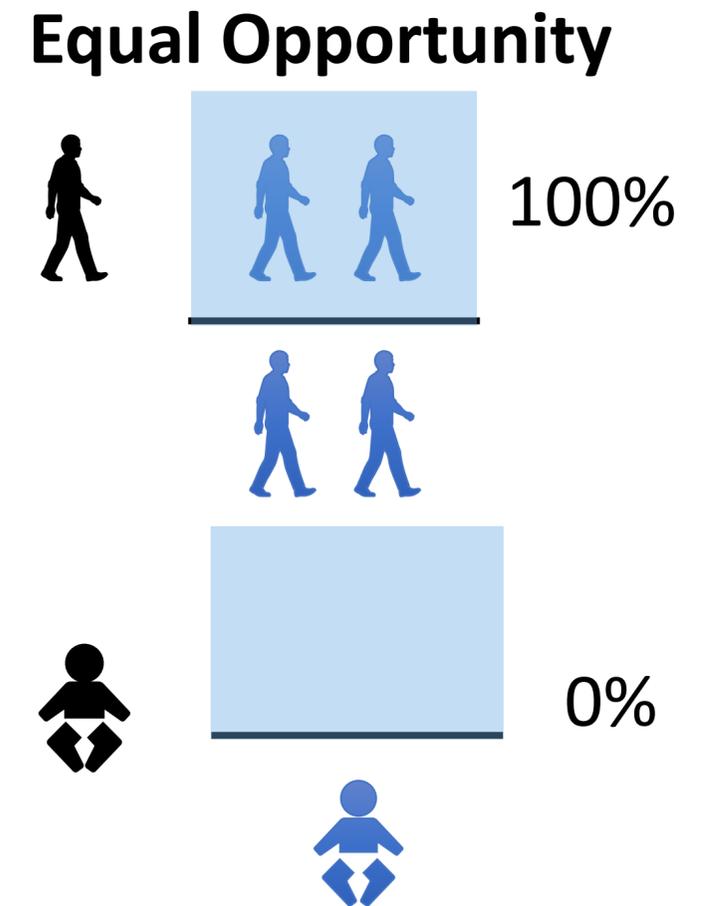
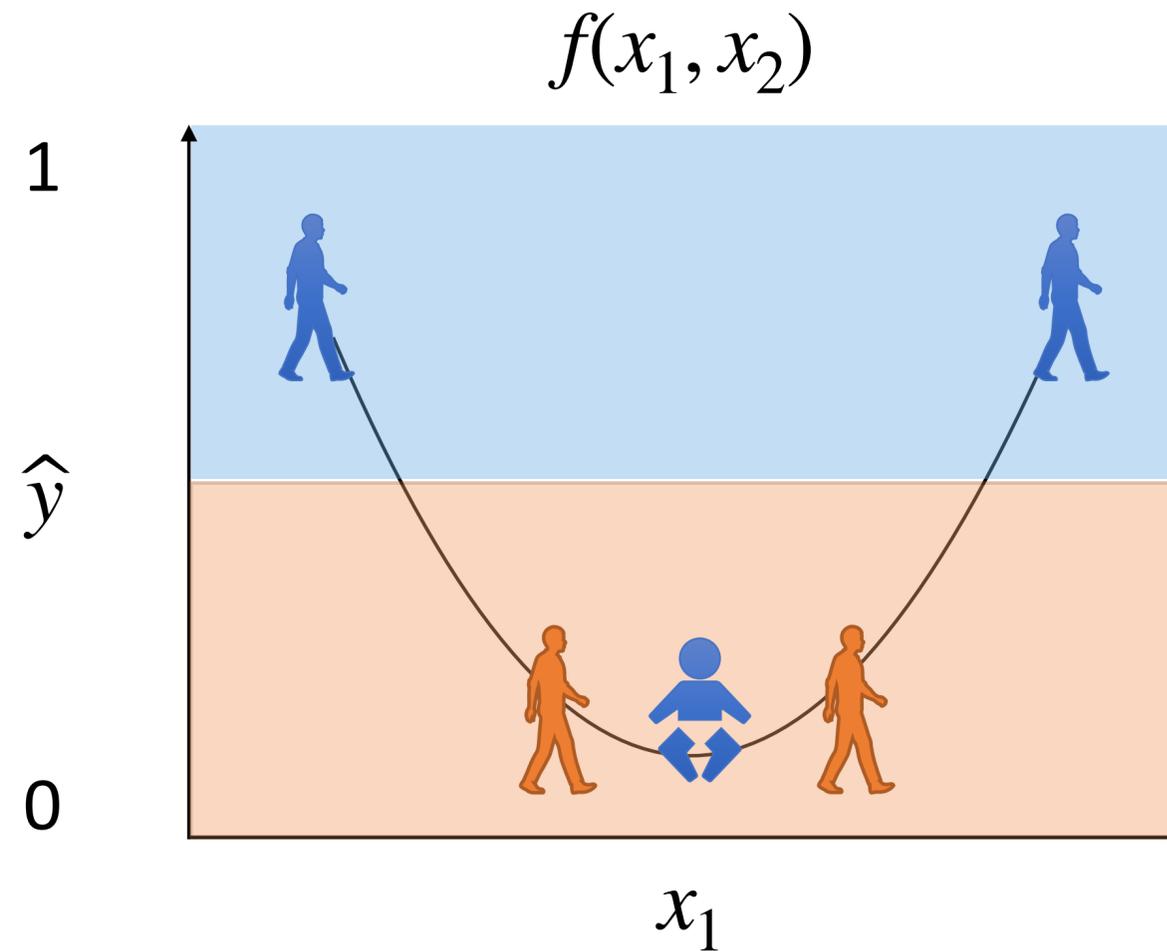
Toy Example

$x_1 \in R$ 

$x_2 \in \{0,1\}$ , 

$y \in \{0,1\}$ , 

$\frac{\partial f(x_1, x_2)}{\partial x_2} = 0$



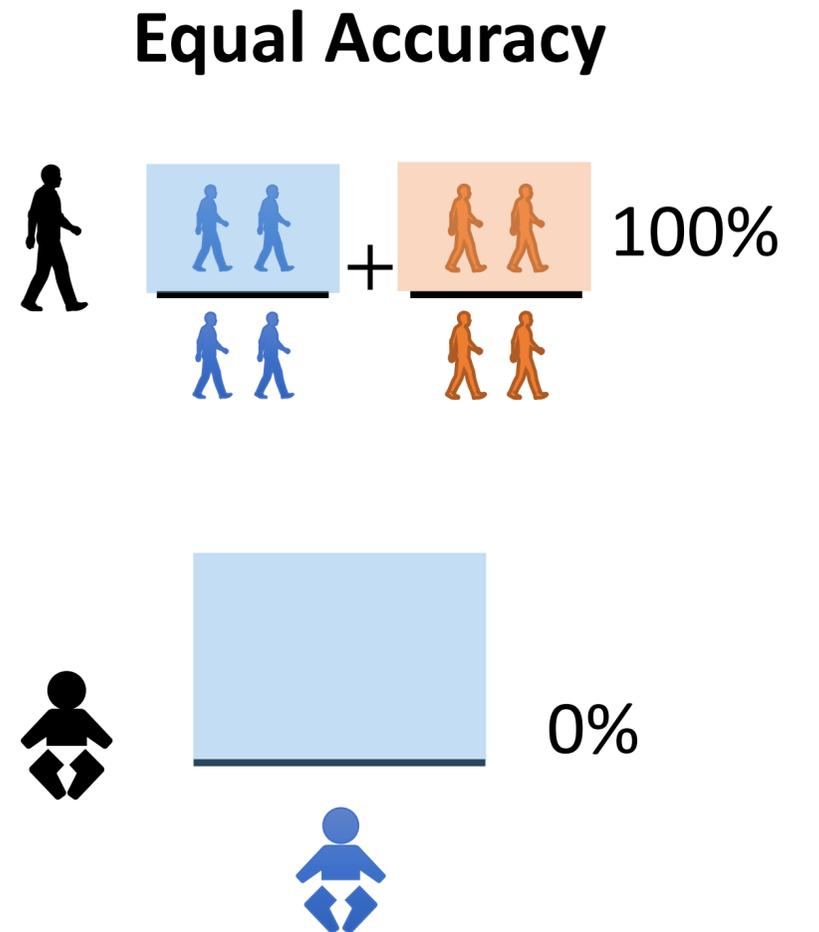
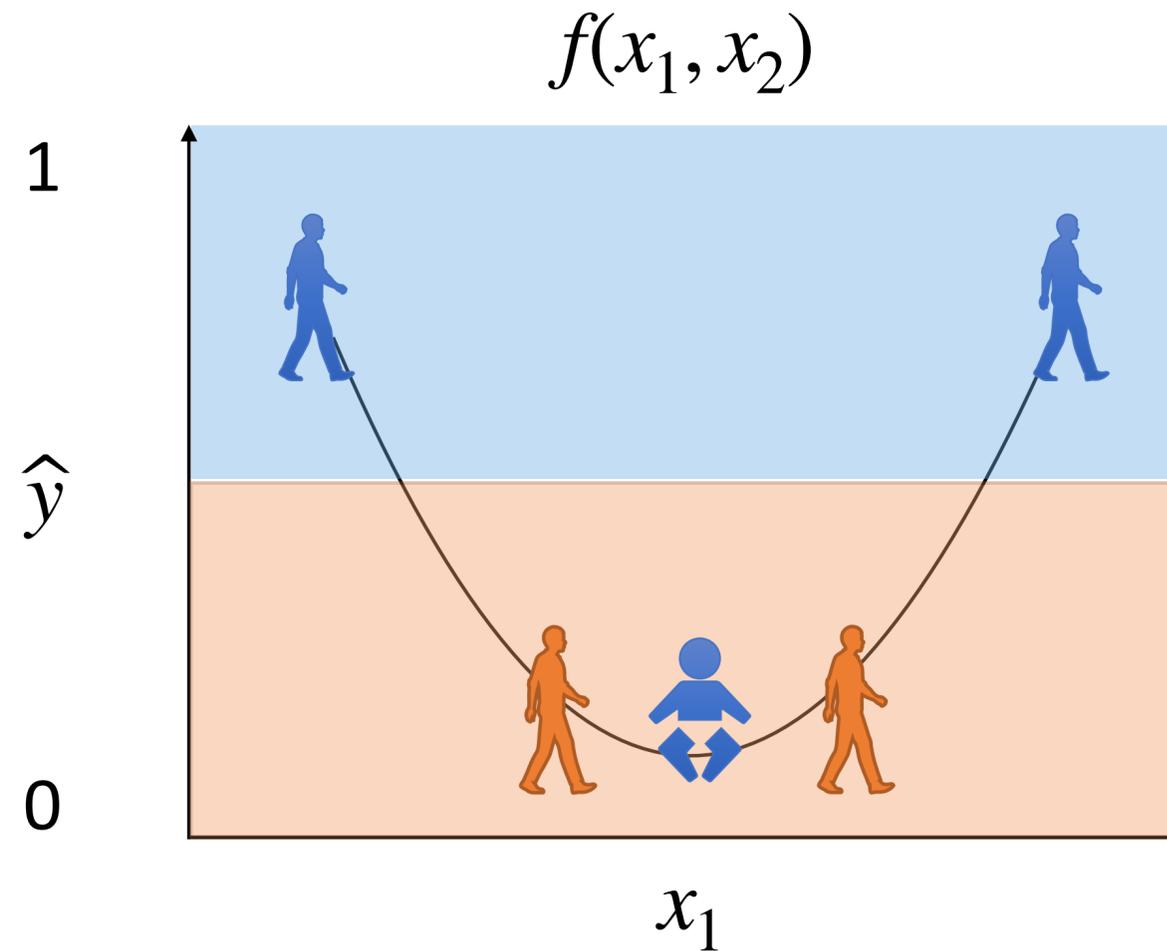
Toy Example

$x_1 \in R$ 

$x_2 \in \{0,1\}$ , 

$y \in \{0,1\}$ , 

$\frac{\partial f(x_1, x_2)}{\partial x_2} = 0$



Can we manipulate explanations?

Modified explanations via adversarial perturbations of **inputs**

- Ghorbani, Abid, and Zou. AAAI 2019
- Dombrowski et al. NeurIPS 2019
- Slack et al. AIES 2019

Control visual explanations via adversarial perturbations of **parameters**

- Heo, Joo, and Moon. NeurIPS 2019

Downgrade explanations via adversarial perturbations of **parameters to hide unfairness**

Our Setup

Classifier $f : \mathbf{X} \mapsto \mathbf{Y}$

Explanation $g(f, \mathbf{x})_j$

Our Goal

$f_\theta \longrightarrow f_{\theta+\delta}$

Desirable Properties

Model Similarity $\forall i, f_{\theta+\delta}(\mathbf{x}^{(i)}) \approx f_\theta(\mathbf{x}^{(i)})$

Low Target Feature Attribution $\forall i, |g(f_{\theta+\delta}, \mathbf{x}^{(i)})_j| \ll |g(f_\theta, \mathbf{x}^{(i)})_j|$

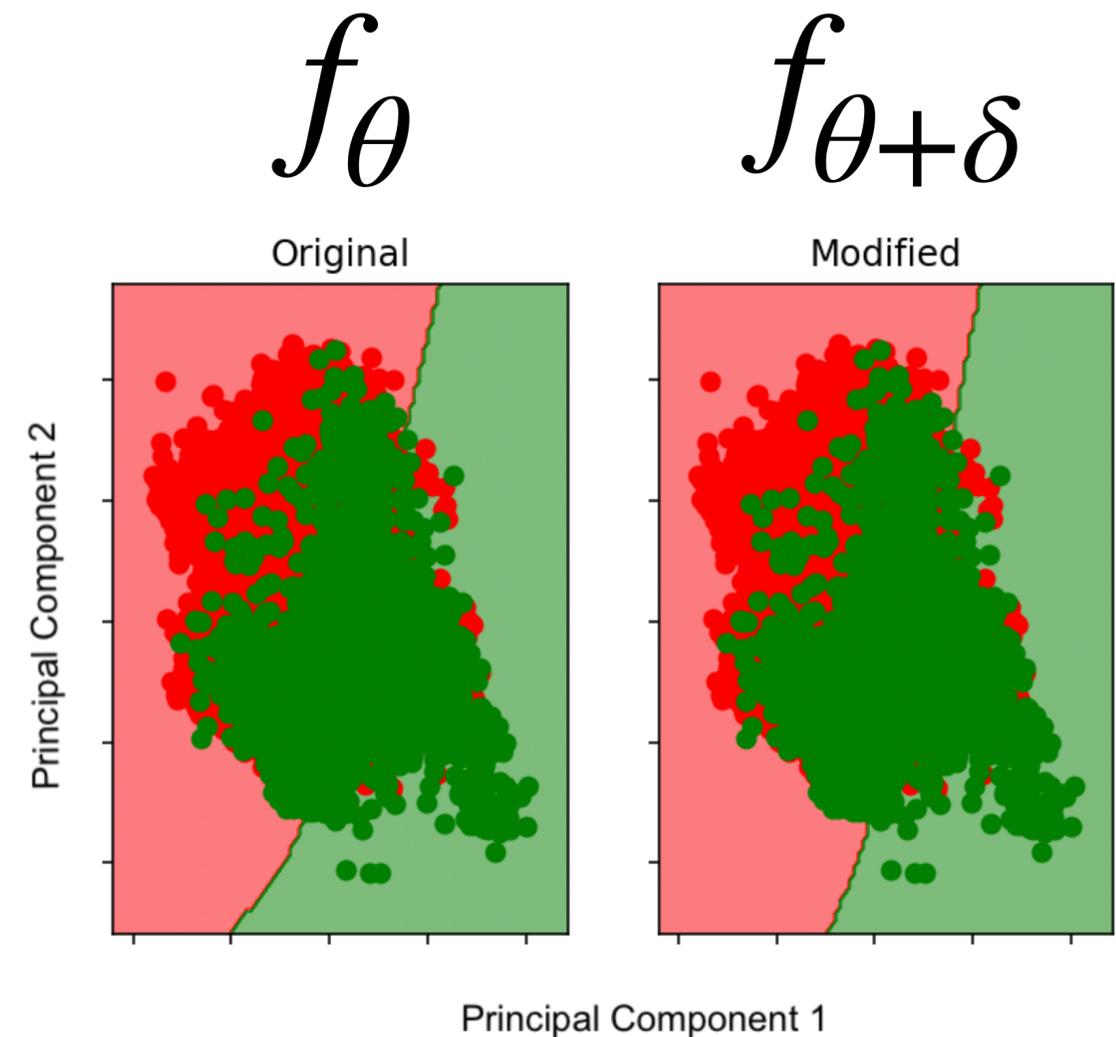
Our Method

Adversarial Explanation Attack

$$\operatorname{argmin}_{\delta} L' = L(f_{\theta+\delta}, \mathbf{x}, \mathbf{y}) + \frac{\alpha}{n} \left\| \left\| \nabla_{\mathbf{x}_{:,j}} L(f_{\theta+\delta}, \mathbf{x}, \mathbf{y}) \right\| \right\|_p$$

Findings

- Little change in accuracy, but difference in outputs is detectable
- Low attribution achieved with respect to **multiple** explanation methods
- High unfairness across multiple fairness metrics (compared to holding feature constant)



Takeaway

Feature importance reveals nothing reliable about model fairness.

Practitioner-Driven Questions

1. Can we provide methodology for practitioners to evaluate explanations?
2. Can existing explainability tools be used to identify model **unfairness**?
3. Can we quantify how much uncertainty is associated with given explanations?

Practitioner-Driven Questions

1. Can we provide methodology for practitioners to evaluate explanations?
2. Can existing explainability tools be used to identify model unfairness?
3. Can we quantify how much **uncertainty** is associated with given explanations?

Effects of Uncertainty on the Quality of Feature Importance Explanations

Torgyn Shaikhina,* **Umang Bhatt,*** Roxanne Zhang, Konstantinos
Georgatzis, Alice Xiang, and Adrian Weller

Appeared at the AAAI Workshop on Explainable Agency in Artificial Intelligence 2021



Motivation

- A. **Vulnerability** of feature-based model explanations
- B. Practitioners may use such explanations for **model selection**, which leads to confirmation bias, over-reliance, and potential manipulation

Our Assumptions

- [Assumption] **Fixed model assumption is restrictive:** Existing methods for generating explanations tend to explain a point estimate, which assumes that only **one** suitable model exists.
- [Remark] Obtain uncertainty estimates by **subsampling** from the posterior of models (from the **same** model class and with the **same** model specification) given training samples and initial states.

Shapley Value

$$\phi_i(v, x, f) = \frac{1}{|T|} \sum_{S \subseteq T \setminus \{i\}} \binom{T-1}{S}^{-1} (v_x(S \cup \{i\}, f) - v_x(S, f))$$

$$v_x(S, f) = \mathbb{E} [f(z) \mid z = \bar{x}_{[\bar{x}_s = x_s]}]$$

Our Approach

$$\phi_i(v, x)^* = \int_{f \in \mathcal{F}} P(f | \mathcal{D}) \phi_i(v, x, f) df$$

$$\tilde{\phi}_i(v, x) = \sum_{f \in \mathcal{B}} w_f \phi_i(v, x, f) = \mathbb{E}_{f \in \mathcal{B}}[\phi_i(v, x, f)]$$

Our Approach

$$\phi_i(v, x)^* = \int_{f \in \mathcal{F}} P(f | \mathcal{D}) \phi_i(v, x, f) df$$

$$\tilde{\phi}_i(v, x) = \sum_{f \in \mathcal{B}} w_f \phi_i(v, x, f) = \mathbb{E}_{f \in \mathcal{B}}[\phi_i(v, x, f)]$$

Variance

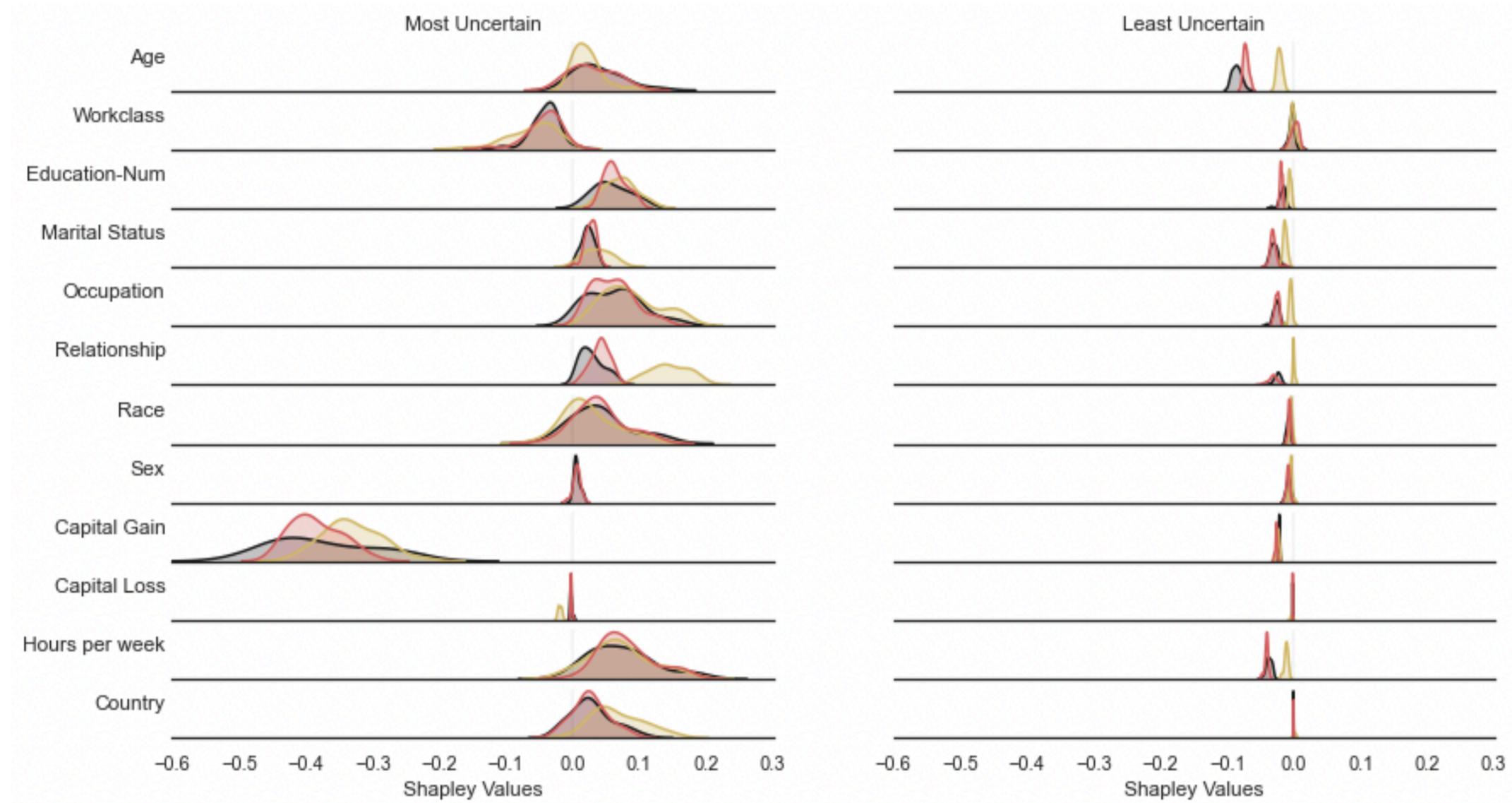
$$s_i(\mathcal{B}, x) = \mathbb{V}_{f \in \mathcal{B}}[\phi_i(v, x, f)]$$

Experimental Setup

- Model is uncertain on a point if its prediction lies 1.5 times the IQR
- Set of models found by (i) changing initialization or (ii) subsampling from the training data
- Multiple Evaluation Criteria: Complexity, Monotonicity, Efficiency, Faithfulness...

Finding #1

Predictive uncertainty leads to high variance in explanations



Feature importance scores of an OOD, most uncertain (left) and least uncertain (right) samples of the Adult dataset computed using Exact (red), Tree (black) and Kernel (yellow) SHAP. Each feature is presented as a distribution of its Shapley values associated with the ensemble of 30 Gradient Boosted Trees.

*Similar empirical evidence for other tree-based models and neural nets

Finding #2

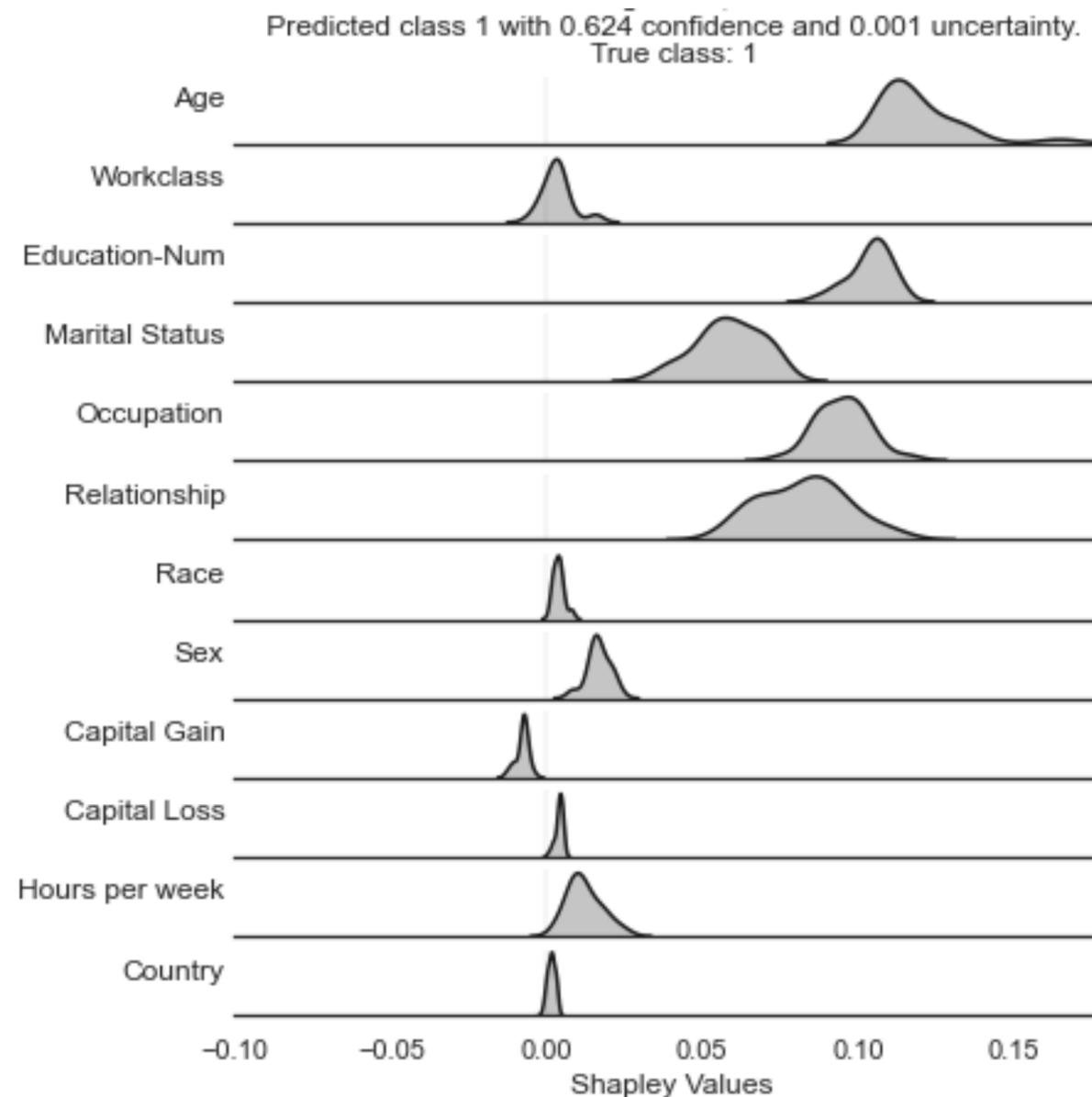
Uncertainty degrades quality of explanations across w.r.t. all metrics:

		Variance	Complexity	Monotonicity	Efficiency	Faithfulness
Exact Shapley	non-OOD	0.07 ± 0.00	1.88 ± 0.22	0.88 ± 0.04	1.00 ± 0.00	0.00 ± 0.24
	OOD	0.39 ± 0.00	1.96 ± 0.14	0.85 ± 0.05	1.00 ± 0.00	-0.01 ± 0.23
TreeSHAP	non-OOD	0.08 ± 0.00	1.90 ± 0.20	0.88 ± 0.04	0.99 ± 0.09	0.00 ± 0.24
	OOD	0.44 ± 0.00	1.96 ± 0.13	0.85 ± 0.04	0.97 ± 0.18	-0.01 ± 0.23
KernelSHAP	non-OOD	0.19 ± 0.00	1.81 ± 0.23	0.89 ± 0.05	1.00 ± 0.00	0.00 ± 0.24
	OOD	0.809 ± 0.00	1.85 ± 0.16	0.82 ± 0.05	1.00 ± 0.00	0.00 ± 0.23
Integrated Gradients	non-OOD	20.00 ± 5.34	1.53 ± 0.22	0.96 ± 0.03	0.99 ± 0.10	0.00 ± 0.16
	OOD	46.48 ± 18.61	1.55 ± 0.18	0.96 ± 0.03	1.00 ± 0.03	0.00 ± 0.13
Smoothgrad	non-OOD	20.35 ± 0.52	1.92 ± 0.11	0.92 ± 0.04	1.00 ± 0.10	0.00 ± 0.13
	OOD	19.97 ± 0.43	1.98 ± 0.06	0.92 ± 0.03	1.00 ± 0.09	-0.02 ± 0.10

Comparison of the explanation quality on OOD and in-distribution samples of the Adult dataset. Presented are mean and standard deviation of quality scores averaged across the ensemble of models. Values in bold indicate where the degradation in quality was statistically significant ($p < 0.05$).

Implication #1

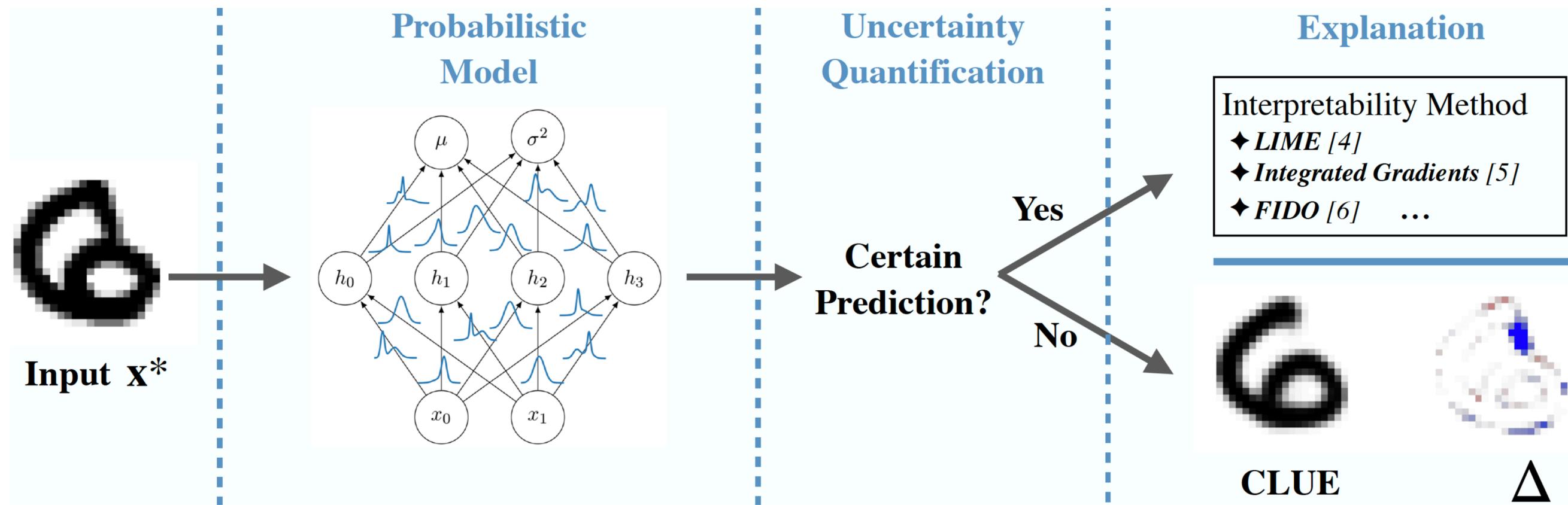
When interpreting feature importance scores, data practitioners must consider the entire distribution (not just point estimates).



Implication #2

Don't use feature importance on uncertain data — it's unreliable...

Use counterfactual latent uncertainty explanations (CLUEs) instead!



Getting a CLUE: A Method for Explaining Uncertainty Estimates

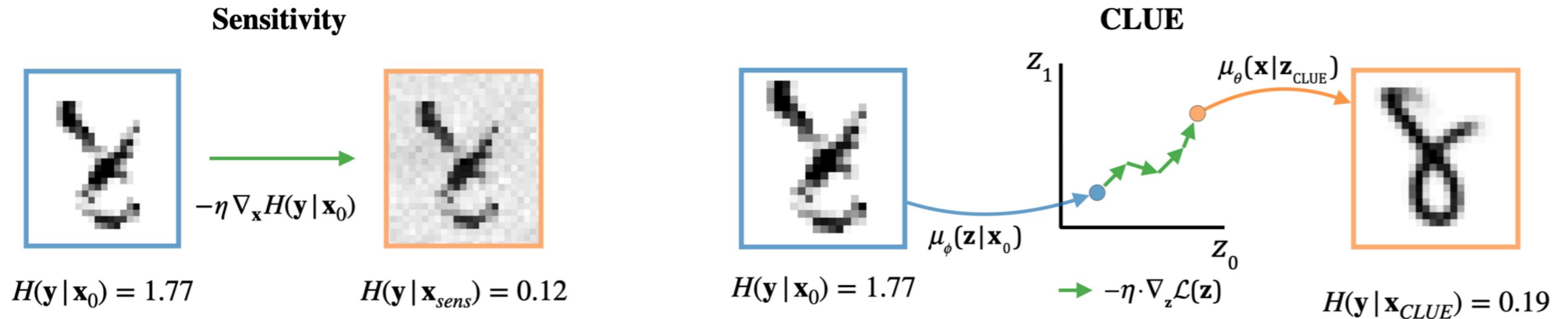
Javier Antoran, **Umang Bhatt**, Tameem Adel, Adrian Weller, and Jose Miguel Hernandez-Lobato

To appear at the International Conference on Learning Representations 2021
<https://arxiv.org/abs/2006.06848>



The
Alan Turing
Institute

Optimization in the Latent Space of a Deep Generative Model



Practitioner-Driven Questions

1. Can we provide methodology for practitioners to evaluate explanations?
2. Can existing explainability tools be used to identify model unfairness?
3. Can we quantify how much **uncertainty** is associated with given explanations?

Practitioner-Driven Questions

1. Can we provide methodology for practitioners to **evaluate** explanations?
2. Can existing explainability tools be used to identify model **unfairness**?
3. Can we quantify how much **uncertainty** is associated with given explanations?

Practical Approaches to Explainable Machine Learning

Thanks for listening!

Umang Bhatt
usb20@cam.ac.uk
[@umangsbhatt](#)