# Sustainability in Explainable AI
Public lecture series Sustainability in Computer Science

Moritz Grosse-Wentrup
Research Group Neuroinformatics
Faculty of Computer Science
University of Vienna

November 27, 2023

# The General Data Protection Regulation (GDPR)

*The data subject should have the right [...] to obtain an **explanation of the decision reached** [...] and to **challenge the decision**.*

# Correctional Offender Management Profiling for Alternative Sanctions

*Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) is a case management and decision support tool developed and owned by Northpointe (now Equivant) used by U.S. courts to assess the likelihood of a defendant becoming a recidivist. COMPAS has been used by the U.S. states of New York, Wisconsin, California, Florida's Broward County, and other jurisdictions.*

```
https://en.wikipedia.org/wiki/COMPAS_(software)
```

# Correctional Offender Management Profiling for Alternative Sanctions

*Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) is a case management and decision support tool developed and owned by Northpointe (now Equivant) used by U.S. courts to assess the likelihood of a defendant becoming a recidivist. COMPAS has been used by the U.S. states of New York, Wisconsin, California, Florida's Broward County, and other jurisdictions.*

`https://en.wikipedia.org/wiki/COMPAS_(software)`



# Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

*Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) is a case management and decision support tool developed and owned by Northpointe (now Equivant) used by U.S. courts to assess the likelihood of a defendant becoming a recidivist. COMPAS has been used by the U.S. states of New York, Wisconsin, California, Florida's Broward County, and other jurisdictions.*

`https://en.wikipedia.org/wiki/COMPAS_(software)`

## Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

*by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica*
*May 23, 2016*

MONKEY CAGE

## A computer program used for bail and sentencing decisions was labeled biased against blacks. It's actually not that clear.

By Sam Corbett-Davies, Emma Pierson, Avi Feller and Sharad Goel
October 17, 2016 at 5:00 a.m. EDT

The Washington Post
*Democracy Dies in Darkness*

# A Machine Learning Primer

- Features $\boldsymbol{X} = [X_1, \ldots, X_d]$; $\boldsymbol{x} \in \mathcal{X}$, e.g., $\mathcal{X} = \mathbb{R}^d$

- Features $\boldsymbol{X} = [X_1, \ldots, X_d]$; $\boldsymbol{x} \in \mathcal{X}$, e.g., $\mathcal{X} = \mathbb{R}^d$
- Measurements $\hat{\boldsymbol{X}} = [\hat{X}_1, \ldots, \hat{X}_d]$; $\boldsymbol{x} \in \mathcal{X}$

# A Machine Learning Primer

- Features $\boldsymbol{X} = [X_1, \ldots, X_d]$; $\boldsymbol{x} \in \mathcal{X}$, e.g., $\mathcal{X} = \mathbb{R}^d$
- Measurements $\hat{\boldsymbol{X}} = [\hat{X}_1, \ldots, \hat{X}_d]$; $\boldsymbol{x} \in \mathcal{X}$
- Labels $Y$ with $y \in \mathcal{Y}$, e.g., $\mathcal{Y} = \{0, 1\}$

# A Machine Learning Primer

- Features $\boldsymbol{X} = [X_1, \ldots, X_d]$; $\boldsymbol{x} \in \mathcal{X}$, e.g., $\mathcal{X} = \mathbb{R}^d$
- Measurements $\hat{\boldsymbol{X}} = [\hat{X}_1, \ldots, \hat{X}_d]$; $\boldsymbol{x} \in \mathcal{X}$
- Labels $Y$ with $y \in \mathcal{Y}$, e.g., $\mathcal{Y} = \{0, 1\}$
- Hypothesis class $\mathcal{H} = \{h_1, \ldots, h_M\}$, $h_m : \mathcal{X} \mapsto \mathcal{Y}$

# A Machine Learning Primer

- Features $\boldsymbol{X} = [X_1, \ldots, X_d]$; $\boldsymbol{x} \in \mathcal{X}$, e.g., $\mathcal{X} = \mathbb{R}^d$
- Measurements $\hat{\boldsymbol{X}} = [\hat{X}_1, \ldots, \hat{X}_d]$; $\boldsymbol{x} \in \mathcal{X}$
- Labels $Y$ with $y \in \mathcal{Y}$, e.g., $\mathcal{Y} = \{0, 1\}$
- Hypothesis class $\mathcal{H} = \{h_1, \ldots, h_M\}$, $h_m : \mathcal{X} \mapsto \mathcal{Y}$
- Loss function $L(h, \boldsymbol{x}, y)$, e.g., zero-one loss.

# A Machine Learning Primer

- Features $\boldsymbol{X} = [X_1, \ldots, X_d]$; $\boldsymbol{x} \in \mathcal{X}$, e.g., $\mathcal{X} = \mathbb{R}^d$
- Measurements $\hat{\boldsymbol{X}} = [\hat{X}_1, \ldots, \hat{X}_d]$; $\boldsymbol{x} \in \mathcal{X}$
- Labels $Y$ with $y \in \mathcal{Y}$, e.g., $\mathcal{Y} = \{0, 1\}$
- Hypothesis class $\mathcal{H} = \{h_1, \ldots, h_M\}$, $h_m : \mathcal{X} \mapsto \mathcal{Y}$
- Loss function $L(h, \boldsymbol{x}, y)$, e.g., zero-one loss.
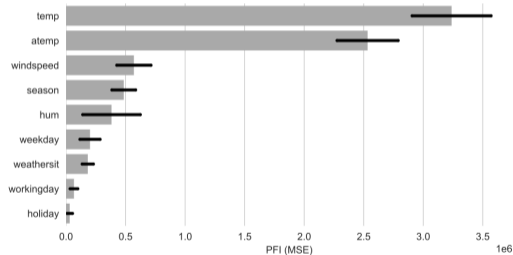- Training set $\mathcal{S} = \{(\boldsymbol{x}^i, y^i)\}_{i=1}^{N}$ sampled iid from $P(Y, \boldsymbol{X})$

- Features $\boldsymbol{X} = [X_1, \ldots, X_d]$; $\boldsymbol{x} \in \mathcal{X}$, e.g., $\mathcal{X} = \mathbb{R}^d$
- Measurements $\hat{\boldsymbol{X}} = [\hat{X}_1, \ldots, \hat{X}_d]$; $\boldsymbol{x} \in \mathcal{X}$
- Labels $Y$ with $y \in \mathcal{Y}$, e.g., $\mathcal{Y} = \{0, 1\}$
- Hypothesis class $\mathcal{H} = \{h_1, \ldots, h_M\}$, $h_m : \mathcal{X} \mapsto \mathcal{Y}$
- Loss function $L(h, \boldsymbol{x}, y)$, e.g., zero-one loss.
- Training set $\mathcal{S} = \{(\boldsymbol{x}^i, y^i)\}_{i=1}^{N}$ sampled iid from $P(Y, \boldsymbol{X})$
- Learning algorithm $A(\mathcal{S}, \mathcal{H}, l)$

# A Machine Learning Primer

- Features $\boldsymbol{X} = [X_1, \ldots, X_d]$; $\boldsymbol{x} \in \mathcal{X}$, e.g., $\mathcal{X} = \mathbb{R}^d$
- Measurements $\hat{\boldsymbol{X}} = [\hat{X}_1, \ldots, \hat{X}_d]$; $\boldsymbol{x} \in \mathcal{X}$
- Labels $Y$ with $y \in \mathcal{Y}$, e.g., $\mathcal{Y} = \{0, 1\}$
- Hypothesis class $\mathcal{H} = \{h_1, \ldots, h_M\}$, $h_m : \mathcal{X} \mapsto \mathcal{Y}$
- Loss function $L(h, \boldsymbol{x}, y)$, e.g., zero-one loss.
- Training set $\mathcal{S} = \{(\boldsymbol{x}^i, y^i)\}_{i=1}^N$ sampled iid from $P(Y, \boldsymbol{X})$
- Learning algorithm $A(\mathcal{S}, \mathcal{H}, l)$
- Bayes-optimal hypothesis $h^* \in \mathcal{H}$ with (empirical) risk $L_{(\mathcal{S})}(h^*)$

# A Machine Learning Primer

- Features $\boldsymbol{X} = [X_1, \ldots, X_d]$; $\boldsymbol{x} \in \mathcal{X}$, e.g., $\mathcal{X} = \mathbb{R}^d$
- Measurements $\hat{\boldsymbol{X}} = [\hat{X}_1, \ldots, \hat{X}_d]$; $\boldsymbol{x} \in \mathcal{X}$
- Labels $Y$ with $y \in \mathcal{Y}$, e.g., $\mathcal{Y} = \{0, 1\}$
- Hypothesis class $\mathcal{H} = \{h_1, \ldots, h_M\}$, $h_m : \mathcal{X} \mapsto \mathcal{Y}$
- Loss function $L(h, \boldsymbol{x}, y)$, e.g., zero-one loss.
- Training set $\mathcal{S} = \{(\boldsymbol{x}^i, y^i)\}_{i=1}^{N}$ sampled iid from $P(Y, \boldsymbol{X})$
- Learning algorithm $A(\mathcal{S}, \mathcal{H}, l)$
- Bayes-optimal hypothesis $h^* \in \mathcal{H}$ with (empirical) risk $L_{(\mathcal{S})}(h^*)$
- Predictions $\hat{Y} = h^*(\hat{\boldsymbol{X}})$

# A Machine Learning Primer

- Features $\boldsymbol{X} = [X_1, \ldots, X_d]$; $\boldsymbol{x} \in \mathcal{X}$, e.g., $\mathcal{X} = \mathbb{R}^d$
- Measurements $\hat{\boldsymbol{X}} = [\hat{X}_1, \ldots, \hat{X}_d]$; $\boldsymbol{x} \in \mathcal{X}$
- Labels $Y$ with $y \in \mathcal{Y}$, e.g., $\mathcal{Y} = \{0, 1\}$
- Hypothesis class $\mathcal{H} = \{h_1, \ldots, h_M\}$, $h_m : \mathcal{X} \mapsto \mathcal{Y}$
- Loss function $L(h, \boldsymbol{x}, y)$, e.g., zero-one loss.
- Training set $\mathcal{S} = \{(\boldsymbol{x}^i, y^i)\}_{i=1}^N$ sampled iid from $P(Y, \boldsymbol{X})$
- Learning algorithm $A(\mathcal{S}, \mathcal{H}, l)$
- Bayes-optimal hypothesis $h^* \in \mathcal{H}$ with (empirical) risk $L_{(\mathcal{S})}(h^*)$
- Predictions $\hat{Y} = h^*(\hat{\boldsymbol{X}})$
- Interventions $\text{do}\{\boldsymbol{X} = ()\}$ (on features) and $\text{do}\{\hat{\boldsymbol{X}} = ()\}$ (on measurements)

# Example: The Permutation Feature Importance Score

Idea: Assess how removing each individual features affects model performance.

$$\mathsf{PFI}(X_i) = L(h^*, \{\boldsymbol{X} \setminus X_i, \tilde{X}_i\}, Y) - L(h^*, \boldsymbol{X}, Y)$$

with $\tilde{X}_i \sim P(X_i)$ and $\tilde{X} \perp\!\!\!\perp \{\boldsymbol{X} \setminus X_i, Y\}$.

Breiman, Leo. "Random forests." Machine learning 45 (2001): 5-32.

# xAI/IML is a Causal Problem

# Structural Causal Models (SCMs)

J. Pearl, *Causality*. 2000.

## Structural Causal Models (SCMs)

For a given set of variables $\mathcal{X} = \{X_i\}_{i=1}^{N}$ a *structural causal model* (SCM) is defined by

$$X_i = f_i(\text{pa}_i, \epsilon_i)$$

with $\{\epsilon_i\}_{i=1}^{N}$ exogenous noise terms and the *parents* $\text{pa}_i \subset \mathcal{X} \backslash X_i$ chosen such that the corresponding graph contains no cycles.
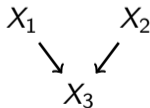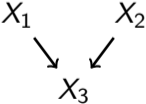
J. Pearl, *Causality*. 2000.

## Structural Causal Models (SCMs)

For a given set of variables $\mathcal{X} = \{X_i\}_{i=1}^N$ a *structural causal model* (SCM) is defined by

$$X_i = f_i(\mathrm{pa}_i, \epsilon_i)$$

with $\{\epsilon_i\}_{i=1}^N$ exogenous noise terms and the *parents* $\mathrm{pa}_i \subset \mathcal{X} \backslash X_i$ chosen such that the corresponding graph contains no cycles.

Example:

J. Pearl, *Causality*. 2000.

# Structural Causal Models (SCMs)

For a given set of variables $\mathcal{X} = \{X_i\}_{i=1}^N$ a *structural causal model* (SCM) is defined by

$$X_i = f_i(\mathrm{pa}_i, \epsilon_i)$$

with $\{\epsilon_i\}_{i=1}^N$ exogenous noise terms and the *parents* $\mathrm{pa}_i \subset \mathcal{X} \backslash X_i$ chosen such that the corresponding graph contains no cycles.
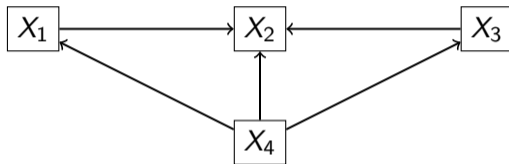
Example:

SCM

$$X_1 = \epsilon_1$$
$$X_2 = \epsilon_2$$
$$X_3 = X_1 \cdot X_2 + \epsilon_3$$

J. Pearl, *Causality*. 2000.

# Structural Causal Models (SCMs)

For a given set of variables $\mathcal{X} = \{X_i\}_{i=1}^N$ a *structural causal model* (SCM) is defined by

$$X_i = f_i(\mathrm{pa}_i, \epsilon_i)$$

with $\{\epsilon_i\}_{i=1}^N$ exogenous noise terms and the *parents* $\mathrm{pa}_i \subset \mathcal{X} \backslash X_i$ chosen such that the corresponding graph contains no cycles.

Example:

SCM            DAG

$$X_1 = \epsilon_1$$
$$X_2 = \epsilon_2$$
$$X_3 = X_1 \cdot X_2 + \epsilon_3$$

$X_1 \qquad X_2$

$X_3$

J. Pearl, *Causality*. 2000.

# Structural Causal Models (SCMs)

For a given set of variables $\mathcal{X} = \{X_i\}_{i=1}^{N}$ a *structural causal model* (SCM) is defined by

$$X_i = f_i(\mathrm{pa}_i, \epsilon_i)$$

with $\{\epsilon_i\}_{i=1}^{N}$ exogenous noise terms and the *parents* $\mathrm{pa}_i \subset \mathcal{X} \backslash X_i$ chosen such that the corresponding graph contains no cycles.

Example:

| SCM | DAG | Data |
|-----|-----|------|
| $X_1 = \epsilon_1$ | $X_1 \qquad X_2$ | $\epsilon \sim p(\epsilon)$ |
| $X_2 = \epsilon_2$ | $\searrow \swarrow$ | $x_i = f_i(\mathrm{pa}_i, \epsilon_i)$ |
| $X_3 = X_1 \cdot X_2 + \epsilon_3$ | $X_3$ | $\boldsymbol{x} \sim p(\boldsymbol{x})$ |

J. Pearl, *Causality*. 2000.

# Structural Causal Models (SCMs)

For a given set of variables $\mathcal{X} = \{X_i\}_{i=1}^N$ a *structural causal model* (SCM) is defined by

$$X_i = f_i(\text{pa}_i, \epsilon_i)$$

with $\{\epsilon_i\}_{i=1}^N$ exogenous noise terms and the *parents* $\text{pa}_i \subset \mathcal{X} \backslash X_i$ chosen such that the corresponding graph contains no cycles.

Example:

| SCM | DAG | Data |
|-----|-----|------|
| $X_1 = \epsilon_1$ | $X_1 \qquad X_2$ | $\epsilon \sim p(\epsilon)$ |
| $X_2 = \epsilon_2$ | $\searrow \swarrow$ | $x_i = f_i(\text{pa}_i, \epsilon_i)$ |
| $X_3 = X_1 \cdot X_2 + \epsilon_3$ | $X_3$ | $\mathbf{x} \sim p(\mathbf{x})$ |

Def.: $X_i$ is a cause of $X_j$, iff there exist values of $X_i$ and $X_j$ such that $p(x_j | \text{do}\{x_i\}) \neq p(x_j)$.

J. Pearl, *Causality*. 2000.

Causal factorization:

Causal factorization:



$$P(\boldsymbol{X}) = P(X_2|X_1, X_3, X_4)P(X_1|X_4)P(X_3|X_4)P(X_4)$$

(J. Pearl, *Causality*. 2000.)

Causal factorization:



$$P(\boldsymbol{X}) = P(X_2|X_1, X_3, X_4)P(X_1|X_4)P(X_3|X_4)P(X_4)$$

Interventions are represented by the do-operator, e.g.,

$$P\left(\boldsymbol{X}|\text{do}(X_1 = x_1)\right) = P(X_2|X_1 = x_1, X_3, X_4)P(X_3|X_4)P(X_4).$$

(J. Pearl, *Causality*. 2000.)

# Causal inference

How can we infer the structure of the DAG from the data it generates? We need concepts and assumptions that link the structural with the observational world:

## Causal inference

How can we infer the structure of the DAG from the data it generates? We need concepts and assumptions that link the structural with the observational world:

Causal Markov Condition (CMC): Every node in $\mathcal{X}$ is conditionally independent of its nondescendents given its parents.

# Causal inference

How can we infer the structure of the DAG from the data it generates? We need concepts and assumptions that link the structural with the observational world:

Causal Markov Condition (CMC): Every node in $\mathcal{X}$ is conditionally independent of its nondescendents given its parents.

Faithfulness: There are no further independence relations among the nodes in $\mathcal{X}$ beyond those implied by d-separation.

# Causal inference

How can we infer the structure of the DAG from the data it generates? We need concepts and assumptions that link the structural with the observational world:

Causal Markov Condition (CMC): Every node in $\mathcal{X}$ is conditionally independent of its nondescendents given its parents.

Faithfulness: There are no further independence relations among the nodes in $\mathcal{X}$ beyond those implied by d-separation.

d-separation: Let $A, B, C$ non-intersecting subsets of $\mathcal{X}$. $A$ and $B$ are d-separated given $C$ iff

- for all nodes on the path where the arrows meet head-to-tail ($\rightarrow . \rightarrow$) or tail-to-tail ($\leftarrow . \rightarrow$) the node is in $C$,
- for all nodes where the arrows meet head-to-head ($\rightarrow . \leftarrow$) neither the node or any of its descendants are in $C$.

# Causal inference

How can we infer the structure of the DAG from the data it generates? We need concepts and assumptions that link the structural with the observational world:

Causal Markov Condition (CMC): Every node in $\mathcal{X}$ is conditionally independent of its nondescendents given its parents.

Faithfulness: There are no further independence relations among the nodes in $\mathcal{X}$ beyond those implied by d-separation.

d-separation: Let $A, B, C$ non-intersecting subsets of $\mathcal{X}$. $A$ and $B$ are d-separated given $C$ iff

- for all nodes on the path where the arrows meet head-to-tail ($\rightarrow . \rightarrow$) or tail-to-tail ($\leftarrow . \rightarrow$) the node is in $C$,
- for all nodes where the arrows meet head-to-head ($\rightarrow . \leftarrow$) neither the node or any of its descendants are in $C$.

Assuming the CMC and faithfulness, $\text{dSep}(A, B | C) \Leftrightarrow A \perp\!\!\!\perp B | C$.

The chain
$$X_1 \to X_2 \to X_3$$

## Example

The chain

$X_1 \rightarrow X_2 \rightarrow X_3$

The fork

$X_1 \leftarrow X_2 \rightarrow X_3$

The chain $\quad$ The fork $\quad$ The collider

$X_1 \rightarrow X_2 \rightarrow X_3 \qquad X_1 \leftarrow X_2 \rightarrow X_3 \qquad X_1 \rightarrow X_2 \leftarrow X_3$

The chain | The fork | The collider
$X_1 \to X_2 \to X_3$ $\quad X_1 \leftarrow X_2 \to X_3$ $\quad X_1 \to X_2 \leftarrow X_3$
$X_1 \not\!\perp X_3$

The chain

$X_1 \to X_2 \to X_3$

$X_1 \not\perp\!\!\!\perp X_3$

$X_1 \perp\!\!\!\perp X_3 | X_2$

The fork

$X_1 \leftarrow X_2 \to X_3$

The collider

$X_1 \to X_2 \leftarrow X_3$

The chain

$X_1 \rightarrow X_2 \rightarrow X_3$

$X_1 \not\perp\!\!\!\perp X_3$

$X_1 \perp\!\!\!\perp X_3 | X_2$

The fork

$X_1 \leftarrow X_2 \rightarrow X_3$

$X_1 \not\perp\!\!\!\perp X_3$

The collider

$X_1 \rightarrow X_2 \leftarrow X_3$

# Example

|  The chain | The fork | The collider |
|:---:|:---:|:---:|
| $X_1 \rightarrow X_2 \rightarrow X_3$ | $X_1 \leftarrow X_2 \rightarrow X_3$ | $X_1 \rightarrow X_2 \leftarrow X_3$ |
| $X_1 \not\perp\!\!\!\perp X_3$ | $X_1 \not\perp\!\!\!\perp X_3$ | |
| $X_1 \perp\!\!\!\perp X_3 \mid X_2$ | $X_1 \perp\!\!\!\perp X_3 \mid X_2$ | |

# Example

|  The chain | The fork | The collider |
|:---:|:---:|:---:|
| $X_1 \rightarrow X_2 \rightarrow X_3$ | $X_1 \leftarrow X_2 \rightarrow X_3$ | $X_1 \rightarrow X_2 \leftarrow X_3$ |
| $X_1 \not\perp\!\!\!\perp X_3$ | $X_1 \not\perp\!\!\!\perp X_3$ | $X_1 \perp\!\!\!\perp X_3$ |
| $X_1 \perp\!\!\!\perp X_3 \mid X_2$ | $X_1 \perp\!\!\!\perp X_3 \mid X_2$ | |

# Example

|  The chain | The fork | The collider |
| :---: | :---: | :---: |
| $X_1 \rightarrow X_2 \rightarrow X_3$ | $X_1 \leftarrow X_2 \rightarrow X_3$ | $X_1 \rightarrow X_2 \leftarrow X_3$ |
| $X_1 \not\perp\!\!\!\perp X_3$ | $X_1 \not\perp\!\!\!\perp X_3$ | $X_1 \perp\!\!\!\perp X_3$ |
| $X_1 \perp\!\!\!\perp X_3 \mid X_2$ | $X_1 \perp\!\!\!\perp X_3 \mid X_2$ | $X_1 \not\perp\!\!\!\perp X_3 \mid X_2$ |

# xAI/IML is a Causal Problem

# xAI/IML is a Causal Problem

# xAI/IML is a Causal Problem

# xAI/IML is a Causal Problem



$Y \perp\!\!\!\perp X_2 | X_3 \Rightarrow \hat{Y} \perp\!\!\!\perp \hat{X}_2$ (if $h$ is optimal) but $\hat{Y} \not\!\perp\!\!\!\perp X_2$

# xAI/IML is a Causal Problem



World                                    Model

$Y \perp\!\!\!\perp X_2 | X_3 \Rightarrow \hat{Y} \perp\!\!\!\perp \hat{X}_2$ (if $h$ is optimal) but $\hat{Y} \not\perp\!\!\!\perp X_2$

$P(\hat{Y} | \text{do}\{\hat{X}_1 = \hat{x}_1 + \delta_{x_1}\}) \neq P(Y | \text{do}\{X_1 = x_1 + \delta_{x_1}\})$

# xAI/IML is a Causal Problem



World               Model

$Y \perp\!\!\!\perp X_2 | X_3 \Rightarrow \hat{Y} \perp\!\!\!\perp \hat{X}_2$ (if $h$ is optimal) but $\hat{Y} \not\!\perp\!\!\!\perp X_2$

$P(\hat{Y}|\text{do}\{\hat{X}_1 = \hat{x}_1 + \delta_{x_1}\}) \neq P(Y|\text{do}\{X_1 = x_1 + \delta_{x_1}\})$

$P(\hat{Y}|\text{do}\{\hat{X}_3 = \hat{x}_3\}) \neq P(Y|\text{do}\{X_3 = x_3\})$

What object is explained?

What object is explained?
- The prediction $\hat{Y}$?

What object is explained?

- The prediction $\hat{Y}$?
- The target variable $Y$?

What object is explained?

- The prediction $\hat{Y}$?
- The target variable $Y$?
- Their relationship $R$ (as measured by the loss)?

What object is explained?

- The prediction $\hat{Y}$?
- The target variable $Y$?
- Their relationship $R$ (as measured by the loss)?

On what level is the object explained?

What object is explained?

- The prediction $\hat{Y}$?
- The target variable $Y$?
- Their relationship $R$ (as measured by the loss)?

On what level is the object explained?

- Associations $X = x$?

What object is explained?

- The prediction $\hat{Y}$?
- The target variable $Y$?
- Their relationship $R$ (as measured by the loss)?

On what level is the object explained?

- Associations $X = x$?
- Model-level interventions $\text{do}(\hat{X} = x)$?

# A Taxonomy of xAI/IML

What object is explained?

- The prediction $\hat{Y}$?
- The target variable $Y$?
- Their relationship $R$ (as measured by the loss)?

On what level is the object explained?

- Associations $X = x$?
- Model-level interventions $do(\hat{X} = x)$?
- World-level interventions $do(X = x)$?

# Nine Perspectives on Model and Data

# Contribution I: Improvement-Focused Causal Recourse (ICR)

*König, G.,* *Freiesleben, T., & Grosse-Wentrup, M. (2023).*
*Improvement-focused causal recourse (ICR).*
*In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 37, No. 10, pp. 11847-11855).*

# Recourse Setting

# Recourse Setting



- suppose a ML model rejects your request (job application, loan application, hospital admission, ....)

- **recourse recommendations** tell you what to do to get accepted

# Recourse Setting



- suppose a ML model rejects your request (job application, loan application, hospital admission, ....)

- **recourse recommendations** tell you what to do to get accepted

- existing methods: counterfactual explanations (CE) [Wachter et al.], causal recourse (CR) [Karimi et al.]

# Problem     Terminology

# Problem

# Terminology

- existing methods (CE, CR) may suggest to **game** the predictor

**gaming**: tricking the predictor into falsely believing that one is qualified

$\hat{y} = 1$ but $y = 0$

# Problem

# Terminology

- existing methods (CE, CR) may suggest to **game** the predictor

- recourse should guide towards both **acceptance** *and* **improvement**

**gaming***:* tricking the predictor into falsely believing that one is qualified

$\hat{y} = 1$ but $y = 0$

**acceptance:** reverting the model's decision

$\hat{y} = 1$

**improvement***:* reverting the underlying real-world state

$y = 1$

# Illustrative Example

**Goal:** predict CoViD risk to decide whether you are allowed to enter a hospital (without testing for CoViD)



https://riskcalc.org/COVID19/

# Illustrative Example



**Counterfactual Explanations (CE)**
**[Wachter et al.]**

| | |
|---|---|
| *viewpoint* | (diagram: vacc. → $\hat{Y}$, sympt. → $\hat{Y}$, CoViD $Y$) |
| *idea* | What is the minimal $do(\underline{X}_I = x')$ such that $\hat{y} = 1$? |
| *exemplary explanation* | *"If you plug lower symptom values into the model, the prediction is favorable."* |

# Illustrative Example



**Counterfactual Explanations (CE)**
**[Wachter et al.]**

| | |
|---|---|
| *viewpoint* | vacc. → $\hat{Y}$ ← sympt. (CoViD $Y$) |
| *idea* | What is the minimal $do(\underline{X}_I = x')$ such that $\hat{y} = 1$? |
| *exemplary explanation* | *"If you plug lower symptom values into the model, the prediction is favorable."* |
| *acceptance* ($\hat{y} = 1$) | not guaranteed |

# Illustrative Example



|  | Counterfactual Explanations (CE) [Wachter et al.] | Causal Recourse (CR) [Karimi et al.] |
|---|---|---|
| *viewpoint* | | |
| *idea* | What is the minimal $do(\underline{X}_I = x')$ such that $\hat{y} = 1$? | What is the most cost-efficient $do(X_I = x')$ such that $\hat{y} = 1$? |
| *exemplary explanation* | *"If you plug lower symptom values into the model, the prediction is favorable."* | *"If you treat your symptoms (take cough syrup), the model will accept you."* |
| *acceptance* $(\hat{y} = 1)$ | not guaranteed | yes |

# Illustrative Example



|  | Counterfactual Explanations (CE) [Wachter et al.] | Causal Recourse (CR) [Karimi et al.] |
|---|---|---|
| *viewpoint* | | |
| *idea* | What is the minimal $do(\underline{X}_I = x')$ such that $\hat{y} = 1$? | What is the most cost-efficient $do(X_I = x')$ such that $\hat{y} = 1$? |
| *exemplary explanation* | *"If you plug lower symptom values into the model, the prediction is favorable."* | *"If you treat your symptoms (take cough syrup), the model will accept you."* |
| *acceptance* ($\hat{y} = 1$) | not guaranteed | yes |
| *improvement* ($y = 1$) | not guaranteed | not guaranteed |

# Illustrative Example

| | Counterfactual Explanations (CE) [Wachter et al.] | Causal Recourse (CR) [Karimi et al.] | Improvement-Focused CR (ICR) [Koenig et al.] |
|---|---|---|---|
| *viewpoint* |  |  |  |
| *idea* | What is the minimal $do(\underline{X}_I = x')$ such that $\hat{y} = 1$? | What is the most cost-efficient $do(X_I = x')$ such that $\hat{y} = 1$? | What is the most cost-efficient $do(X_I = x')$ such that $y = 1$? |
| *exemplary explanation* | *"If you plug lower symptom values into the model, the prediction is favorable."* | *"If you treat your symptoms (take cough syrup), the model will accept you."* | *"If you get vaccinated, you will decrease your CoViD risk and thus be accepted."* |
| *acceptance* ($\hat{y} = 1$) | not guaranteed | yes | guaranteed |
| *improvement* ($y = 1$) | not guaranteed | not guaranteed | yes |

# The Methods & The 9 Perspectives

|  | $\hat{Y}$ | $R$ | $Y$ |
|---|---|---|---|
| $do(\underline{X} = x')$ | CEs |  |  |
| $X = x'$ |  |  |  |
| $do(X = x')$ | CR |  | ICR |

# ICR: Optimization Problem

For a target improvement probability $\bar{\gamma}$, some cost function $c$ and pre-recourse observation $x^{pre}$, the ICR action $a := do(X_{I_a} := \theta_{I_a})$ is found by optimizing

$$\text{argmin}_a \quad c(a; x^{pre}) \quad s.t. \quad \gamma(a; x^{pre}) \geq \bar{\gamma};$$

with $\gamma(a, x^{pre})$ being the improvement probability for action $a$ and an individual with characteristics $x^{pre}$.

# How to Define Improvement Confidence $\gamma$?

**Goal:** For an action $do(a)$, estimate probability of improvement while taking as many features as possible into account.

# How to Define Improvement Confidence $\gamma$?

**Goal:** For an action $do(a)$, estimate probability of improvement while taking as many features as possible into account.

|  | **_individualized_**<br>**(structural causal model known):** |
| --- | --- |
| _estimation based on_ | structural counterfactuals |
| _precision_ | takes all features into account |
| _definition_ | $\gamma^{ind} := P(Y^{post} = 1 \mid x^{pre}, do(a))$ |

# How to Define Improvement Confidence $\gamma$?

**Goal:** For an action $do(a)$, estimate probability of improvement while taking as many features as possible into account.

|  | *individualized*<br>**(structural causal model known):** | *subpopulation-based*<br>**(only causal graph known):** |
|---|---|---|
| *estimation based on* | structural counterfactuals | conditional average treatment effect |
| *precision* | takes all features into account | only takes features that are not affected by the action into account |
| *definition* | $\gamma^{ind} := P(Y^{post} = 1 \mid x^{pre}, do(a))$ | $\gamma^{sub} := P(Y^{post} = 1 \mid x_{G_a}^{pre}, do(a))$ |

# How to Achieve Acceptance ($\hat{y} = 1$)?

# How to Achieve Acceptance ($\hat{y} = 1$)?

**Intuition:** Classifiers remain accurate under ICR actions → acceptance ensures improvement.

# How to Achieve Acceptance ($\hat{y} = 1$)?

**Intuition:** Classifiers remain accurate under ICR actions → acceptance ensures improvement.

Why?

# How to Achieve Acceptance ($\hat{y} = 1$)?

**Intuition:** Classifiers remain accurate under ICR actions → acceptance ensures improvement.

Why?

- ICR only recommends interventions on causes

# How to Achieve Acceptance ($\hat{y} = 1$)?

**Intuition:** Classifiers remain accurate under ICR actions → acceptance ensures improvement.

Why?

- ICR only recommends interventions on causes

- intervening on causes does not affect $P(Y \mid X)$

# How to Achieve Acceptance ($\hat{y} = 1$)?

**Intuition:** Classifiers remain accurate under ICR actions → acceptance ensures improvement.

Why?

- ICR only recommends interventions on causes

- intervening on causes does not affect $P(Y \mid X)$

→ the classifier is *stable* w.r.t. ICR actions

# How to Achieve Acceptance ($\hat{y} = 1$)?

**Intuition:** Classifiers remain accurate under ICR actions → acceptance ensures improvement.

Why?

- ICR only recommends interventions on causes

- intervening on causes does not affect $P(Y \mid X)$

→ the classifier is *stable* w.r.t. ICR actions



data generating process

Gunnar König

Alex Markham