



Interpretable Off-Policy Learning via Hyperbox Search

International Conference on Machine Learning (ICML) 2022

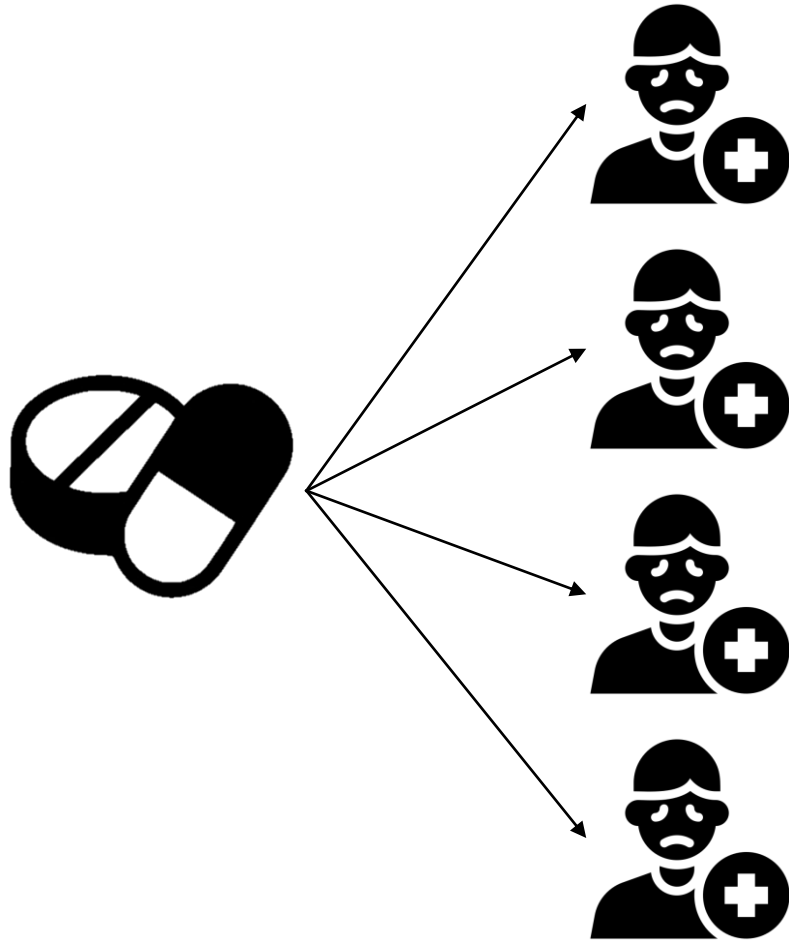
Baltimore, Maryland, USA

Daniel Tschernutter

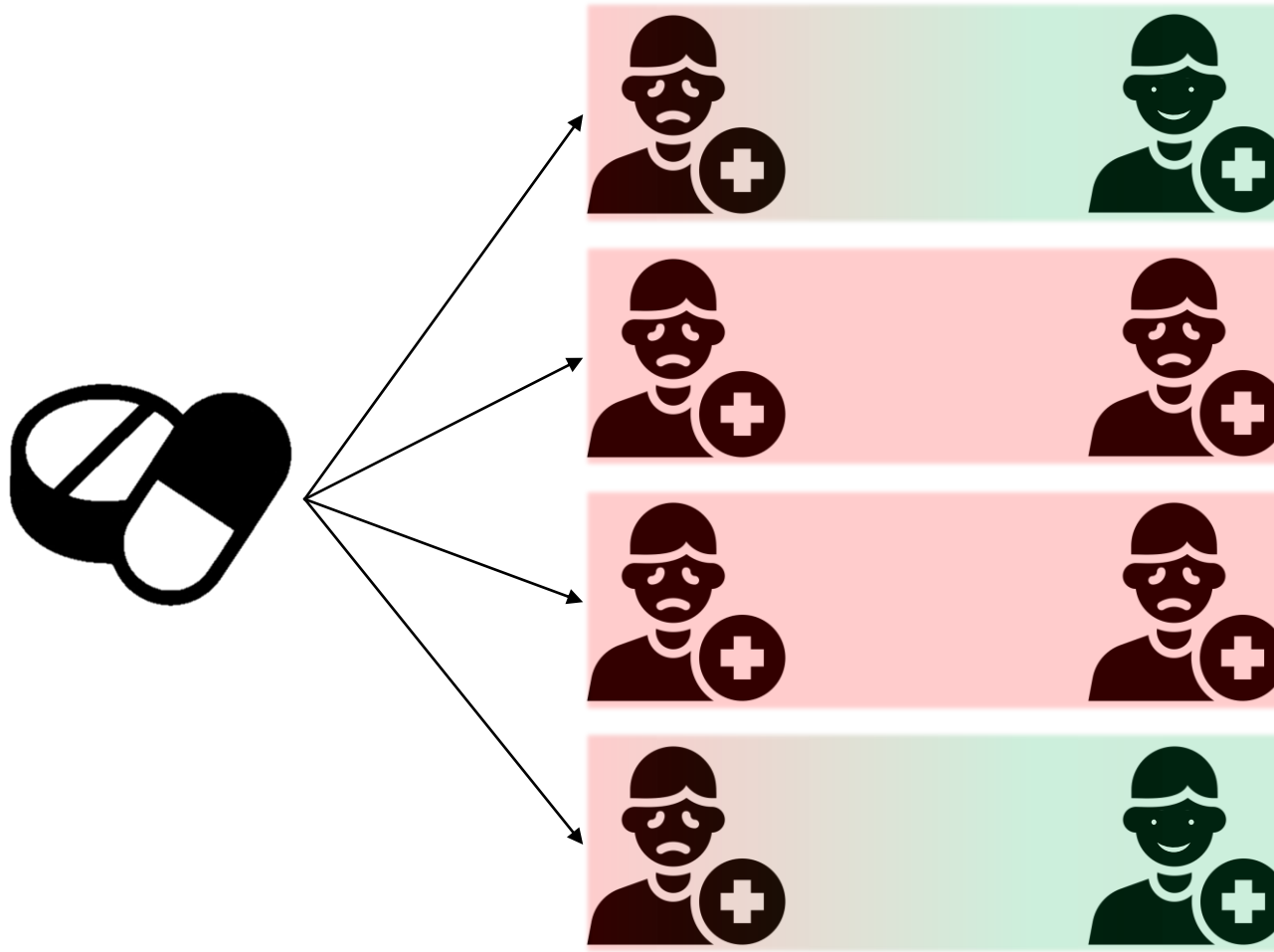
Tobias Hatt

Stefan Feuerriegel

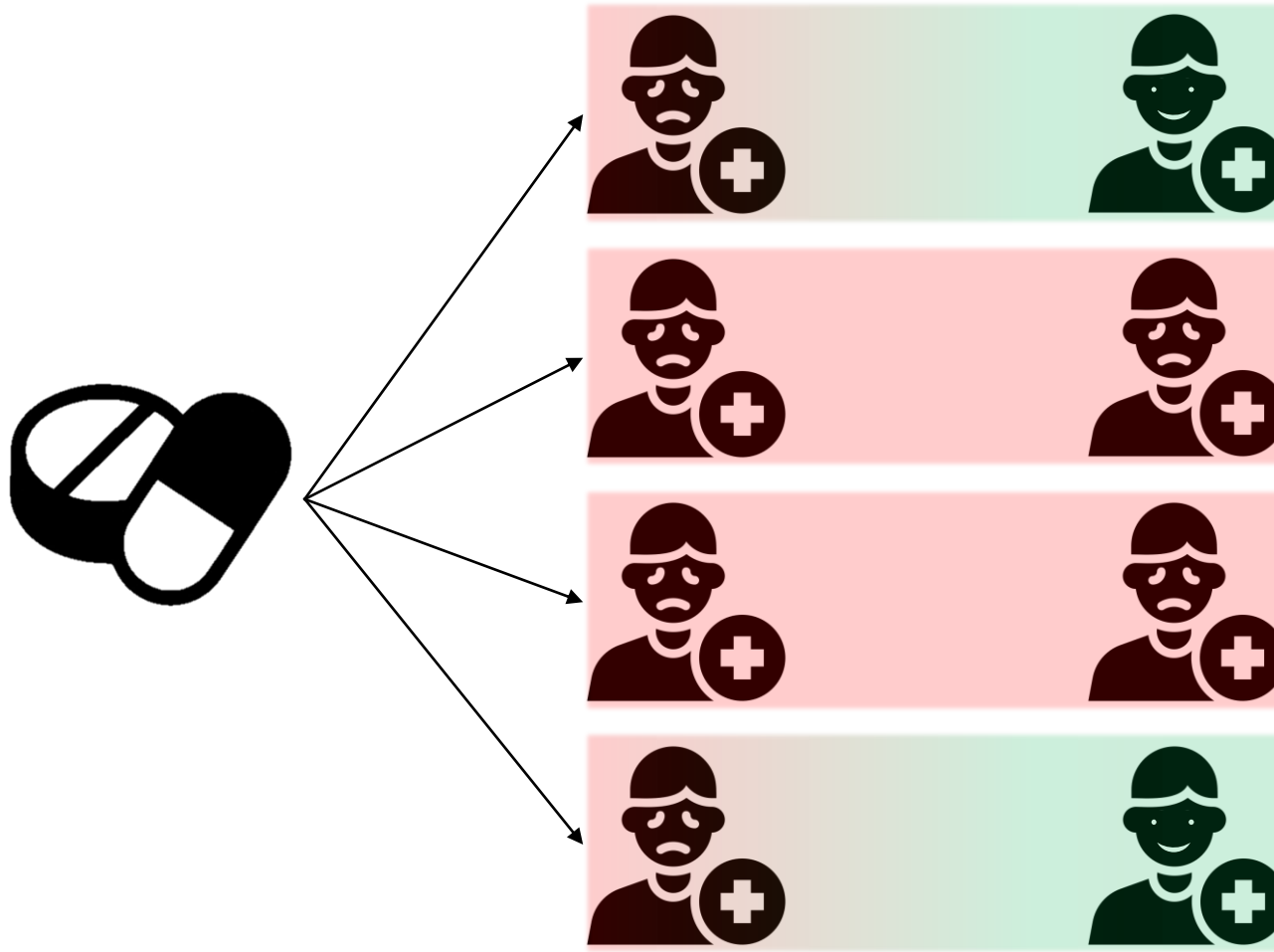
Why Off-Policy Learning?



Why Off-Policy Learning?

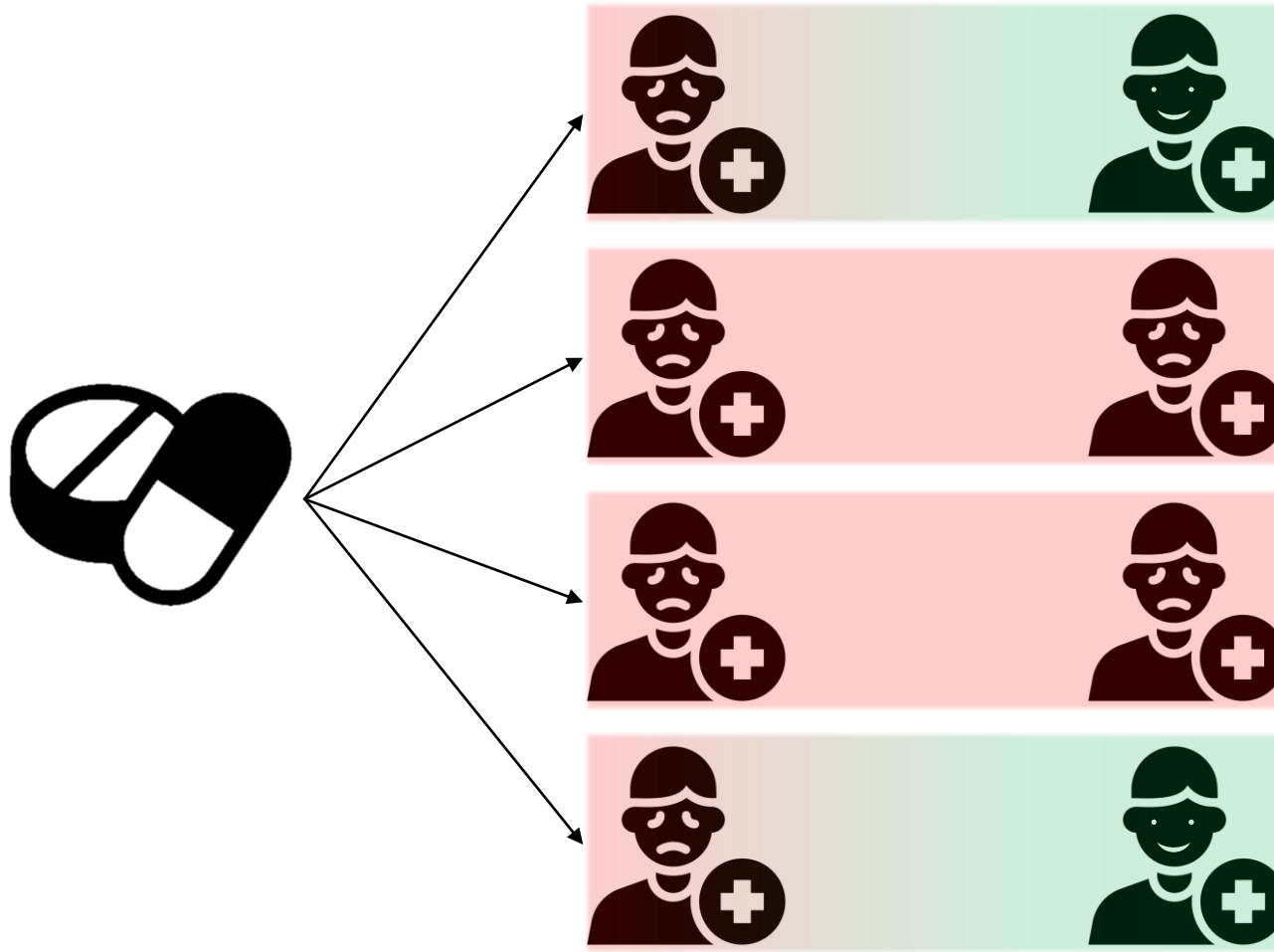


Why Off-Policy Learning?



- Aim is to **select treatments** that are effective for **individual patients**

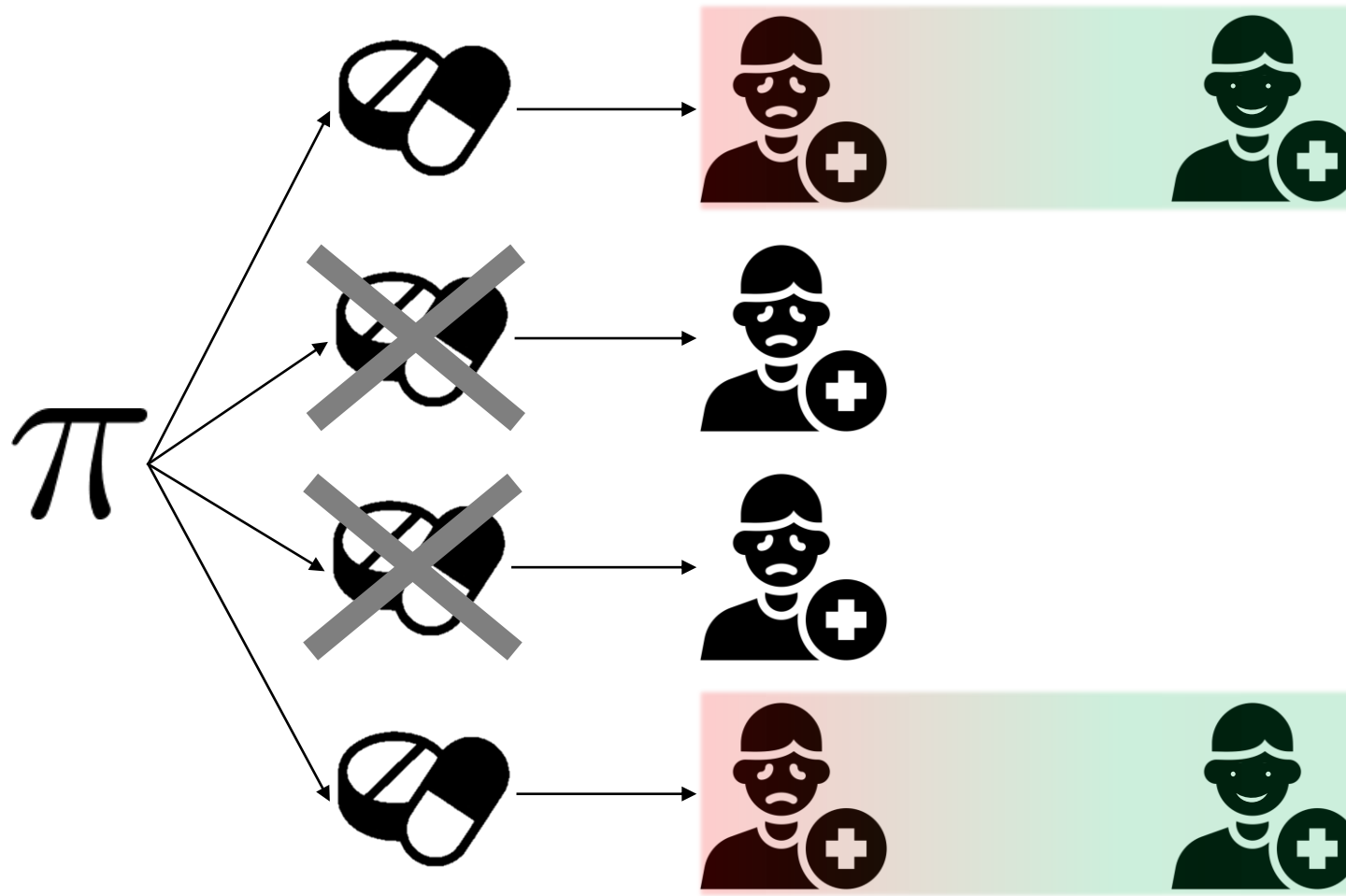
Why Off-Policy Learning?



- Aim is to **select treatments** that are effective for **individual patients**
- Personalized decision-making formalized via so called **policies**

$$\pi : \mathcal{X} \subseteq \mathbb{R}^d \rightarrow \{-1, 1\}$$

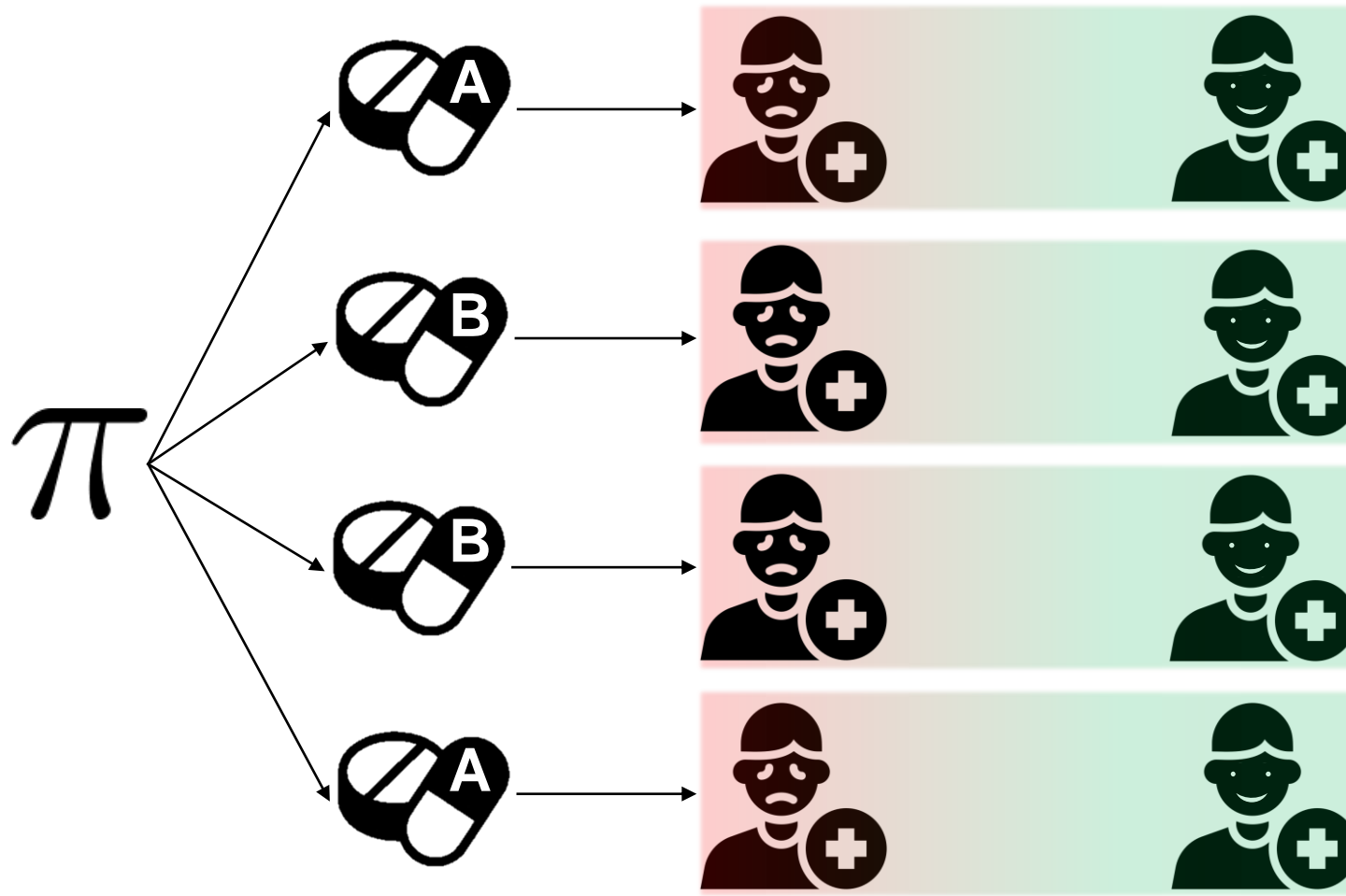
Why Off-Policy Learning?



- Aim is to **select treatments** that are effective for **individual patients**
- Personalized decision-making formalized via so called **policies**

$$\pi : \mathcal{X} \subseteq \mathbb{R}^d \rightarrow \{-1, 1\}$$

Why Off-Policy Learning?



- Aim is to **select treatments** that are effective for **individual patients**
- Personalized decision-making formalized via so called **policies**

$$\pi : \mathcal{X} \subseteq \mathbb{R}^d \rightarrow \{-1, 1\}$$

Off-Policy Learning

- 1 Formalize objective via policy value

$$V(\pi) = \mathbb{E} [Y(1) I(\pi(X) = 1) + Y(-1) I(\pi(X) = -1)]$$

Off-Policy Learning

1 Formalize objective via policy value

$$V(\pi) = \mathbb{E} [Y(1) I(\pi(X) = 1) + Y(-1) I(\pi(X) = -1)]$$

Diagram illustrating the formalization of the objective via policy value:

- $V(\pi)$ is labeled as the **Policy value**.
- $Y(1)$ and $Y(-1)$ are labeled as **Potential outcomes**.
- $I(\pi(X) = 1)$ and $I(\pi(X) = -1)$ are labeled as the **Indicator function for treatment decisions**.

Off-Policy Learning

1 Formalize objective via policy value



2 Estimate policy value from data

$$V(\pi) = \mathbb{E} [Y(1) I(\pi(X) = 1) + Y(-1) I(\pi(X) = -1)]$$

$$\mathcal{J}(\pi) = \mathbb{E} [\psi I(T \neq \pi(X))]$$

$$\psi^{\text{DM}} = T(\mu_{-1}(X) - \mu_1(X))$$

$$\psi^{\text{IPS}} = \frac{-Y}{e_T(X)}$$

$$\psi^{\text{DR}} = \psi^{\text{DM}} + \psi^{\text{IPS}} + \frac{\mu_T(X)}{e_T(X)}$$

$$\mu_t(x) = \mathbb{E} [Y(t) \mid X = x]$$

$$e_t(x) = \mathbb{P}(T = t \mid X = x)$$

Standard methods:

DM – direct method

IPS – inverse propensity score

DR – doubly robust method

Off-Policy Learning

1 Formalize objective via policy value



2 Estimate policy value from data

$$V(\pi) = \mathbb{E} [Y(1) I(\pi(X) = 1) + Y(-1) I(\pi(X) = -1)]$$

$$\mathcal{J}(\pi) = \mathbb{E} [\psi I(T \neq \pi(X))]$$

$$\psi^{\text{DM}} = T(\mu_{-1}(X) - \mu_1(X))$$

$$\psi^{\text{IPS}} = \frac{-Y}{e_T(X)}$$

$$\psi^{\text{DR}} = \psi^{\text{DM}} + \psi^{\text{IPS}} + \frac{\mu_T(X)}{e_T(X)}$$

$$\mu_t(x) = \mathbb{E} [Y(t) \mid X = x]$$

$$e_t(x) = \mathbb{P}(T = t \mid X = x)$$

$$\left. \begin{array}{l} \mathcal{J}(\pi) = \mathbb{E} [\psi I(T \neq \pi(X))] \\ \psi^{\text{DM}} = T(\mu_{-1}(X) - \mu_1(X)) \\ \psi^{\text{IPS}} = \frac{-Y}{e_T(X)} \\ \psi^{\text{DR}} = \psi^{\text{DM}} + \psi^{\text{IPS}} + \frac{\mu_T(X)}{e_T(X)} \\ \mu_t(x) = \mathbb{E} [Y(t) \mid X = x] \\ e_t(x) = \mathbb{P}(T = t \mid X = x) \end{array} \right\} \min_{\pi} V(\pi) \Leftrightarrow \min_{\pi} \mathcal{J}(\pi)$$

Off-Policy Learning

1 Formalize objective via policy value



2 Estimate policy value from data

$$V(\pi) = \mathbb{E} [Y(1) I(\pi(X) = 1) + Y(-1) I(\pi(X) = -1)]$$

$$\mathcal{J}(\pi) = \mathbb{E} [\psi I(T \neq \pi(X))]$$

$$\psi^{\text{DM}} = T(\mu_{-1}(X) - \mu_1(X))$$

$$\psi^{\text{IPS}} = \frac{-Y}{e_T(X)}$$

$$\psi^{\text{DR}} = \psi^{\text{DM}} + \psi^{\text{IPS}} + \frac{\mu_T(X)}{e_T(X)}$$

$$\mu_t(x) = \mathbb{E} [Y(t) \mid X = x]$$

$$e_t(x) = \mathbb{P} (T = t \mid X = x)$$

$$\min_{\pi} V(\pi) \Leftrightarrow \min_{\pi} \mathcal{J}(\pi)$$

$$\min_{\pi} \mathcal{J}_n(\pi)$$

Estimate
from data

Off-Policy Learning

1 Formalize objective via policy value



2 Estimate policy value from data



3 Optimize over pre-specified policy class Π

$$V(\pi) = \mathbb{E} [Y(1) I(\pi(X) = 1) + Y(-1) I(\pi(X) = -1)]$$

$$\mathcal{J}(\pi) = \mathbb{E} [\psi I(T \neq \pi(X))]$$

$$\psi^{\text{DM}} = T(\mu_{-1}(X) - \mu_1(X))$$

$$\psi^{\text{IPS}} = \frac{-Y}{e_T(X)}$$

$$\psi^{\text{DR}} = \psi^{\text{DM}} + \psi^{\text{IPS}} + \frac{\mu_T(X)}{e_T(X)}$$

$$\mu_t(x) = \mathbb{E} [Y(t) \mid X = x]$$

$$e_t(x) = \mathbb{P}(T = t \mid X = x)$$

$$\min_{\pi} V(\pi) \Leftrightarrow \min_{\pi} \mathcal{J}(\pi)$$

$$\min_{\pi} \mathcal{J}_n(\pi)$$

$$\min_{\pi \in \Pi} \mathcal{J}_n(\pi)$$

Off-Policy Learning

1 Formalize objective via policy value



2 Estimate policy value from data



3 Optimize over pre-specified policy class Π

$$V(\pi) = \mathbb{E} [Y(1) I(\pi(X) = 1) + Y(-1) I(\pi(X) = -1)]$$

$$\mathcal{J}(\pi) = \mathbb{E} [\psi I(T \neq \pi(X))]$$

$$\psi^{\text{DM}} = T(\mu_{-1}(X) - \mu_1(X))$$

$$\psi^{\text{IPS}} = \frac{-Y}{e_T(X)}$$

$$\psi^{\text{DR}} = \psi^{\text{DM}} + \psi^{\text{IPS}} + \frac{\mu_T(X)}{e_T(X)}$$

$$\mu_t(x) = \mathbb{E} [Y(t) \mid X = x]$$

$$e_t(x) = \mathbb{P} (T = t \mid X = x)$$

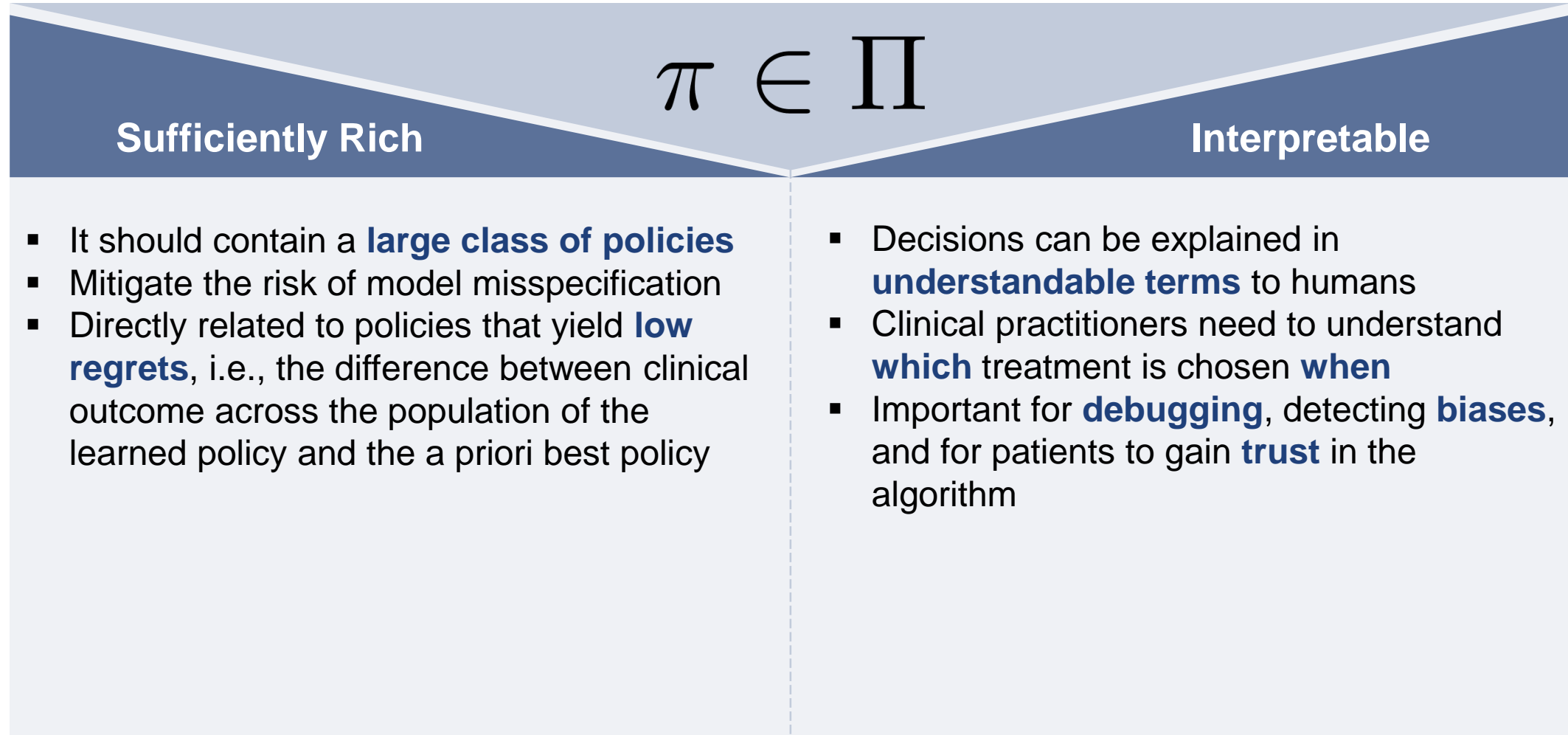
$$\min_{\pi} V(\pi) \Leftrightarrow \min_{\pi} \mathcal{J}(\pi)$$

$$\min_{\pi} \mathcal{J}_n(\pi)$$

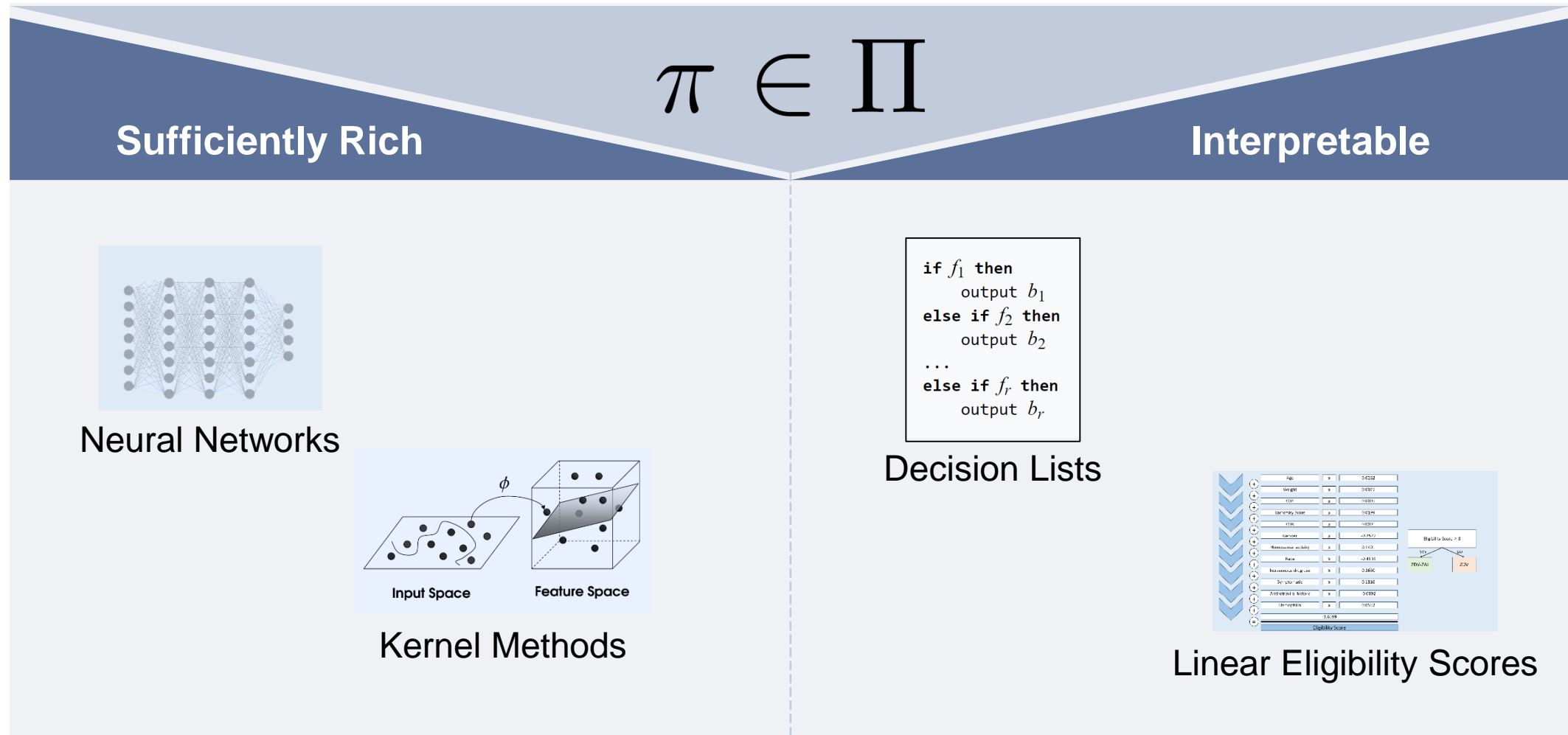
$$\min_{\pi \in \Pi} \mathcal{J}_n(\pi)$$

Requirements from clinical practice

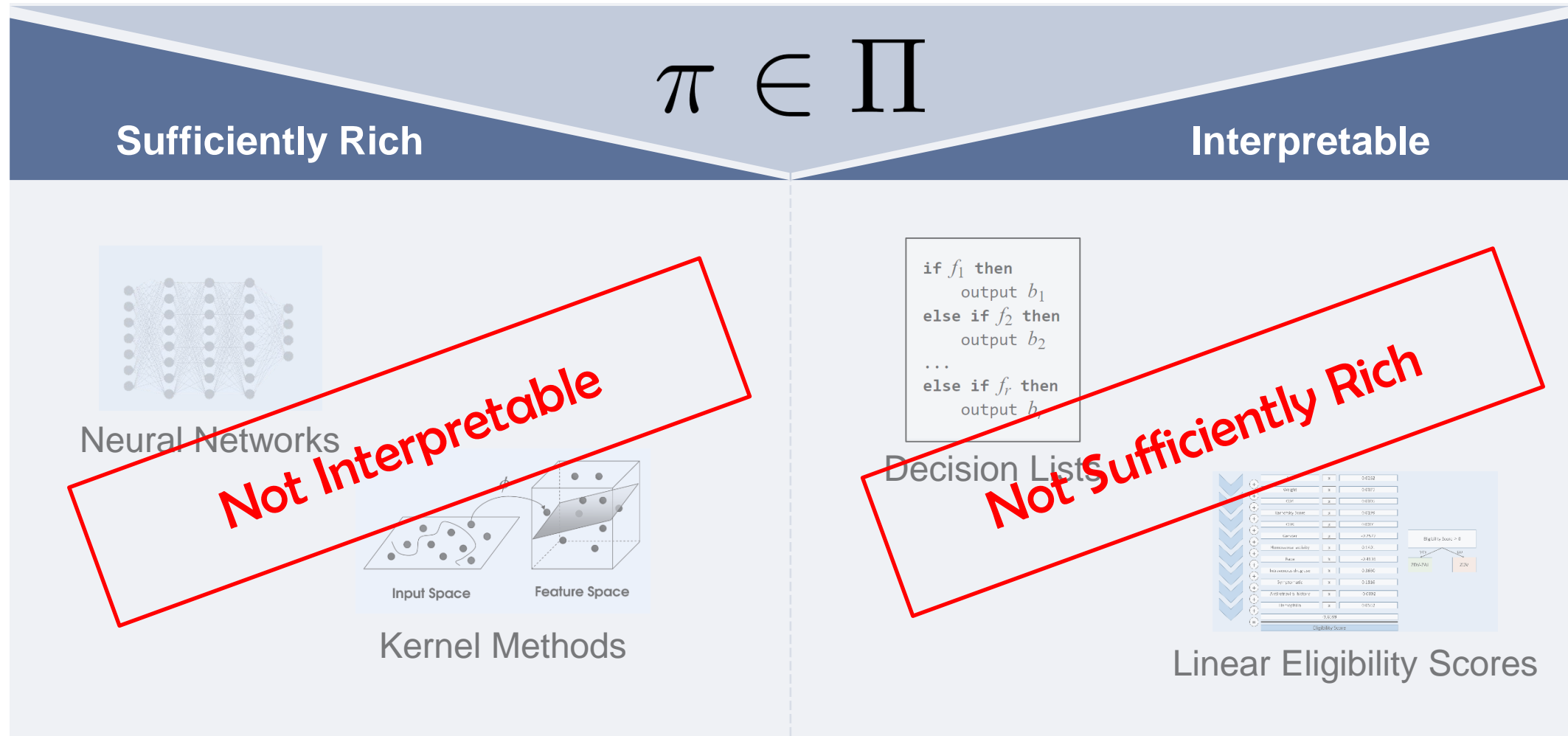
Requirements in Clinical Practice



Requirements in Clinical Practice



Requirements in Clinical Practice



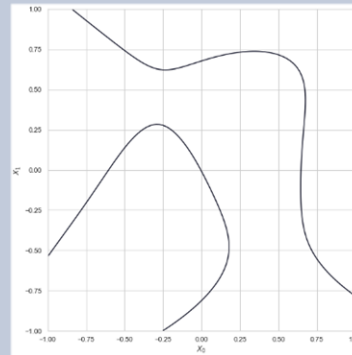
Idea: Interpretable Policy Class via Hyperboxes

$$\min_{\pi} V(\pi) = \min_{\pi} \mathbb{E} [Y(1) I(\pi(X) = 1) + Y(-1) I(\pi(X) = -1)]$$

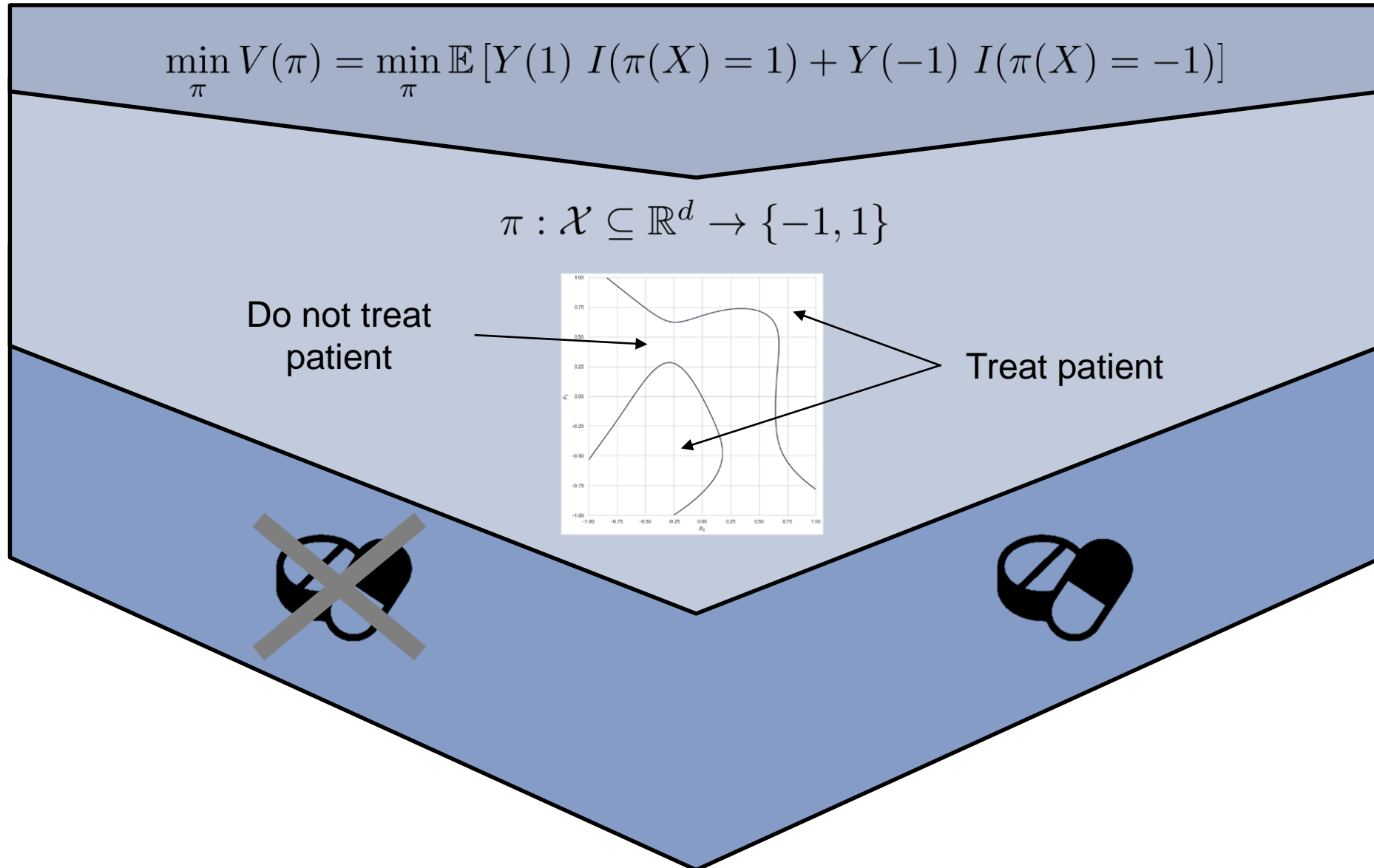
Idea: Interpretable Policy Class via Hyperboxes

$$\min_{\pi} V(\pi) = \min_{\pi} \mathbb{E} [Y(1) I(\pi(X) = 1) + Y(-1) I(\pi(X) = -1)]$$

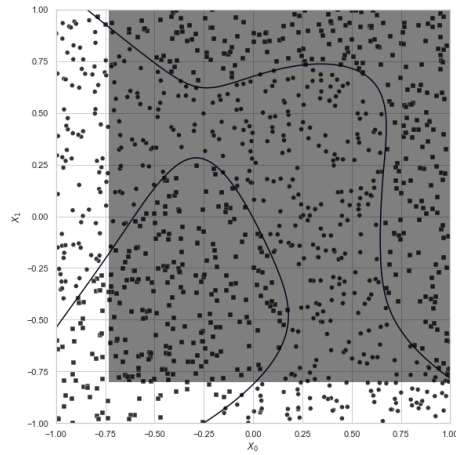
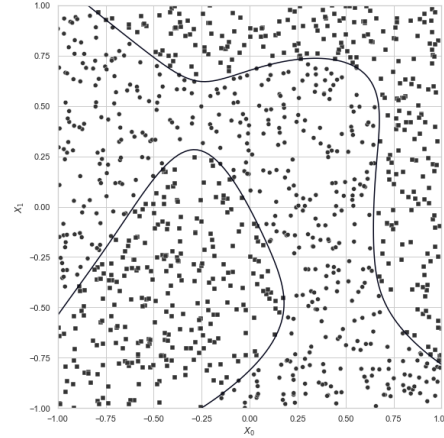
$$\pi : \mathcal{X} \subseteq \mathbb{R}^d \rightarrow \{-1, 1\}$$



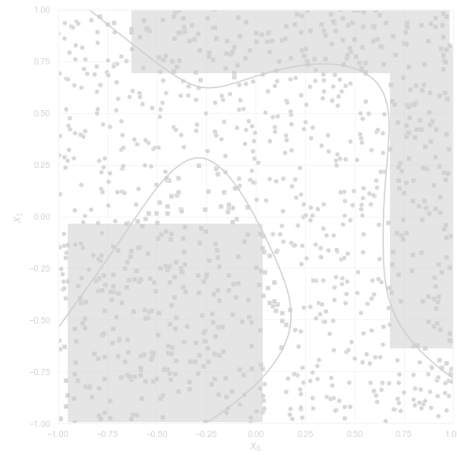
Idea: Interpretable Policy Class via Hyperboxes



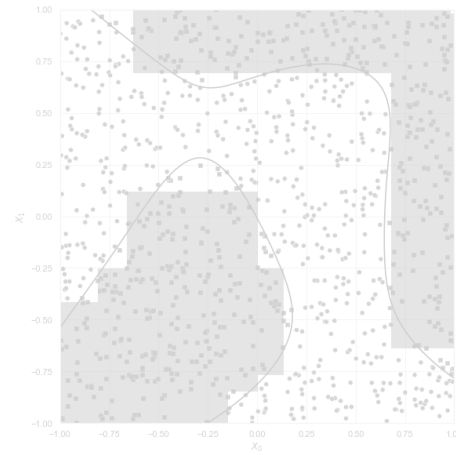
Policy Class Π_H^M : Why Hyperboxes?



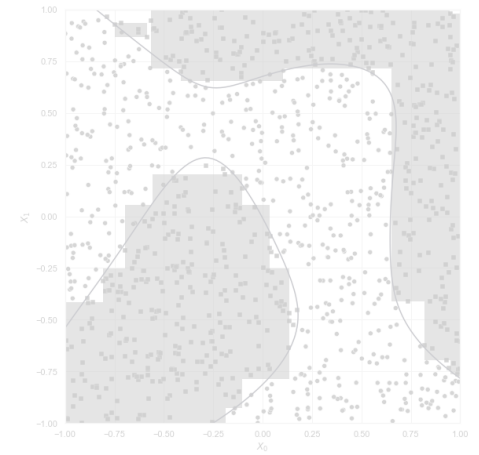
$M = 1$



$M = 3$

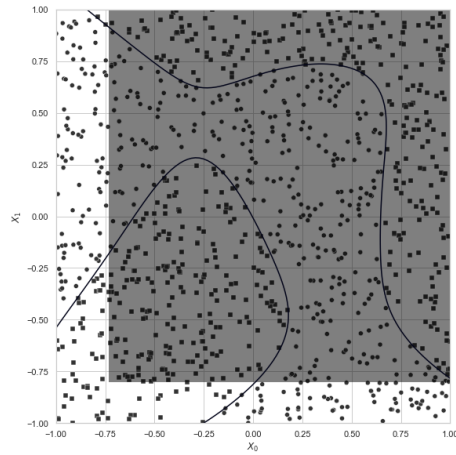
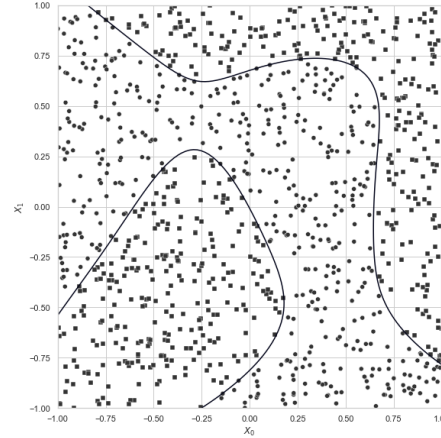
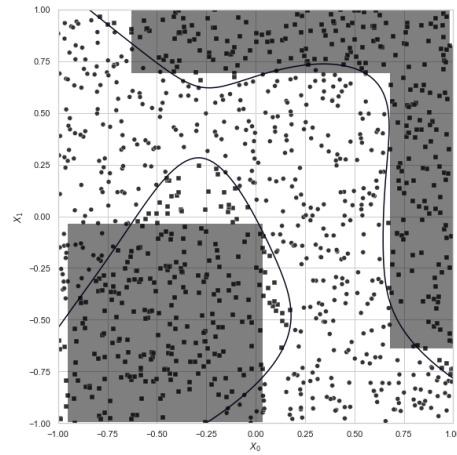
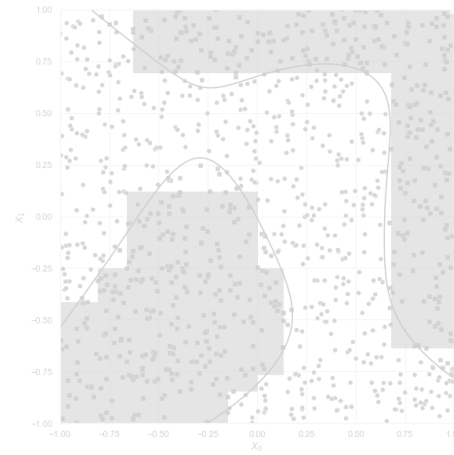
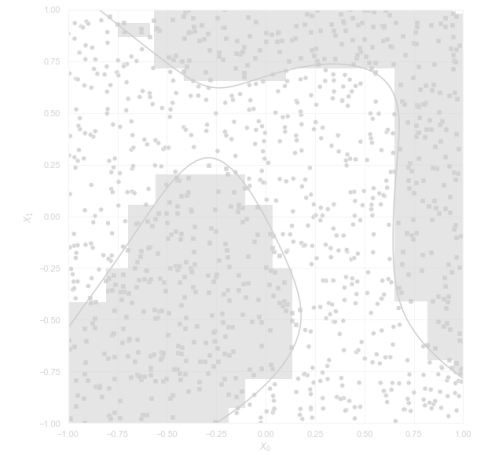


$M = 5$

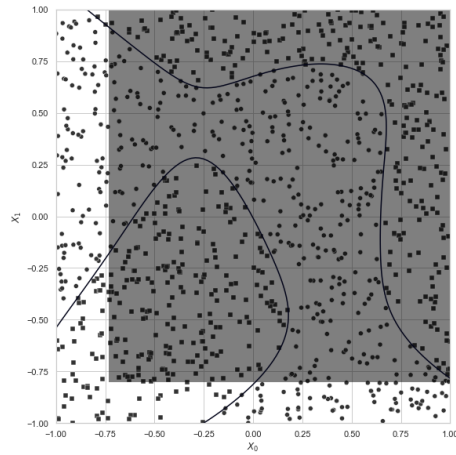
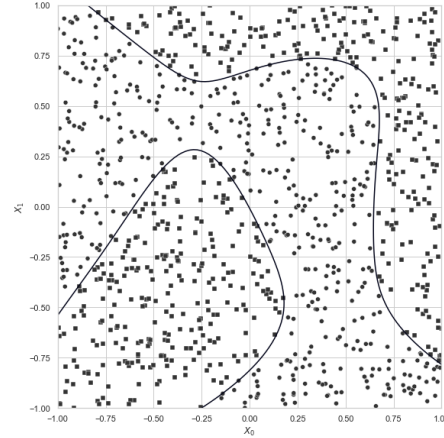


$M = 10$

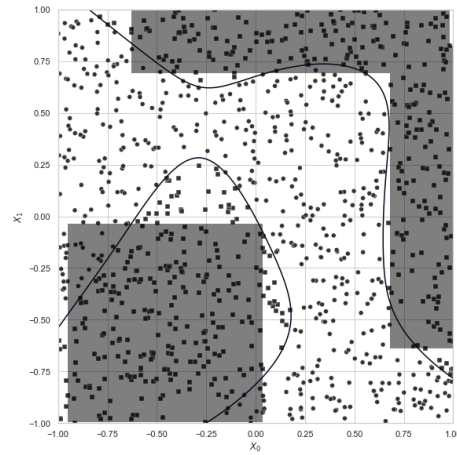
Policy Class Π_H^M : Why Hyperboxes?

 $M = 1$  $M = 3$  $M = 5$  $M = 10$

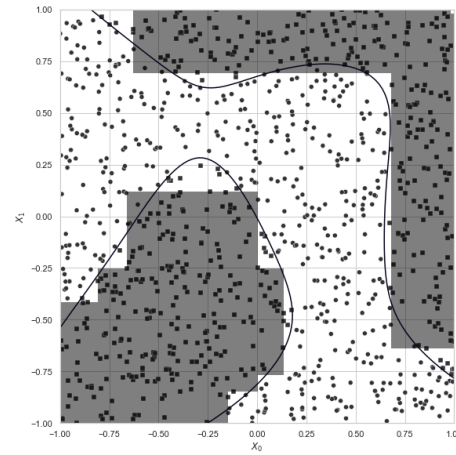
Policy Class Π_H^M : Why Hyperboxes?



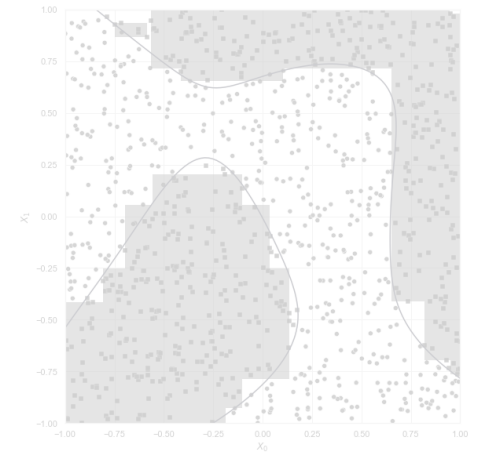
$M = 1$



$M = 3$

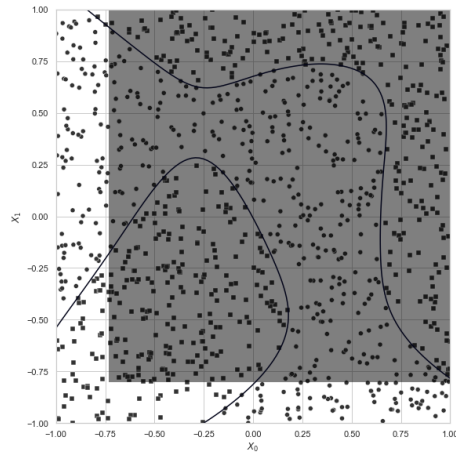
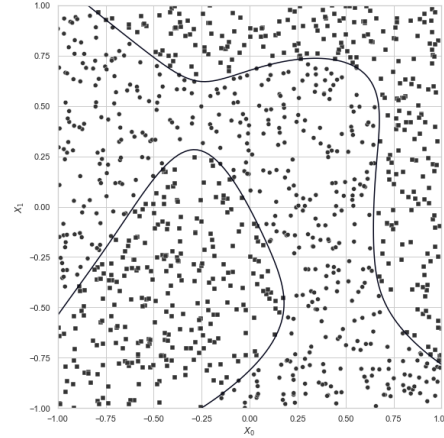
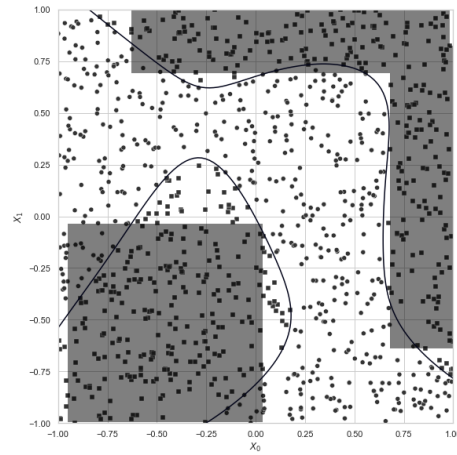
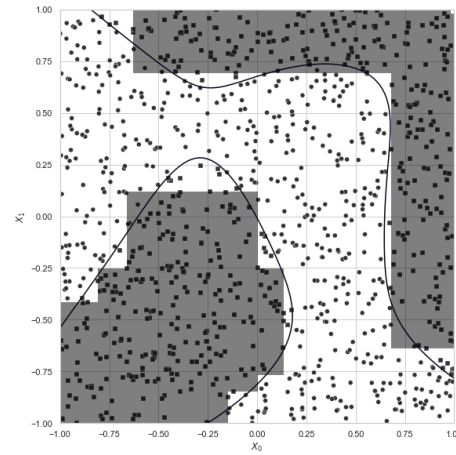
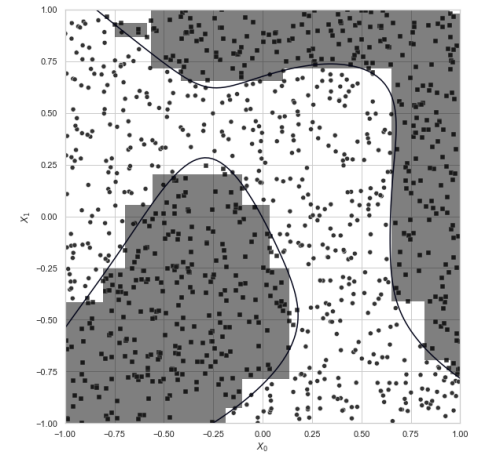


$M = 5$

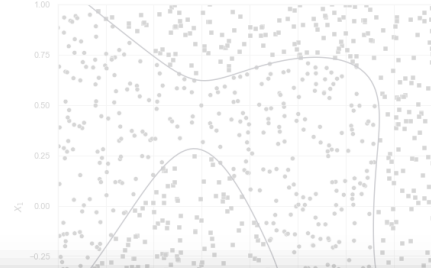


$M = 10$

Policy Class Π_H^M : Why Hyperboxes?

 $M = 1$  $M = 3$  $M = 5$  $M = 10$

Policy Class Π_H^M : Why Hyperboxes?

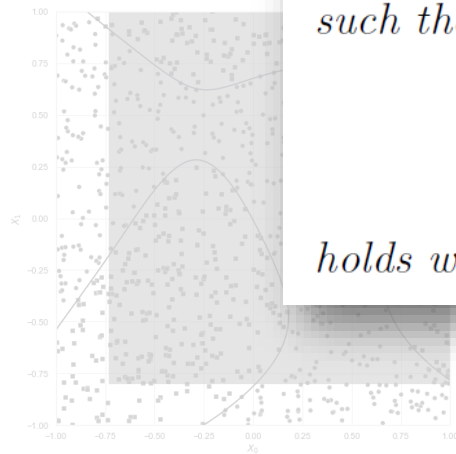


Theorem 1. *Let $1 \leq p < \infty$ and $\pi^* : \mathcal{X} \rightarrow \{-1, 1\}$ be any Lebesgue measurable function⁴. Then, for every $\delta \in (0, 1)$ and $\epsilon > 0$, there exists a sample size $n_{\delta, \epsilon} \in \mathbb{N}$ and $M \in \mathbb{N}$ sufficiently large, as well as, a policy $\pi_{\mathcal{D}^*} \in \Pi_H^M$ as defined in (2.3), such that*

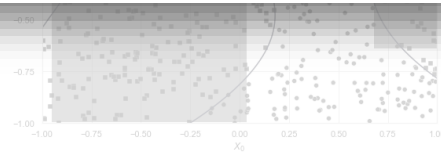
$$\|\pi^* - \pi_{\mathcal{D}^*}\|_p < \epsilon \quad (2.15)$$

holds with probability at least $1 - \delta$.

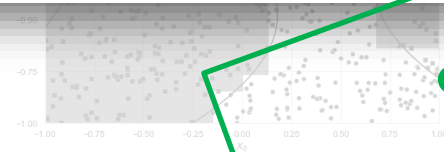
Sufficiently Rich



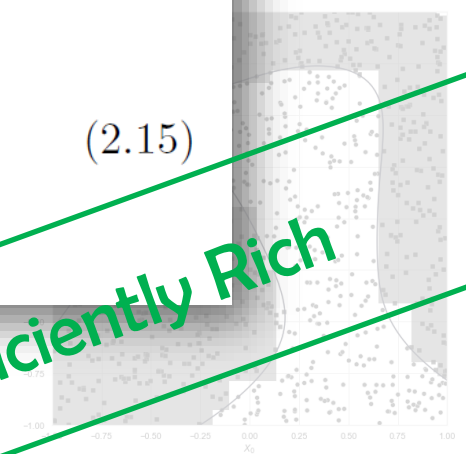
$M = 1$



$M = 3$

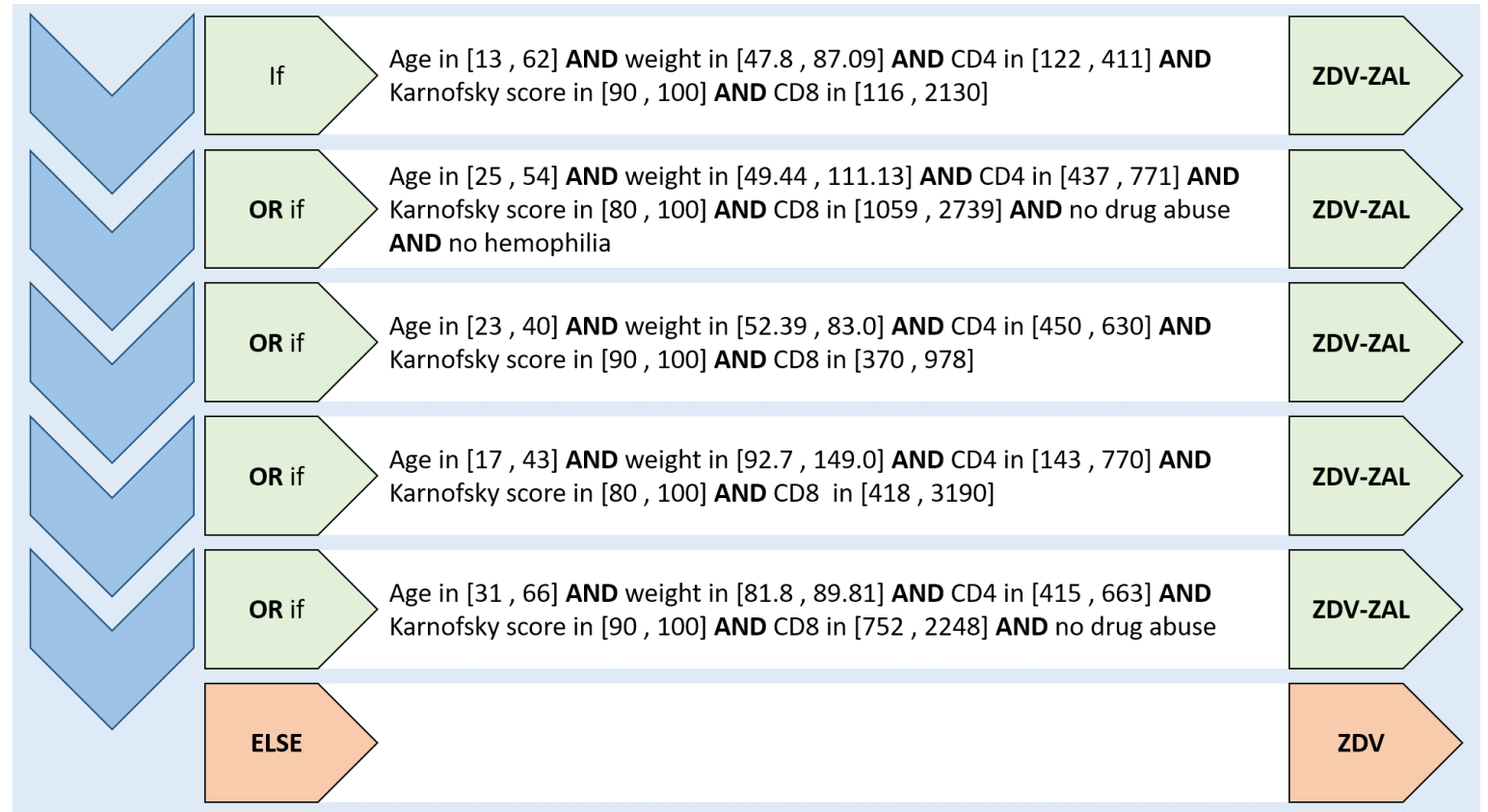
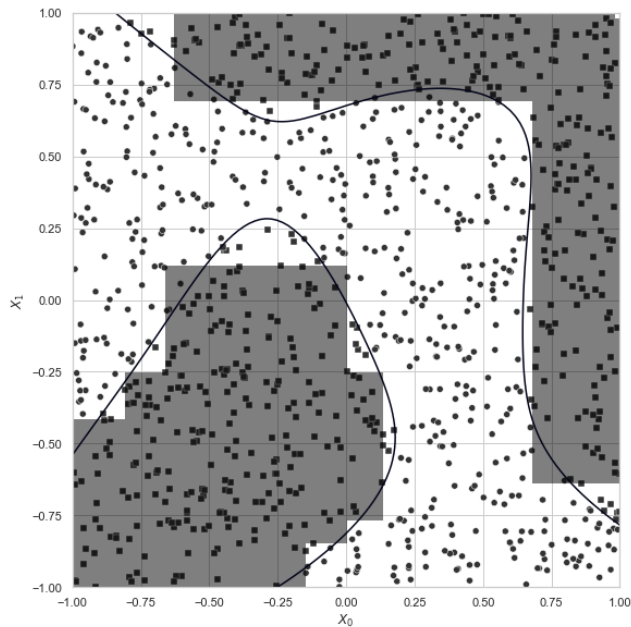


$M = 5$

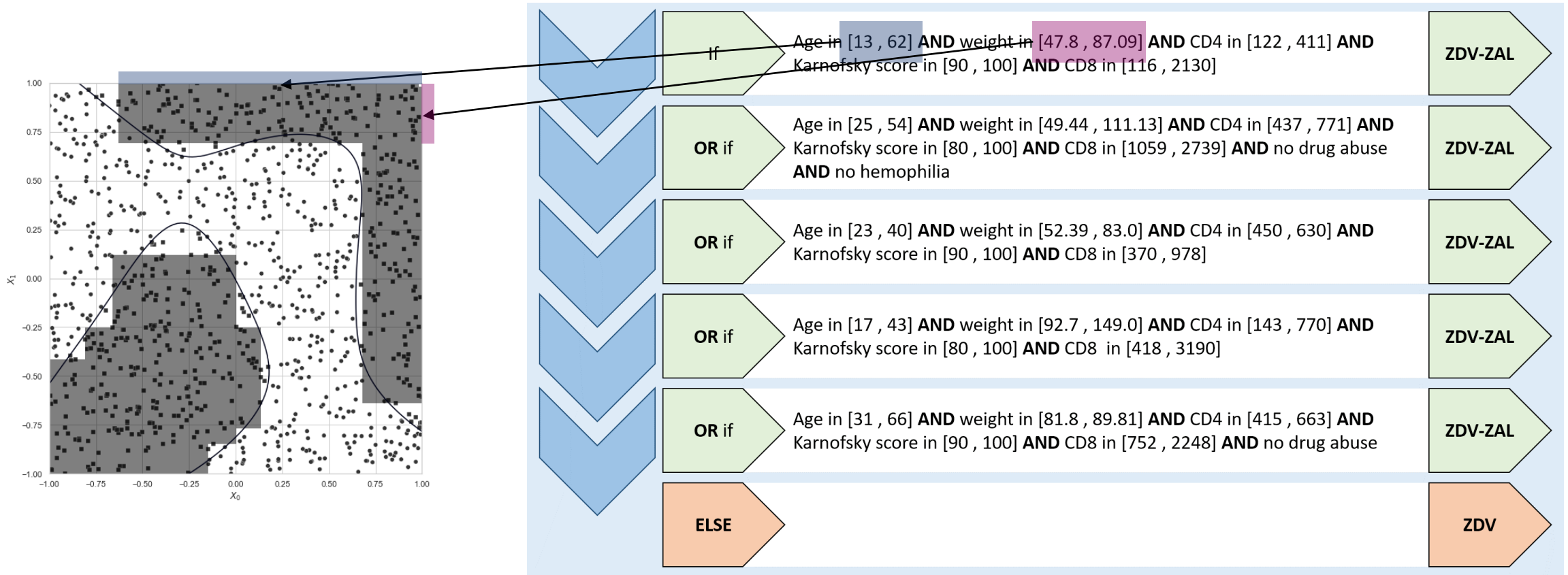


$M = 10$

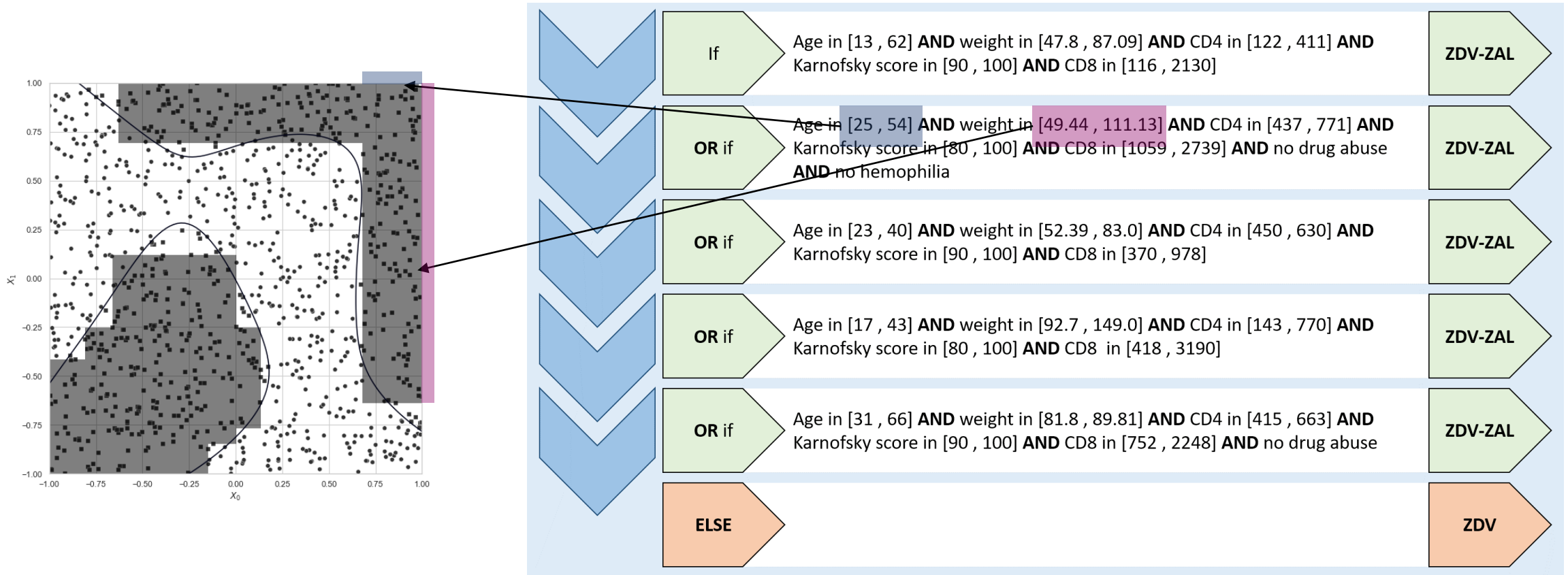
Policy Class Π_H^M : Why are Hyperboxes Interpretable?



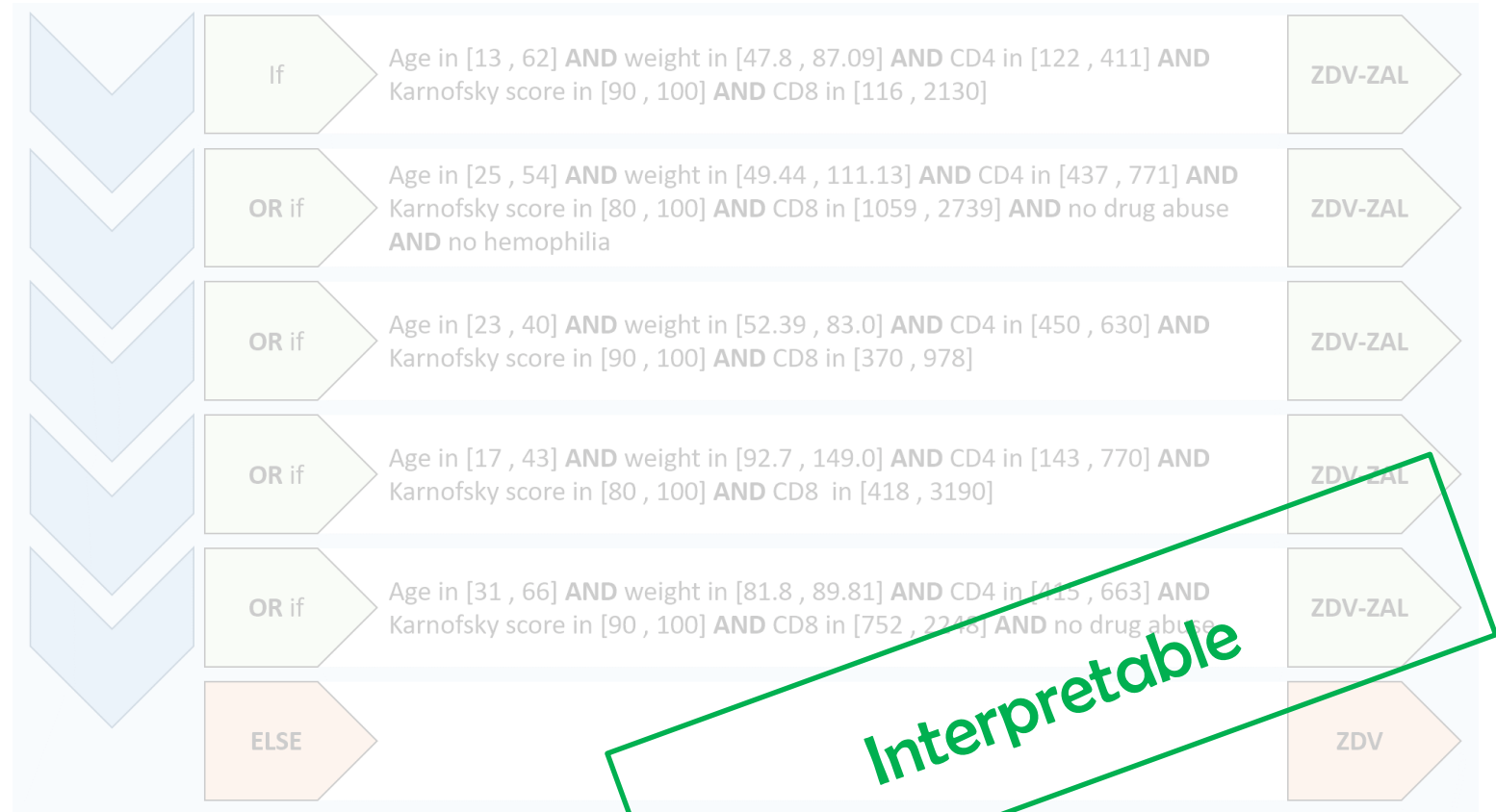
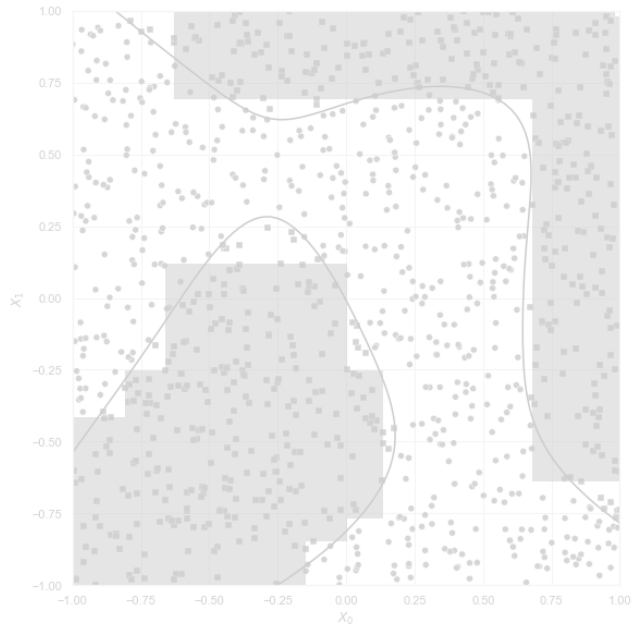
Policy Class Π_H^M : Why are Hyperboxes Interpretable?



Policy Class Π_H^M : Why are Hyperboxes Interpretable?



Policy Class Π_H^M : Why are Hyperboxes Interpretable?



IOPL Algorithm

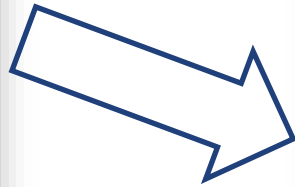
MILP Formulation of Off-Policy Learning

$$\begin{aligned}
 & \min \frac{1}{n} \sum_{i=1}^n \psi_i \xi_i \\
 \text{s.t. } & \xi_i + \sum_{j \in \mathcal{K}_i} s_j \geq 1 && \text{for } i \in I_1 \cap \mathcal{P}, \\
 & \xi_i \geq s_j && \text{for } i \in I_{-1} \cap \mathcal{P} \text{ and } j \in \mathcal{K}_i, \\
 & \xi_i \leq 1 - s_j && \text{for } i \in I_1 \cap \mathcal{N} \text{ and } j \in \mathcal{K}_i, \\
 & \xi_i \leq \sum_{j \in \mathcal{K}_i} s_j && \text{for } i \in I_{-1} \cap \mathcal{N}, \\
 & \sum_{j=1}^N s_j \leq M, \\
 & s_j \in \{0, 1\}, \xi_i \in [0, 1].
 \end{aligned}$$

IOPL Algorithm

MILP Formulation of Off-Policy Learning

$$\begin{aligned}
 & \min \frac{1}{n} \sum_{i=1}^n \psi_i \xi_i \\
 \text{s.t. } & \xi_i + \sum_{j \in \mathcal{K}_i} s_j \geq 1 && \text{for } i \in I_1 \cap \mathcal{P}, \\
 & \xi_i \geq s_j && \text{for } i \in I_{-1} \cap \mathcal{P} \text{ and } j \in \mathcal{K}_i, \\
 & \xi_i \leq 1 - s_j && \text{for } i \in I_1 \cap \mathcal{N} \text{ and } j \in \mathcal{K}_i, \\
 & \xi_i \leq \sum_{j \in \mathcal{K}_i} s_j && \text{for } i \in I_{-1} \cap \mathcal{N}, \\
 & \sum_{j=1}^N s_j \leq M, \\
 & s_j \in \{0, 1\}, \xi_i \in [0, 1].
 \end{aligned}$$



IOPL is a highly efficient **branch-and-price algorithm**, i.e., a column generation procedure within a branch-and-bound framework

Interpretable Off-Policy Learning

Algorithm 1: IOPL

```

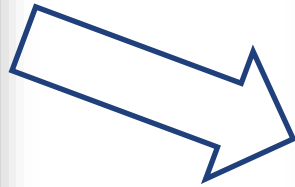
Input: Initial working set  $\mathcal{W}_0$ 
Output: Optimal subset  $\mathcal{K}^* \subseteq \mathcal{K}$ , optimal solution  $s^*$ 
1 Initialize the list of active subproblems  $\mathcal{L} \leftarrow \{\text{RMILP}(\mathcal{W}_0, \emptyset)\}$ 
2 Initialize iteration counter  $l \leftarrow 0$ 
3 while  $\mathcal{L}$  not empty do
4   Select and remove first subproblem  $\text{RMILP}(\mathcal{W}, \mathcal{C})$  from  $\mathcal{L}$ 
5   Perform column generation to get the current relaxed solution and the new working set
      $s', v', \mathcal{W}' \leftarrow \text{ColumnGeneration}(\text{RMILP}(\mathcal{W}, \mathcal{C}))$ 
6   if  $l = 0$  then
7     Solve restricted integer problem to get current optimal integer solution and objective value
        $s^*, v^* \leftarrow \text{Solve}(\text{MILP}(\mathcal{W}'))$ 
8     Update optimal subset  $\mathcal{W}^* \leftarrow \mathcal{W}'$ 
9   end
10  if  $v' \leq v^*$  then
11    if  $s'$  integral then
12      Update new optimal integer solution  $(\mathcal{W}^*, s^*) \leftarrow (\mathcal{W}', s')$ 
13    else
14      Set  $j'$  according to branching rule
15      Branch by updating  $\mathcal{L} \leftarrow \mathcal{L} \cup \{\text{RMILP}(\mathcal{W}', \mathcal{C} \cup \{(j', 1)\}), \text{RMILP}(\mathcal{W}', \mathcal{C} \cup \{(j', 0)\})\}$ 
16      Solve restricted problem  $(s', v') \leftarrow \text{Solve}(\text{MILP}(\mathcal{W}'))$ 
17      if  $v' \leq v^*$  then
18        Update new optimal integer solution  $(\mathcal{W}^*, s^*) \leftarrow (\mathcal{W}', s')$ 
19      end
20    end
21  end
22  Increase counter  $l \leftarrow l + 1$ 
23 end

```

IOPL Algorithm

MILP Formulation of Off-Policy Learning

$$\begin{aligned}
 & \min \frac{1}{n} \sum_{i=1}^n \psi_i \xi_i \\
 \text{s.t. } & \xi_i + \sum_{j \in \mathcal{K}_i} s_j \geq 1 && \text{for } i \in I_1 \cap \mathcal{P}, \\
 & \xi_i \geq s_j && \text{for } i \in I_{-1} \cap \mathcal{P} \text{ and } j \in \mathcal{K}_i, \\
 & \xi_i \leq 1 - s_j && \text{for } i \in I_1 \cap \mathcal{N} \text{ and } j \in \mathcal{K}_i, \\
 & \xi_i \leq \sum_{j \in \mathcal{K}_i} s_j && \text{for } i \in I_{-1} \cap \mathcal{N}, \\
 & \sum_{j=1}^N s_j \leq M, \\
 & s_j \in \{0, 1\}, \xi_i \in [0, 1].
 \end{aligned}$$



IOPL is a highly efficient **branch-and-price algorithm**, i.e., a column generation procedure within a branch-and-bound framework

Interpretable Off-Policy Learning

Algorithm 1: IOPL

```

Input: Initial working set  $\mathcal{W}_0$ 
Output: Optimal subset  $\mathcal{K}^* \subseteq \mathcal{K}$ , optimal solution  $s^*$ 
1 Initialize the list of active subproblems  $\mathcal{L} \leftarrow \{\text{RMILP}(\mathcal{W}_0, \emptyset)\}$ 
2 Initialize iteration counter  $l \leftarrow 0$ 
3 while  $\mathcal{L}$  not empty do
4   Select and remove first subproblem  $\text{RMILP}(\mathcal{W}, \mathcal{C})$  from  $\mathcal{L}$ 
5   Perform column generation to get the current relaxed solution and the new working set
      $s', v', \mathcal{W}' \leftarrow \text{ColumnGeneration}(\text{RMILP}(\mathcal{W}, \mathcal{C}))$ 
6   if  $l = 0$  then
7     Solve restricted integer problem to get current optimal integer solution and objective value
        $s^*, v^* \leftarrow \text{Solve}(\text{MILP}(\mathcal{W}'))$ 
8     Update optimal subset  $\mathcal{W}^* \leftarrow \mathcal{W}'$ 
9   end
10  if  $v' \leq v^*$  then
11    if  $s'$  integral then
12      Update new optimal integer solution  $(\mathcal{W}^*, s^*) \leftarrow (\mathcal{W}', s')$ 
13    else
14      Set  $j'$  according to branching rule
15      Branch by updating  $\mathcal{L} \leftarrow \mathcal{L} \cup \{\text{RMILP}(\mathcal{W}', \mathcal{C} \cup \{(j', 1)\}), \text{RMILP}(\mathcal{W}', \mathcal{C} \cup \{(j', 0)\})\}$ 
16      Solve restricted problem  $(s', v') \leftarrow \text{Solve}(\text{MILP}(\mathcal{W}'))$ 
17      if  $v' \leq v^*$  then
18        Update new optimal integer solution  $(\mathcal{W}^*, s^*) \leftarrow (\mathcal{W}', s')$ 
19      end
20    end
21  end
22  Increase counter  $l \leftarrow l + 1$ 
23 end

```

For more details on our algorithm, experiments with baselines, more theoretical results, and proofs, see the paper

Thank you!

Interpretable Off-Policy Learning via Hyperbox Search

Daniel Tschernutter¹



Tobias Hatt¹



Stefan Feuerriegel^{1,2}



1

ETH zürich

2

