$$X_i^e = g(X_{S_i}^e) + \epsilon^e, \text{ for some } i \in \{1, \ldots, m\}, \qquad (1)$$

(3)

$$\hat{Y}_i^{\text{test}} =: \quad \frac{\mathsf{E}_{\mathcal{P}_e}[X_i|X_1] - \mathsf{E}_{\mathcal{P}_e}[X_i|X_1, Y=0]}{}$$

The noise variables $\epsilon_Y$ and $\epsilon_3$ are i.i.d. $\mathcal{N}(0, \sigma^2)$. Suppose

particular $e \in \mathcal{E}$ becomes difficult as $\beta_1^e$ and $\mu_2^e$ vary with



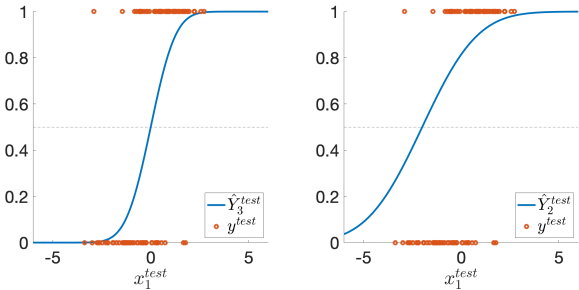Fig. 1: Comparisons of $\hat{Y}_3^{\text{test}}$ (left) and $\hat{Y}_2^{\text{test}}$ (right), where

$$\mathsf{E}_{\mathcal{P}_e}[X_3|X_1 = x_1] - \mathsf{E}_{\mathcal{P}_e}[X_3|X_1 = x_1, Y = 0]$$

**Definition 1.** *For $k \in \{1, \ldots, m\}$, $S \subseteq \{1, \ldots, m\}\backslash k$, and $h(X_S, Y) := \mathsf{E}_{\mathcal{P}_e}[X_k|X_S, Y]$, the pair $(k, S)$ satisfies the*

$$\frac{}{h(X_S, 1) - h(X_S, 0)}, \tag{8}$$

1) $X_k^e = g(X_R^e, Y^e) + \epsilon^e$ *as in* (1) ,
2) $X_Q^e \perp\!\!\!\perp X_k^e \mid (X_R^e, \; Y^e)$ .

$$\mathsf{E}_{\mathcal{P}_e}[Y | \phi_e(X) = (x_Q, x_R, z)]$$

(9)

$\square$

$$\begin{cases} X_1^e := f_1^e(X_{PA(X_1^e)}^e, \; \epsilon_1^e), \end{cases}$$

(10)

all $i \in \{0, \ldots, m\}$.

some $i \in \{0, \ldots, m\}$, let $f_i^e(X_{PA(X^e)}^e, \epsilon_i^e) = g(X_{PA(X^e)}^e) +$

$\}$

**Input:** $\boldsymbol{Y}^e$, for each $e \in \mathcal{E}_{\text{train}}$, and $\boldsymbol{X}^e$, for each $e \in \mathcal{E}_{\text{obs}}$

---

$$\frac{}{|\ ,_{\text{inv}}|} \sum_{(k,S) \in \mathcal{T}_{\text{inv}}} \hat{\boldsymbol{Y}}_{k,S}^{\text{test}}$$

---

**for** each $e \in \mathcal{E}_{\text{train}}$ and $i \in \{0, 1\}$ **do**

$\quad \boldsymbol{R}_i^e = \boldsymbol{X}_{k.Y=i}^e - \hat{g}_i(\boldsymbol{X}_{S.Y=i}^e)$

$\quad \boldsymbol{R}$

$\quad \text{pval}_i^e = t\text{-test}(\boldsymbol{R}_i^e, \boldsymbol{R}_i^{-e})$

$\{3, \dots, 7\}$. For each $i \in \{2, \dots, m\}$ and $e \in \mathcal{E}_{\text{train}}$, $X_i^e \sim$

for $e = e^1$, $[0, 2]$ for $e = e^2$, and $[0, 3]$ for $e = e^{\text{test}}$. Then, where $S_1 = \{2, \dots, m\}$, $Y^e | X_{S_1}^e$ follows a logistic model such that $\mathcal{P}_e(Y = 1 | X_{S_1}) = 1/(1 + e^{-X_{S_1} \beta^e})$ for $e \in \mathcal{E}_{\text{train}}$. For $e^{\text{test}}$, $Y^{\text{test}} | X_{S_1}^{\text{test}}$ follows a probit model such that $Y^{\text{test}} = 1$, if $X_{S,}^{\text{test}} \beta^{\text{test}} + \epsilon < 0$, where $\epsilon \sim \mathcal{N}(0, 1)$. For all $e \in \mathcal{E}_{\text{obs}}$,

then scaled such that they sum to one. For all $e \in \mathcal{E}_{\text{obs}}$,

Specifically, $g_1(X_{S,}^e) = X_{S,}^e \eta_1$ and $g_0(X_{S,}^e) = X_{S,}^e \eta_0$. The

**Two real-world data.** We also include experiments on two real datasets: *census* [18] and *mushroom* [19]. The census

14 societal and demographic variables such as age, education,

whether or not an individual's income exceeded 50k/yr. The

naturally growing mushrooms' size, shape, and color and

or paths. Results in Table II indicate that bIMP outperforms

## VII. ACKNOWLEDGEMENTS

$e \in \mathcal{E}_{\text{obs}}$. Without loss of generality, let $X_i^e$ be continuous for all $i \in \{1, \dots, m\}$. The pdf of $X_k^e | X_S^e$ for any $e \in \mathcal{E}_{\text{obs}}$ is

$$
\begin{aligned}
f_{X_k^e | X_S^e}(x_k | x) \\
&= f_{X_k^e | X_S^e, Y^e}(x_k | x, 1) \cdot p_{Y^e | X_S^e}(1 | x) \\
&\quad + f_{X_k^e | X_S^e, Y^e}(x_k | x, 0) \cdot p_{Y^e | X_S^e}(0 | x) \\
&= f_{X_k^e | X_S^e, Y^e}(x_k | x, 1) \cdot p_{Y^e | X_S^e}(1 | x) \\
&\quad + f_{X_k^e | X_S^e, Y^e}(x_k | x, 0) \cdot \left[ 1 - p_{Y^e | X_S^e}(1 | x) \right] \\
&= p_{Y^e | X_S^e}(1 | x) \left[ f_{X_k^e | X_S^e, Y^e}(x_k | x, 1) - f_{X_k^e | X_S^e, Y^e}(x_k | x, 0) \right]
\end{aligned}
\tag{12}
$$

$$
\int_{-\infty}^{\infty} x_k \cdot f_{X_k^e | X_S^e}(x_k | x) \, dx_k
$$

$$
- \mathsf{E}_{\mathcal{P}_e}[Y | X_S = x] \cdot \mathsf{E}_{\mathcal{P}_e}[X_k | X_S = x, Y = 0]
\tag{13}
$$

$$
\frac{\mathsf{E}_{\mathcal{P}_e}[X_k | X_S] - \mathsf{E}_{\mathcal{P}_e}[X_k | X_S, Y = 0]}{\mathsf{E}_{\mathcal{P}_e}[X_k | X_S, Y = 1] - \mathsf{E}_{\mathcal{P}_e}[X_k | X_S, Y = 0]}.
\tag{14}
$$

on $e$ and (II) the denominator of (14) is non-zero. Since $X_S^e = (X_R^e, X_Q^e)$,

$$
\mathsf{E}_{\mathcal{P}_e}[X_k | X_S, Y] = \mathsf{E}_{\mathcal{P}_e}[X_k | X_R, X_Q, Y] \overset{(a)}{=} \mathsf{E}_{\mathcal{P}_e}[X_k | X_R, Y]
$$

$$
\overset{(b)}{=} \mathsf{E}_{\mathcal{P}_e}[g(X_R, Y) + \epsilon | X_R, Y] = g(X_R^e, Y^e),
\tag{15}
$$

where $(a)$ follows since $X_Q^e \perp\!\!\!\perp X_k^e | X_R^e, Y^e$, $(b)$ follows from

$\epsilon$ has zero mean. Thus, the $\mathsf{E}_{\mathcal{P}_e}[X_k | X_S = (x_Q, x_R), Y = y]$ does not depend on $e$ as $\mathsf{E}_{\mathcal{P}_e}[X_k | X_S = (x_Q, x_R), Y = y] =$

*mushroom* data below.

| | | | |
|---|---|---|---|
| meadows | 76.0 | 87.5 | 46.2 |
| paths | 88.1 | 90.9 | 11.8 |

□

## REFERENCES

.

[12] M. Rojas-Carulla, B. Schölkopf, R. Turner, and J. Peters, "Invariant

[13] D. Rothenhäusler, N. Meinshausen, P. Bühlmann, and J. Peters, "Anchor

[6] B. Schölkopf, D. Janzing, J. Peters, E. Sgouritsa, K. Zhang, and J. Mooij,

.