## A. CATE Interval

**Lemma 2.** *The unbiased estimate of the expected potential outcome under hidden confounding, given in Equation* (3) *has the following equivalent characterization:*

$$\mathbb{E}\left[Y^t \mid \mathbf{X} = \mathbf{x}\right] = \mu(\mathbf{x}, t) + \frac{\int (y - \mu(\mathbf{x}, t)) w_t(y \mid \mathbf{x}) f(y \mid \mathbf{x}, t) dy}{\int w_t(y \mid \mathbf{x}) f(y \mid \mathbf{x}, t) dy}. \tag{12}$$

*Proof.*

$$\mathbb{E}\left[Y^t \mid \mathbf{X} = \mathbf{x}\right] = \mu(w_t; \mathbf{x}, t) \tag{13a}$$

$$= \frac{\int y w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x}) dy}{\int w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x}) dy} \tag{13b}$$

$$= \int y \frac{w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x})}{\int w_t(y' \mid \mathbf{x}) f_t(y' \mid \mathbf{x}) dy'} dy \tag{13c}$$

$$= \mu(\mathbf{x}, t) + \int y \frac{w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x})}{\int w_t(y' \mid \mathbf{x}) f_t(y' \mid \mathbf{x}) dy'} dy - \mu(\mathbf{x}, t) \tag{13d}$$

$$= \mu(\mathbf{x}, t) + \int y \frac{w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x})}{\int w_t(y' \mid \mathbf{x}) f_t(y' \mid \mathbf{x}) dy'} dy - \mu(\mathbf{x}, t) \int \frac{w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x})}{\int w_t(y' \mid \mathbf{x}) f_t(y' \mid \mathbf{x}) dy'} dy \tag{13e}$$

$$= \mu(\mathbf{x}, t) + \int y \frac{w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x})}{\int w_t(y' \mid \mathbf{x}) f_t(y' \mid \mathbf{x}) dy'} dy - \int \mu(\mathbf{x}, t) \frac{w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x})}{\int w_t(y' \mid \mathbf{x}) f_t(y' \mid \mathbf{x}) dy'} dy \tag{13f}$$

$$= \mu(\mathbf{x}, t) + \int (y - \mu(\mathbf{x}, t)) \frac{w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x})}{\int w_t(y' \mid \mathbf{x}) f_t(y' \mid \mathbf{x}) dy'} dy \tag{13g}$$

$$= \mu(\mathbf{x}, t) + \frac{\int (y - \mu(\mathbf{x}, t)) w_t(y \mid \mathbf{x}) e_t(\mathbf{x}) f(y \mid \mathbf{x}, t) dy}{\int w_t(y \mid \mathbf{x}) e_t(\mathbf{x}) f(y \mid \mathbf{x}, t) dy} \tag{13h}$$

$$= \mu(\mathbf{x}, t) + \frac{\int (y - \mu(\mathbf{x}, t)) w_t(y \mid \mathbf{x}) f(y \mid \mathbf{x}, t) dy}{\int w_t(y \mid \mathbf{x}) f(y \mid \mathbf{x}, t) dy}. \tag{13i}$$

$\square$

**Lemma 3.** *The bounds for the conditional expected potential outcomes $\underline{\mu}^\Gamma(\mathbf{x}, t)$ and $\overline{\mu}^\Gamma(\mathbf{x}, t)$ defined in equations* (4) *have the following equivalent characterization:*

$$\underline{\mu}^\Gamma(\mathbf{x}, t) = \inf_{y^* \in \mathcal{Y}} \mu(\mathbf{x}, t) + \frac{\int_{-\infty}^{y^*} r_t(y) f(y \mid \mathbf{x}, t) dy}{\alpha_t'^\Gamma(\mathbf{x}) + P(Y \le y^* \mid \mathbf{x}, t)},$$

$$\overline{\mu}^\Gamma(\mathbf{x}, t) = \sup_{y^* \in \mathcal{Y}} \mu(\mathbf{x}, t) + \frac{\int_{y^*}^{\infty} r_t(y) f(y \mid \mathbf{x}, t) dy}{\alpha_t'^\Gamma(\mathbf{x}) + P(Y > y^* \mid \mathbf{x}, t)},$$

*where $r_t(y) = (y - \mu(\mathbf{x}, t))$ and $\alpha_t'^\Gamma(\mathbf{x}) = \frac{\alpha_t(\mathbf{x}; \Gamma)}{\beta_t(\mathbf{x}; \Gamma) - \alpha_t(\mathbf{x}; \Gamma)}$.*

*Proof.* We prove the result for $\underline{\mu}^\Gamma(\mathbf{x}, t)$ and the result for $\overline{\mu}^\Gamma(\mathbf{x}, t)$ can be proved analogously. From Kallus et al. (2019) Lemma 1,

$$\underline{\mu}(\mathbf{x}, t) = \inf_{w_t(y \mid \mathbf{x}) \in [\alpha_t(\mathbf{x}; \Gamma), \beta_t(\mathbf{x}; \Gamma)]} \frac{\int y w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x}) dy}{\int w_t(y \mid \mathbf{x}) f_t(y \mid \mathbf{x}) dy}$$

$$= \inf_{u \in \mathcal{U}^{ni}} \frac{\alpha_t(\mathbf{x}; \Gamma) \int y f_t(y \mid \mathbf{x}) dy + (\beta_t(\mathbf{x}; \Gamma) - \alpha_t(\mathbf{x}; \Gamma)) \int u(y) y f_t(y \mid \mathbf{x}) dy}{\alpha_t(\mathbf{x}; \Gamma) \int f_t(y \mid \mathbf{x}) dy + (\beta_t(\mathbf{x}; \Gamma) - \alpha_t(\mathbf{x}; \Gamma)) \int u(y) f_t(y \mid \mathbf{x}) dy}, \tag{15}$$

$$= \underline{\mu}^\Gamma(\mathbf{x}, t)$$

where $\mathcal{U}^{ni} = \{u : \mathcal{Y} \to [0, 1] \mid u(y) \text{ is non-increasing}\}$. Therefore, from the equivalence in Equation (13),

$$
\begin{aligned}
\underline{\mu}^{\Gamma}(\mathbf{x}, \mathrm{t}) &= \inf_{u \in \mathcal{U}^{ni}} \mu(\mathbf{x}, \mathrm{t}) + \frac{\alpha_{\mathrm{t}}(\mathbf{x}) \int (y - \mu(\mathbf{x}, \mathrm{t})) f(y \mid \mathbf{x}, \mathrm{t}) dy + (\beta_{\mathrm{t}}(\mathbf{x}) - \alpha_{\mathrm{t}}(\mathbf{x})) \int u(y)(y - \mu(\mathbf{x}, \mathrm{t})) f(y \mid \mathbf{x}, \mathrm{t}) dy}{\alpha_{\mathrm{t}}(\mathbf{x}) \int f(y \mid \mathbf{x}, \mathrm{t}) dy + (\beta_{\mathrm{t}}(\mathbf{x}) - \alpha_{\mathrm{t}}(\mathbf{x})) \int u(y) f(y \mid \mathbf{x}, \mathrm{t}) dy} \\
&= \inf_{u \in \mathcal{U}^{ni}} \mu(\mathbf{x}, \mathrm{t}) + \frac{\alpha_{\mathrm{t}}(\mathbf{x})(\mu(\mathbf{x}, \mathrm{t}) - \mu(\mathbf{x}, \mathrm{t})) + (\beta_{\mathrm{t}}(\mathbf{x}) - \alpha_{\mathrm{t}}(\mathbf{x})) \int u(y)(y - \mu(\mathbf{x}, \mathrm{t})) f(y \mid \mathbf{x}, \mathrm{t}) dy}{\alpha_{\mathrm{t}}(\mathbf{x}) \int f(y \mid \mathbf{x}, \mathrm{t}) dy + (\beta_{\mathrm{t}}(\mathbf{x}) - \alpha_{\mathrm{t}}(\mathbf{x})) \int u(y) f(y \mid \mathbf{x}, \mathrm{t}) dy} \\
&= \inf_{u \in \mathcal{U}^{ni}} \mu(\mathbf{x}, \mathrm{t}) + \frac{(\beta_{\mathrm{t}}(\mathbf{x}; \Gamma) - \alpha_{\mathrm{t}}(\mathbf{x}; \Gamma)) \int u(y)(y - \mu(\mathbf{x}, \mathrm{t})) f(y \mid \mathbf{x}, \mathrm{t}) dy}{\alpha_{\mathrm{t}}(\mathbf{x}; \Gamma) + (\beta_{\mathrm{t}}(\mathbf{x}; \Gamma) - \alpha_{\mathrm{t}}(\mathbf{x}; \Gamma)) \int u(y) f(y \mid \mathbf{x}, \mathrm{t}) dy} \\
&= \inf_{u \in \mathcal{U}^{ni}} \mu(\mathbf{x}, \mathrm{t}) + \frac{\int u(y)(y - \mu(\mathbf{x}, \mathrm{t})) f(y \mid \mathbf{x}, \mathrm{t}) dy}{\alpha_{\mathrm{t}}'^{\Gamma}(\mathbf{x}) + \int u(y) f(y \mid \mathbf{x}, \mathrm{t}) dy} \\
&= \inf_{y^* \in \mathcal{Y}} \mu(\mathbf{x}, \mathrm{t}) + \frac{\int_{-\infty}^{y^*} (y - \mu(\mathbf{x}, \mathrm{t})) f(y \mid \mathbf{x}, \mathrm{t}) dy}{\alpha_{\mathrm{t}}'^{\Gamma}(\mathbf{x}) + \int_{-\infty}^{y^*} f(y \mid \mathbf{x}, \mathrm{t}) dy} \\
&= \inf_{y^* \in \mathcal{Y}} \mu(\mathbf{x}, \mathrm{t}) + \frac{\int_{-\infty}^{y^*} (y - \mu(\mathbf{x}, \mathrm{t})) f(y \mid \mathbf{x}, \mathrm{t}) dy}{\alpha_{\mathrm{t}}'^{\Gamma}(\mathbf{x}) + \mathrm{P}(Y \le y^* \mid \mathbf{x}, \mathrm{t})} \\
&= \inf_{y^* \in \mathcal{Y}} \mu(\mathbf{x}, \mathrm{t}) + \frac{\int_{-\infty}^{y^*} \mathrm{r}_{\mathrm{t}}(y) f(y \mid \mathbf{x}, \mathrm{t}) dy}{\alpha_{\mathrm{t}}'^{\Gamma}(\mathbf{x}) + \mathrm{P}(Y \le y^* \mid \mathbf{x}, \mathrm{t})}.
\end{aligned}
$$

$\square$

## B. CATE Interval Estimator

*Proof for Theorem 1.* Here we prove that $\widehat{\underline{\mu}}_{\boldsymbol{\omega}}^{\Gamma}(\mathbf{x}, \mathrm{t}) \xrightarrow{p} \underline{\mu}^{\Gamma}(\mathbf{x}, \mathrm{t})$, from which $\widehat{\overline{\mu}}_{\boldsymbol{\omega}}^{\Gamma}(\mathbf{x}, \mathrm{t}) \xrightarrow{p} \overline{\mu}^{\Gamma}(\mathbf{x}, \mathrm{t})$ can be proved analogously. Note that $\xrightarrow{p}$ indicates convergence in probability. As a reminder

$$
\widehat{\underline{\mu}}_{\boldsymbol{\omega}}^{\Gamma}(\mathbf{x}, \mathrm{t}) = \inf_{y^* \in \mathcal{Y}} \widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x}, \mathrm{t}) + \frac{\int_{-\infty}^{y^*} (y - \widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x}, \mathrm{t})) f_{\boldsymbol{\omega}}(y \mid \mathbf{x}, \mathrm{t}) dy}{\alpha_{\boldsymbol{\omega}}'^{\Gamma}(\mathbf{x}, \mathrm{t}) + \int_{-\infty}^{y^*} f_{\boldsymbol{\omega}}(y \mid \mathbf{x}, \mathrm{t}) dy}, \tag{16}
$$

and

$$
\underline{\mu}^{\Gamma}(\mathbf{x}, \mathrm{t}) = \inf_{y^* \in \mathcal{Y}} \mu(\mathbf{x}, \mathrm{t}) + \frac{\int_{-\infty}^{y^*} (y - \mu(\mathbf{x}, \mathrm{t})) f(y \mid \mathbf{x}, \mathrm{t}) dy}{\alpha_{\mathrm{t}}'^{\Gamma}(\mathbf{x}) + \mathrm{P}(Y \le y^* \mid \mathbf{x}, \mathrm{t})}. \tag{17}
$$

Further, our assumptions are

1. $n \to \infty$, and $\mathbf{x} \in \mathcal{D}$.

2. Y is a bounded random variable.

3. $f_{\boldsymbol{\omega}}(y \mid \mathbf{x}, \mathrm{t})$ converges in measure to $f(y \mid \mathbf{x}, \mathrm{t})$. Specifically, $\lim_{n \to \infty} P(\{y \in \mathcal{Y} : |f(y \mid \mathbf{x}, \mathrm{t}) - f_{\boldsymbol{\omega}|\mathcal{D}_n}(y \mid \mathbf{x}, \mathrm{t})| \ge \epsilon\}) = 0$, for every $\epsilon \ge 0$, where $\mathcal{D}_n$ is a dataset of size $n$. Convergence in measure is a generalization of convergence in probability.

4. $\widehat{e}_{\mathrm{t}, \boldsymbol{\omega}}(\mathbf{x})$ and $\widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x})$ are consistent estimators of $\mathbb{E}[\mathrm{T} = \mathrm{t} \mid \mathbf{X} = \mathbf{x}]$ and $\mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}, \mathrm{T} = \mathrm{t}]$.

5. $e_{\mathrm{t}}(\mathbf{x}, y)$ is bounded away from 0 and 1 uniformly over $\mathbf{x} \in \mathcal{X}$, $y \in \mathcal{Y}$, and $\mathrm{t} \in \{0, 1\}$.

We need to show that $\lim_{n \to \infty} P(|\widehat{\underline{\mu}}_{\boldsymbol{\omega}}^{\Gamma}(\mathbf{x}, \mathrm{t}) - \underline{\mu}^{\Gamma}(\mathbf{x}, \mathrm{t})| \ge \epsilon) = 0$, for all $\epsilon > 0$, where the parameters $\boldsymbol{\omega}$ are dependent on the size of the dataset $n$. First, define the following quantities:

$$
\kappa_{\mathrm{y}}^{y^*}(\mathbf{x}, \mathrm{t}; n) = \int_{-\infty}^{y^*} (y - \widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x}, \mathrm{t})) f_{\boldsymbol{\omega}}(y \mid \mathbf{x}, \mathrm{t}) dy, \quad I_{\mathrm{y}}^{y^*}(\mathbf{x}, \mathrm{t}) = \int_{-\infty}^{y^*} (y - \mu(\mathbf{x}, \mathrm{t})) f(y \mid \mathbf{x}, \mathrm{t}) dy,
$$

$$
\kappa^{y^*}(\mathbf{x}, \mathrm{t}; n) = \int_{-\infty}^{y^*} f_{\boldsymbol{\omega}}(y \mid \mathbf{x}, \mathrm{t}) dy, \qquad\qquad I^{y^*}(\mathbf{x}, \mathrm{t}) = \int_{-\infty}^{y^*} f(y \mid \mathbf{x}, \mathrm{t}) dy,
$$

so that

$$\widehat{\underline{\mu}}_{\boldsymbol{\omega}}^{\Gamma}(\mathbf{x}, t) = \inf_{y^* \in \mathcal{Y}} \widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x}, t) + \frac{\kappa_y^{y^*}(\mathbf{x}, t; n)}{\alpha_{\boldsymbol{\omega}}'^{\Gamma}(\mathbf{x}, t) + \kappa^{y^*}(\mathbf{x}, t; n)}, \quad \underline{\mu}^{\Gamma}(\mathbf{x}, t) = \inf_{y^* \in \mathcal{Y}} \mu(\mathbf{x}, t) + \frac{I_y^{y^*}(\mathbf{x}, t)}{\alpha_t'^{\Gamma}(\mathbf{x}) + I^{y^*}(\mathbf{x}, t)}$$

To start, we want to express $|\widehat{\underline{\mu}}_{\boldsymbol{\omega}}^{\Gamma}(\mathbf{x}, t) - \underline{\mu}^{\Gamma}(\mathbf{x}, t)|$ in terms of the following 4 terms: $\Delta_1(\mathbf{x}, t; n) = |\widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x}, t) - \mu(\mathbf{x}, t)|$; $\Delta_2(\mathbf{x}, t; n) = |\alpha_{\boldsymbol{\omega}}'^{\Gamma} - \alpha_t'^{\Gamma}|$; $\Delta_3(\mathbf{x}, t; n) = \sup_{y^* \in \mathcal{Y}} |\delta_3|$, where $\delta_3 = \kappa_y^{y^*} - I_y^{y^*}$; and $\Delta_4(\mathbf{x}, t; n) = \sup_{y^* \in \mathcal{Y}} |\delta_4|$, where $\delta_4 = \kappa^{y^*} - I^{y^*}$. To this end, we use Lemma 3 from Kallus et al. (2019) in line 1 below and define the following inequality:

$$|\widehat{\underline{\mu}}_{\boldsymbol{\omega}}^{\Gamma}(\mathbf{x}, t) - \underline{\mu}^{\Gamma}(\mathbf{x}, t)| \leq \sup_{y^* \in \mathcal{Y}} \left| \widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x}, t) - \mu(\mathbf{x}, t) + \frac{\kappa_y^{y^*}(\mathbf{x}, t; n)}{\alpha_{\boldsymbol{\omega}}'^{\Gamma}(\mathbf{x}, t) + \kappa^{y^*}(\mathbf{x}, t; n)} - \frac{I_y^{y^*}(\mathbf{x}, t)}{\alpha_t'^{\Gamma}(\mathbf{x}) + I^{y^*}(\mathbf{x}, t)} \right|,$$

$$\leq \Delta_1(\mathbf{x}, t; n) + \sup_{y^* \in \mathcal{Y}} \left| \frac{\kappa_y^{y^*}}{\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}} - \frac{I_y^{y^*}}{\alpha_t'^{\Gamma} + I^{y^*}} \right|,$$

$$\leq \Delta_1(\mathbf{x}, t; n) + \sup_{y^* \in \mathcal{Y}} \left\{ |\kappa_y^{y^*}| \frac{|(\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}) - (\alpha_t'^{\Gamma} + I^{y^*})|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} - \frac{|\kappa_y^{y^*} - I_y^{y^*}|}{|\alpha_t'^{\Gamma} + I^{y^*}|} \right\},$$

$$\leq \Delta_1(\mathbf{x}, t; n) + \sup_{y^* \in \mathcal{Y}} \left\{ \frac{|\kappa_y^{y^*}||\alpha_{\boldsymbol{\omega}}'^{\Gamma} - \alpha_t'^{\Gamma}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \frac{|\kappa_y^{y^*}||\kappa^{y^*} - I^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \frac{|\kappa_y^{y^*} - I_y^{y^*}|}{|\alpha_t'^{\Gamma} + I^{y^*}|} \right\},$$

$$\leq \Delta_1(\mathbf{x}, t; n) + \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*}||\alpha_{\boldsymbol{\omega}}'^{\Gamma} - \alpha_t'^{\Gamma}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*}||\kappa^{y^*} - I^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*} - I_y^{y^*}|}{|\alpha_t'^{\Gamma} + I^{y^*}|},$$

$$= \Delta_1(\mathbf{x}, t; n) + |\alpha_{\boldsymbol{\omega}}'^{\Gamma} - \alpha_t'^{\Gamma}| \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*}||\kappa^{y^*} - I^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*} - I_y^{y^*}|}{|\alpha_t'^{\Gamma} + I^{y^*}|},$$

$$= \Delta_1(\mathbf{x}, t; n) + \Delta_2(\mathbf{x}, t; n) \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*}||\kappa^{y^*} - I^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*} - I_y^{y^*}|}{|\alpha_t'^{\Gamma} + I^{y^*}|},$$

$$\leq \Delta_1(\mathbf{x}, t; n) + \Delta_2(\mathbf{x}, t; n) \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*} - I_y^{y^*} + I_y^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*} - I_y^{y^*} + I_y^{y^*}||\kappa^{y^*} - I^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\kappa_y^{y^*} - I_y^{y^*}|}{|\alpha_t'^{\Gamma} + I^{y^*}|},$$

$$= \Delta_1(\mathbf{x}, t; n) + \Delta_2(\mathbf{x}, t; n) \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3 + I_y^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3 + I_y^{y^*}||\kappa^{y^*} - I^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma} + I^{y^*}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3|}{|\alpha_t'^{\Gamma} + I^{y^*}|},$$

$$\leq \Delta_1(\mathbf{x}, t; n) + \Delta_2(\mathbf{x}, t; n) \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3 + I_y^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3 + I_y^{y^*}||\kappa^{y^*} - I^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3|}{|\alpha_t'^{\Gamma}|},$$

$$\leq \Delta_1(\mathbf{x}, t; n) + \Delta_2(\mathbf{x}, t; n) \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3| + |I_y^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma}|} + \sup_{y^* \in \mathcal{Y}} \frac{(|\delta_3| + |I_y^{y^*}|)|\kappa^{y^*} - I^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3|}{|\alpha_t'^{\Gamma}|},$$

$$\leq \Delta_1(\mathbf{x}, t; n) + \Delta_2(\mathbf{x}, t; n) \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3||\kappa^{y^*} - I^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma} + \kappa^{y^*}||\alpha_t'^{\Gamma}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3|}{|\alpha_t'^{\Gamma}|},$$

$$\leq \Delta_1(\mathbf{x}, t; n) + \Delta_2(\mathbf{x}, t; n) \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma}||\alpha_t'^{\Gamma}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3||\kappa^{y^*} - I^{y^*}|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma}||\alpha_t'^{\Gamma}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3|}{|\alpha_t'^{\Gamma}|},$$

$$\leq \Delta_1(\mathbf{x}, t; n) + \Delta_2(\mathbf{x}, t; n) \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma}||\alpha_t'^{\Gamma}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3||\delta_4|}{|\alpha_{\boldsymbol{\omega}}'^{\Gamma}||\alpha_t'^{\Gamma}|} + \sup_{y^* \in \mathcal{Y}} \frac{|\delta_3|}{|\alpha_t'^{\Gamma}|},$$

$$= \Delta_1(\mathbf{x}, t; n) + \frac{\Delta_2(\mathbf{x}, t; n)\Delta_3(\mathbf{x}, t; n)}{\alpha_{\boldsymbol{\omega}}'^{\Gamma}\alpha_t'^{\Gamma}} + \frac{\Delta_3(\mathbf{x}, t; n)\Delta_4(\mathbf{x}, t; n)}{\alpha_{\boldsymbol{\omega}}'^{\Gamma}\alpha_t'^{\Gamma}} + \frac{\Delta_3(\mathbf{x}, t; n)}{\alpha_t'^{\Gamma}}.$$

So, we now need only prove that $\Delta_1(\mathbf{x}, t; n) \xrightarrow{p} 0$, $\Delta_2(\mathbf{x}, t; n) \xrightarrow{p} 0$, $\Delta_3(\mathbf{x}, t; n) \xrightarrow{p} 0$, and $\Delta_4(\mathbf{x}, t; n) \xrightarrow{p} 0$, when $n \to \infty$. Note that both $\Delta_1(\mathbf{x}, t; n) \xrightarrow{p} 0$ and $\Delta_2(\mathbf{x}, t; n) \xrightarrow{p} 0$ are covered by Assumption 4 of Theorem 1; namely, $\widehat{e}_{t, \boldsymbol{\omega}}(\mathbf{x})$ and $\widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x})$ are consistent estimators of $\mathbb{E}[T = t \mid \mathbf{X} = \mathbf{x}]$ and $\mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}, T = t]$.

First, we prove that $\Delta_4(\mathbf{x}, t; n) \xrightarrow{p} 0$.

**Prove that** $\sup_{y^* \in \mathcal{Y}} \left| \kappa^{y^*} - I^{y^*} \right| \xrightarrow{p} 0$

$$\sup_{y^* \in \mathcal{Y}} \left| \kappa^{y^*} - I^{y^*} \right| = \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} f_{\boldsymbol{\omega}}(y \mid \mathbf{x}, t) dy - \int_{-\infty}^{y^*} f(y \mid \mathbf{x}, t) dy \right|$$

$$= \sup_{y^* \in \mathcal{Y}} \left| P_{\boldsymbol{\omega}}(y \leq y^* \mid \mathbf{x}, t) - P(y \leq y^* \mid \mathbf{x}, t) \right|$$

Convergence in probability implies convergence in distribution ($\lim_{n \to \infty} P_n(\mathbf{X} \leq \mathbf{x}) = P(\mathbf{X} \leq \mathbf{x})$), so by Assumption 3 in Theorem 1

$$\lim_{n \to \infty} P \left( \left| \sup_{y^* \in \mathcal{Y}} \left| \kappa^{y^*} - I^{y^*} \right| \right| \geq \epsilon \right) = \lim_{n \to \infty} P \left( \sup_{y^* \in \mathcal{Y}} \left| P_{\boldsymbol{\omega}}(y \leq y^* \mid \mathbf{x}, t) - P(y \leq y^* \mid \mathbf{x}, t) \right| \geq \epsilon \right)$$

$$= P \left( \sup_{y^* \in \mathcal{Y}} \left| P(y \leq y^* \mid \mathbf{x}, t) - P(y \leq y^* \mid \mathbf{x}, t) \right| \geq \epsilon \right)$$

$$= P \left( \sup_{y^* \in \mathcal{Y}} |0| \geq \epsilon \right)$$

$$= P \left( 0 \geq \epsilon \right)$$

$$= 0$$

Finally, we prove $\Delta_3(\mathbf{x}, t; n) \xrightarrow{p} 0$.

**Prove that** $\sup_{y^* \in \mathcal{Y}} \left| \kappa_y^{y^*} - I_y^{y^*} \right| \xrightarrow{p} 0$

$$\sup_{y^* \in \mathcal{Y}} |\kappa_y^{y^*} - I_y^{y^*}| = \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} (y - \widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x}, t)) f_{\boldsymbol{\omega}}(y \mid \mathbf{x}, t) dy - \int_{-\infty}^{y^*} (y - \mu(\mathbf{x}, t)) f(y \mid \mathbf{x}, t) dy \right|,$$

$$= \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} y f_{\boldsymbol{\omega}}(y \mid \mathbf{x}, t) dy - \int_{-\infty}^{y^*} y f(y \mid \mathbf{x}, t) dy + \mu(\mathbf{x}, t) \int_{-\infty}^{y^*} f(y \mid \mathbf{x}, t) dy - \widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x}, t) \int_{-\infty}^{y^*} f_{\boldsymbol{\omega}}(y \mid \mathbf{x}, t) dy \right|,$$

$$= \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} y f_{\boldsymbol{\omega}}(y) dy - \int_{-\infty}^{y^*} y f(y) dy + \mu \int_{-\infty}^{y^*} f(y) dy - \widehat{\mu}_{\boldsymbol{\omega}} \int_{-\infty}^{y^*} f_{\boldsymbol{\omega}}(y) dy \right|,$$

$$= \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} y f_{\boldsymbol{\omega}}(y) dy - \int_{-\infty}^{y^*} y f(y) dy + (\mu - \widehat{\mu}_{\boldsymbol{\omega}} + \widehat{\mu}_{\boldsymbol{\omega}}) \int_{-\infty}^{y^*} f(y) dy - \widehat{\mu}_{\boldsymbol{\omega}} \int_{-\infty}^{y^*} (f_{\boldsymbol{\omega}}(y) - f(y) + f(y)) dy \right|,$$

$$= \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} y f_{\boldsymbol{\omega}}(y) dy - \int_{-\infty}^{y^*} y f(y) dy + (\mu - \widehat{\mu}_{\boldsymbol{\omega}}) \int_{-\infty}^{y^*} f(y) dy - \widehat{\mu}_{\boldsymbol{\omega}} \int_{-\infty}^{y^*} (f_{\boldsymbol{\omega}}(y) - f(y)) dy \right|,$$

$$= \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} y f_{\boldsymbol{\omega}}(y) dy - \int_{-\infty}^{y^*} y f(y) dy + (\mu - \widehat{\mu}_{\boldsymbol{\omega}}) \int_{-\infty}^{y^*} f(y) dy - \widehat{\mu}_{\boldsymbol{\omega}} (P_{\boldsymbol{\omega}}(y \leq y^* \mid \mathbf{x}, t) - P(y \leq y^* \mid \mathbf{x}, t)) \right|.$$

As a first step, we can use the result for $\Delta_4(\mathbf{x}, t; n)$ to remove the green term from the supremum and now we need to show that

$$\lim_{n \to \infty} P \left( \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} y f_{\boldsymbol{\omega}}(y) dy - \int_{-\infty}^{y^*} y f(y) dy + (\mu - \widehat{\mu}_{\boldsymbol{\omega}}) \int_{-\infty}^{y^*} f(y) dy \right| \geq \epsilon. \right) = 0$$

Next, under assumption 4 of Theorem 1 we have $\mu(\mathbf{x}, t) - \widehat{\mu}_{\boldsymbol{\omega}}(\mathbf{x}, t) \xrightarrow{p} 0$, and we are left finally to show that

$$\lim_{n \to \infty} P \left( \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} y f_{\boldsymbol{\omega}}(y) dy - \int_{-\infty}^{y^*} y f(y) dy \right| \geq \epsilon. \right) = 0$$

Assumption 2 of Theorem 1 states that Y is a bounded random variable. As such, there exists a $g(y)$ such that $|y f_{\boldsymbol{\omega}}(y)| \leq g(y)$ for all $n$ and $y \in \mathcal{Y}$. Therefore, in conjunction with Assumption 3, by Lebesgue's dominated convergence theorem we have $\lim_{n \to \infty} \int_{-\infty}^{y^*} y f_{\boldsymbol{\omega}}(y) dy = \int_{-\infty}^{y^*} y f(y) dy$

$$\lim_{n \to \infty} P \left( \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} y f_{\boldsymbol{\omega}}(y) dy - \int_{-\infty}^{y^*} y f(y) dy \right| \geq \epsilon. \right) = P \left( \sup_{y^* \in \mathcal{Y}} \left| \lim_{n \to \infty} \int_{-\infty}^{y^*} y f_{\boldsymbol{\omega}}(y) dy - \int_{-\infty}^{y^*} y f(y) dy \right| \geq \epsilon. \right)$$

$$= P \left( \sup_{y^* \in \mathcal{Y}} \left| \int_{-\infty}^{y^*} y f(y) dy - \int_{-\infty}^{y^*} y f(y) dy \right| \geq \epsilon. \right)$$

$$= P \left( \sup_{y^* \in \mathcal{Y}} |0| \geq \epsilon. \right)$$

$$= 0$$

Therefore, $\widehat{\underline{\mu}}_{\boldsymbol{\omega}}^{\Gamma}(\mathbf{x}, t) \xrightarrow{P} \underline{\mu}^{\Gamma}(\mathbf{x}, t)$, and $\widehat{\overline{\mu}}_{\boldsymbol{\omega}}^{\Gamma}(\mathbf{x}, t) \xrightarrow{P} \overline{\mu}^{\Gamma}(\mathbf{x}, t)$ can be proved analogously, which concludes our proof that both $\widehat{\underline{\tau}}_{\boldsymbol{\omega}}(\mathbf{x}; \Gamma) \xrightarrow{P} \underline{\tau}(\mathbf{x}; \Gamma)$, and $\widehat{\overline{\tau}}_{\boldsymbol{\omega}}(\mathbf{x}; \Gamma) \xrightarrow{P} \overline{\tau}(\mathbf{x}; \Gamma)$. $\square$

## C. Datasets

### C.1. Simulated Data

The simulated dataset presented by Kallus et al. (2019) is described by the following structural causal model (SCM):

$$u := N_u, \tag{22a}$$
$$x := N_x, \tag{22b}$$
$$t := N_t, \tag{22c}$$
$$y := (2t - 1)x + (2t - 1) - 2\sin(2(2t - 1)x) - 2(2u - 1)(1 + 0.5x) + N_y, \tag{22d}$$

where $N_u \sim \text{Bern}(0.5)$, $N_x \sim \text{Unif}[-2, 2]$, $N_u \perp\!\!\!\perp N_x$, $N_t \sim \text{Bern}(e(x, u))$, $e(x, u) = \frac{u}{\alpha_t(x;\Gamma^*)} + \frac{1-u}{\beta_t(x;\Gamma^*)}$, $e(x) = \text{sigmoid}(0.75x + 0.5)$, and $N_y \sim \mathcal{N}(0, 1)$.

Remember that only x, t, and y are observed. So the bias induced by hidden confounding at x is given by

$$\tilde{\tau}(x) - \tau(x) = 2(2 + x) \left( P(u = 1 \mid T = 1, X = x) - P(u = 1 \mid T = 0, X = x) \right), \tag{23}$$

where $\tilde{\tau}(x)$ is the confounded CATE estimate.

Each random realization of the simulated dataset generates 1000 training examples, 100 validation examples, and 1000 test examples. In the experiments we report results over 50 random realizations. The seeds for the random number generators are $i$, $i + 1$, and $i + 2$; $\{i \in [0, 1, \ldots, 49]\}$, for the training, validation, and test sets, respectively. Code is available in file /library/datasets/synthetic.py on github at `https://github.com/anndvision/quince`.

### C.2. HC-MNIST

HC-MNIST is an extension of the above dataset with high-dimensional covariates $\mathbf{x}$. Specifically, $\mathbf{x}$ are MNIST digits. HC-MNIST is described by the following SCM:

$$u := N_u, \tag{24a}$$
$$\mathbf{x} := N_{\mathbf{x}}, \tag{24b}$$
$$\phi := \left( \text{clip} \left( \frac{\mu_{N_{\mathbf{x}}} - \mu_c}{\sigma_c}; -1.4, 1.4 \right) - \text{Min}_c \right) \frac{\text{Max}_c - \text{Min}_c}{1.4 - \text{-}1.4} \tag{24c}$$
$$t := N_t, \tag{24d}$$
$$y := (2t - 1)\phi + (2t - 1) - 2\sin(2(2t - 1)\phi) - 2(2u - 1)(1 + 0.5\phi) + N_y, \tag{24e}$$

where $N_u$, $N_t$ (swapping x for $\phi$), and $N_y$ are as described in Appendix C.1. $N_{\mathbf{x}}$ is a sample of an MNIST image. The sampled image has a corresponding label $c \in [0, \ldots, 9]$. $\mu_{N_{\mathbf{x}}}$ is the average intensity of the sampled image. $\mu_c$ and $\sigma_c$ are the mean and standard deviation of the average image intensities over all images with label c in the MNIST training set.

In other words, $\mu_{\rm c} = \mathbb{E}[\mu_{N_{\rm x}} \mid {\rm c}]$ and $\sigma_{\rm c}^2 = {\rm Var}[\mu_{N_{\rm x}} \mid {\rm c}]$. To map the high dimensional images $\mathbf{x}$ onto a one-dimensional manifold $\phi$ with the same domain as ${\rm x} \in [-2, 2]$ above, we first clip the standardized average image intensity on the range $(-1.4, 1.4)$. Each digit class has its own domain in $\phi$, so there is a linear transformation of the clipped value onto the range $[{\rm Min}_{\rm c}, {\rm Max}_{\rm c}]$. Finally, ${\rm Min}_{\rm c} = -2 + \frac{4}{10}{\rm c}$, and ${\rm Max}_{\rm c} = -2 + \frac{4}{10}({\rm c} + 1)$.

For each random realization of the dataset, the MNIST training set is split into training ($n = 35000$) and validation ($n = 15000$) subsets using the scikit-learn function train_test_split(). The test set is generated using the MNIST test set ($n = 10000$). The random seeds are $\{i \in [0, 1, \ldots, 19]\}$ for the 20 random realizations generated. Code to generate this dataset is available in file /library/datasets/hcmnist.py on github at `https://github.com/anndvision/quince`.

## C.3. IHDP Hidden Confounding

The experimental data from the Infant Health and Development Program (IHDP) are used by Hill (2011) to generate simulated outcomes. The treatment group reveives "intensive high-quality child care and home visits from a trained provider." Hill (2011) uses "measurements on the child–birth weight, head circumference, weeks born preterm, birth order, first born, neonatal health index, sex, twin status–as well as behaviors engaged in during pregnancy–smoked cigarettes, drank alcohol, took drugs–and measurements on the mother at the time she gave birth–age, marital status, educational attainment (did not graduate from high school, graduated from high school, attended some college but did not graduate, graduated from college), whether she worked during pregnancy, whether she received prenatal care–and the site (8 total) in which the family resided at the start of the intervention. There are 6 continuous covariates and 19 binary covariates." Hill (2011) excludes "a nonrandom portion of the treatment group: all children with nonwhite mothers," in order to simulate an observational study. Table 3 enumerates the included covariates. There are 139 examples in the treatment group and 608 examples in the control group, for a total of 747 examples.

*Table 3.* **IHDP Covariates** Binary covariates $x_9 - x_{18}$ are attributes of the mother. Mother's education level "College" indicated by covariates $x_{10} - x_{12}$ all zero. Site 8 indicated by covariates $x_{19} - x_{25}$ all zero. We show the frequency of occurrence for each binary covariate $p({\rm x} = 1)$, as well as the adjusted mutual information $I({\rm x}; {\rm t})$ between the binary covariate and the treatment variable.

| Continuous Covariate | Description | Binary Covariate | Description | $I({\rm x}; {\rm t})$ | $p({\rm x} = 1)$ |
|---|---|---|---|---|---|
| $x_1$ | birthweight | $x_7$ | child's gender (female=1) | 0.00 | 0.51 |
| $x_2$ | head circumference | $x_8$ | is child a twin | 0.00 | 0.09 |
| $x_3$ | number of weeks pre-term | $x_9$ | married when child born | **0.02** | **0.52** |
| $x_4$ | birth order | $x_{10}$ | left High School | 0.00 | 0.36 |
| $x_5$ | "neo-natal health index" | $x_{11}$ | completed High School | 0.00 | 0.27 |
| $x_6$ | mom's age | $x_{12}$ | some College | 0.00 | 0.22 |
| | | $x_{13}$ | child is first born | 0.00 | 0.36 |
| | | $x_{14}$ | smoked cigarettes when pregnant | **0.01** | **0.48** |
| | | $x_{15}$ | consumed alcohol when pregnant | 0.00 | 0.14 |
| | | $x_{16}$ | used drugs when pregnant | 0.00 | 0.96 |
| | | $x_{17}$ | worked during pregnancy | **0.01** | **0.59** |
| | | $x_{18}$ | received any prenatal care | **0.01** | 0.96 |
| | | $x_{19}$ | site 1 | 0.00 | 0.14 |
| | | $x_{20}$ | site 2 | **0.01** | 0.14 |
| | | $x_{21}$ | site 3 | 0.00 | 0.16 |
| | | $x_{22}$ | site 4 | **0.01** | 0.08 |
| | | $x_{23}$ | site 5 | **0.02** | 0.07 |
| | | $x_{24}$ | site 6 | **0.01** | 0.13 |
| | | $x_{25}$ | site 7 | **0.02** | 0.16 |

Response surface B, designed by Hill (2011), is described by the following SCM:

$$\mathbf{x} := N_{\mathbf{x}}, \tag{25a}$$

$$\mathrm{t} := N_{\mathrm{t}}, \tag{25b}$$

$$\mathrm{y} := (\mathrm{t} - 1)\left(\exp(\beta_{\mathbf{x}}(\mathbf{x} + \mathbf{w})) + N_{\mathrm{Y}^0}\right) + \mathrm{t}\left(\beta_{\mathbf{x}}\mathbf{x} - \omega^s + N_{\mathrm{Y}^1}\right), \tag{25c}$$

where $(N_{\mathbf{x}}, N_{\mathrm{t}}) \sim p_{\mathcal{D}}(\{\mathrm{x}_1, \ldots \mathrm{x}_{25}\}, \mathrm{t})$, $N_{\mathrm{Y}^0} \sim \mathcal{N}(0, 1)$, and $N_{\mathrm{Y}^1} \sim \mathcal{N}(0, 1)$. The coefficients $\beta_{\mathbf{x}}$ are a vector of randomly sampled values $(0.0, 0.1, 0.2, 0.3, 0.4)$ with probabilities $(0.6, 0.1, 0.1, 0.1, 0.1)$. Hill (2011) describes $\omega^s$ as follows: "For the $s$th simulation, $[\omega^s]$ is chosen in the overlap setting, where we estimate the effect of the treatment on the treated [(CATT)], such that CATT equals 4; similarly it was chosen in the incomplete setting, where we estimate the effect of the treatment on the controls [(CATC)], so that CATC equals 4." An offset vector $\mathbf{w}$, equal in dimension to $\mathbf{x}$, with every value set to 0.5, is added to $\mathbf{x}$.

To induce hidden confounding, we need to select a variable $\mathrm{u}$ that is associated with the treatment that will be hidden from the CATE interval estimator, and design a response surface where the outcome will always be affected by $\mathrm{u}$. In Table 3, we list 3 potential candidates for $\mathrm{u}$: $\mathrm{x}_9$, $\mathrm{x}_{14}$, and $\mathrm{x}_{17}$. Each of these variable have a non-negligible association with the treatment, as indicated by the adjusted mutual information score $I(\mathrm{x}; \mathrm{t})$, and have a frequency of taking the value 1 at around 0.5 (increasing the chances that we will have both positive and negative examples in each of the training, validation, and testing splits). Here we select $\mathrm{x}_9$ and define the following SCM:

$$\mathrm{u} := N_{\mathrm{u}}, \tag{26a}$$

$$\mathbf{x} := N_{\mathbf{x}}, \tag{26b}$$

$$\mathrm{t} := N_{\mathrm{t}}, \tag{26c}$$

$$\mathrm{y} := (\mathrm{t} - 1)(\exp(\beta_{\mathbf{x}}(\mathbf{x} + \mathbf{w}) + \beta_{\mathrm{u}}(\mathrm{u} + 0.5)) + N_{\mathrm{Y}^0}) + \mathrm{t}(\beta_{\mathbf{x}}\mathbf{x} + \beta_{\mathrm{u}}\mathrm{u} - \omega^s + N_{\mathrm{Y}^1}), \tag{26d}$$

where $(N_{\mathrm{u}}, N_{\mathbf{x}}, N_{\mathrm{t}}) \sim p_{\mathcal{D}}(\mathrm{x}_9, \{\mathrm{x}_1, \ldots, \mathrm{x}_8, \mathrm{x}_{10}, \ldots, \mathrm{x}_{25}\}, \mathrm{t})$, $N_{\mathrm{Y}^0} \sim \mathcal{N}(0, 1)$, and $N_{\mathrm{Y}^1} \sim \mathcal{N}(0, 1)$. The coefficient $\beta_{\mathrm{u}}$ is randomly sampled from $(0.1, 0.2, 0.3, 0.4, 0.5)$ with probabilities $(0.2, 0.2, 0.2, 0.2, 0.2)$. The remaining parameters–$\beta_{\mathbf{x}}$, $\omega^s$, and $\boldsymbol{\omega}$–are given as above, taking into account $\mathrm{u}$.

For each random realization of the dataset, the IHDP data is split into training ($n = 470$), validation ($n = 202$) and test ($n = 75$) subsets using the scikit-learn function train_test_split(). The random seeds for both splitting and outcome generation are $\{i \in [0, 1, \ldots, 999]\}$ for the 1000 realizations generated. Code to generate this dataset is available in file /library/datasets/ihdp.py on github at `https://github.com/anndvision/quince`.

## D. Implementation Details

Experiments for the Simulated and IHDP datasets were run using a single NVIDIA GeForce GTX 1080 ti, an Intel(R) Core(TM) i7-8700K, on a desktop computer with 16GB of RAM. Experiments for the HCMNIST dataset were run using 4 NVIDIA GeForce RTX 2080 ti GPUs, an Intel(R) Core(TM) i9-9900K, on a server with 64GB of RAM. Code is written in python. Packages used include PyTorch (Paszke et al., 2019), scikit-learn (Pedregosa et al., 2011), Ray (Moritz et al., 2018), NumPy, SciPy, and Matplotlib. We use ray tune (Liaw et al., 2018) with the hyperopt (Bergstra et al., 2013) search algorithm to optimize our network hyper-parameters. The hyper-parameters we consider are accounted for in Table 4. The hyper-paramter optimization objective for each dataset is the expected batch-wise log-likelihood of the validation data for a single dataset realization with random seed 1331.

Each experiment is replicated using the training, validation, and testing datasets described in the previous section.

Code to replicate these experiments is available at `https://github.com/anndvision/quince`.

### D.1. Simulated Data

As a reminder, we need parametric models for the distribution over outcomes $p_{\boldsymbol{\omega}}(Y \mid \mathbf{x}, \mathrm{t})$ and the nominal propensity $\widehat{e}_{\boldsymbol{\omega}}(\mathbf{x})$. For $p_{\boldsymbol{\omega}}(Y \mid \mathbf{x}, \mathrm{t})$, we use a 4 hidden layer mixture density network (MDN) (Bishop, 1994) with 5 mixture components. The 1D treatment variable $\mathrm{t}$ and 1D covariate $\mathbf{x}$ are concatenated to make a 2D network input. Each hidden layer is comprised of a 100 neuron linear transformation, followed by a ReLU activation function. The MDN parameters are inferred using a linear layer to predict the 5 means, a linear layer followed by a softplus activation to predict the square root of the 5 variances, and

| Hyper-parameter | Search Space |
|---|---|
| hidden units | [50, 100, 200, 400] |
| network depth | [1, 2, 3, 4, 5] |
| negative slope | [ReLU, 0.1, 0.2, 0.3, 0.4, 0.5, ELU] |
| dropout rate | [0.00, 0.10, 0.15, 0.20, 0.25, 0.50] |
| spectral norm | [0.95, 1.0, 2.5, 3.0, 6.0, 12.0, 24.0] |
| batch size | [16, 32, 64, 100, 200] |
| learning rate | [5e-5, 1e-4, 2e-4, 5e-4, 1e-3] |

*Table 4.* Hyper-parameter search space

| Hyper-parameter | Simulated | HCMNIST | IHDP |
|---|---|---|---|
| hidden units | 200 | 200 | 200 |
| network depth | 4 | 2 | 4 |
| negative slope | ReLU | ReLU | LeakyReLU 0.3 |
| dropout rate | 0.10 | 0.15 | 0.5 |
| spectral norm | 6.0 | 3.0 | 6.0 |
| batch size | 32 | 200 | 200 |
| learning rate | 1e-3 | 5e-4 | 5e-4 |

*Table 5.* Final hyper-parameters for each dataset

a linear layer to predict the logits of the 5 mixture components. We use the pytorch MixtureSameFamily distribution, with mixture_distribution=Categorical(.), and component_distribution=Normal(.) (Paszke et al., 2019). The objective function optimized is the negative log likelihood for the label $y$ of the above distribution with the parameters predicted from $(x, t)$. Dropout is applied to the inputs of each layer after the input layer with a rate of 0.1. Spectral normalization is applied to the weights of the networks with value 6.0. For $\hat{e}_{\omega}(\mathbf{x})$, we use a 4 hidden layer neural network with Bernoulli likelihood. Each hidden layer is comprised of a 200 neuron linear transformation, followed by a ReLU activation function. Spectral normalization is applied to the weights of the networks with value 6.0. The objective function optimized is the negative log likelihood for the observed treatment $t$ of the Bernoulli distribution with the logits predicted from $x$. For both models, We use Adam optimization with default pytorch parameters (Kingma & Ba, 2017). We use a batch size of 32. We use early stopping based on the objective function value on the validation set with a patience of 20 epochs and train for a maximum of 500 epochs. We train an ensemble of 10 models as an estimation of Bayesian model averaging. At test time, we do 10 MC samples, corresponding to a forward pass of each model in the ensemble for $\omega$ and 100 MC samples for $y$, for each model under $t = 0$ and $t = 1$.

**Hyper-parameter selection**    The hyper parameter search space is given in Table 4 and a summary of the final hyper-parameters used are given in Table 5 under the column Simulated. Because the hidden confounding is a binary variable, it induces a bi-modal distribution in $y$ at $\mathbf{x}$, as shown in Figure 5. In practice, we would not know the form of the distribution of $y$ at $\mathbf{x}$. To this end we select 5 mixture components for the MDN to show that we can over estimate the true modality, and still obtain sensible results. Alternatively, the validation set could be used to find the number of components that minimizes negative log likelihood of the data. The number of MC samples are chosen based on the stability of network predictions, i.e. we increase the number of MC samples until the variances with respect to $\omega$ or $y$ no longer change significantly.

### D.2. HC-MNIST

For $p_{\omega}(Y \mid \mathbf{x}, t)$, we use a ResNet CNN feature extractor with 2 residual blocks, followed by a 2 hidden layer MDN with 5 mixture components. The 1D treatment variable $t$ and ResNet output are concatenated to make a 49 dimensional MDN input. Each hidden layer in the MDN is comprised of a 200 neuron linear transformation, followed by a ReLU activation function. The MDN parameters are inferred using in the same manner as for the Simulated data above. The objective function optimized is the negative log likelihood for the label $y$ of the above distribution with the parameters predicted from $(x, t)$. Dropout is applied to the inputs of each layer after the input layer with a rate of 0.15. Spectral normalization is applied to the weights of the network with value 3.0. For $\hat{e}_{\omega}(\mathbf{x})$, we use a ResNet CNN feature extractor with 2 residual
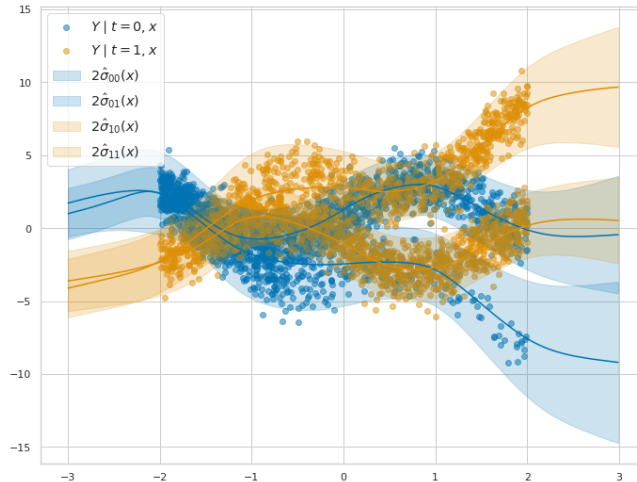
*Figure 5.* **Hidden confounding induces a multi-modal distribution in** $y$ **at** $x$

blocks, followed by a 2 hidden layer neural network with Bernoulli likelihood. Each hidden layer of the Neural Network is comprised of a 200 neuron linear transformation, followed by a ReLU activation function. Dropout is applied to the inputs of each layer after the input layer with a rate of 0.15. Spectral normalization is applied to the weights of the network with value 3.0. The objective function optimized is the negative log likelihood for the observed treatment $t$ of the Bernoulli distribution with the logits predicted from $x$. For both models, We use Adam optimization with a learning rate of 5e-4 (Kingma & Ba, 2017). We use a batch size of 200. We use early stopping based on the objective function value on the validation set with a patience of 20 epochs and train for a maximum of 500 epochs. We train an ensemble of 5 models as an estimation of Bayesian model averaging. At test time, we do 5 MC samples, corresponding to a forward pass of each model in the ensemble for $\omega$ and 100 MC samples for $y$, for each model under $t = 0$ and $t = 1$.

### D.3. IHDP Hidden Confounding

For $p_{\omega}(Y \mid \mathbf{x}, t)$, we use a neural network feature extractor with 4 hidden layers, followed by a 2 hidden layer MDN with 5 mixture components. The 1D treatment variable $t$ and feature extractor output are concatenated to make a 201 dimensional MDN input. Each hidden layer in the feature extractor and MDN is comprised of a 200 neuron linear transformation, followed by an LeakyReLU activation function. The MDN parameters are inferred using in the same manner as for the Simulated data above. The objective function optimized is the negative log likelihood for the label $y$ of the above distribution with the parameters predicted from $(x, t)$. Dropout is applied to the inputs of each layer after the input layer with a rate of 0.5. Spectral normalization is applied to the weights of the network with value 6.0. For $\widehat{e}_{\omega}(\mathbf{x})$, we use a neural network feature extractor with 3 hidden layers, followed by a 2 hidden layer neural network with Bernoulli likelihood. Each hidden layer of the Neural Network is comprised of a 200 neuron linear transformation, followed by a ELU activation function. Dropout is applied to the inputs of each layer after the input layer with a rate of 0.5. Spectral normalization is applied to the weights of the network with value 6.0. The objective function optimized is the negative log likelihood for the observed treatment $t$ of the Bernoulli distribution with the logits predicted from $x$. For both models, We use Adam optimization with a learning rate of 5e-4 (Kingma & Ba, 2017). We use a batch size of 200. We use early stopping based on the objective function value on the validation set with a patience of 20 epochs and train for a maximum of 500 epochs. We train an ensemble of 10 models as an estimation of Bayesian model averaging. At test time, we do 10 MC samples, corresponding to a forward pass of each model in the ensemble for $\omega$ and 100 MC samples for $y$, for each model under $t = 0$ and