

## A Breaking down the Continuous Treatment Marginal Sensitivity Model

Let's go deeper into the Continuous Treatment Marginal Sensitivity Model (CMSM).

### A.1 MSM for binary treatment values

This section details the Marginal Sensitivity Model of [Tan06]. For binary treatments,  $\mathcal{T}_B = \{0, 1\}$ , the (nominal) propensity score,  $e(\mathbf{x}) \equiv \Pr(T = 1 \mid \mathbf{X} = \mathbf{x})$ , states how the treatment status,  $t$ , depends on the covariates,  $\mathbf{x}$ , and is identifiable from observational data. The potential outcomes,  $Y_0$  and  $Y_1$ , conditioned on the covariates,  $\mathbf{x}$ , are distributed as  $P(Y_0 \mid \mathbf{X} = \mathbf{x})$  and  $P(Y_1 \mid \mathbf{X} = \mathbf{x})$ . Each of these conditional distributions can be written as mixtures with weights based on the propensity score:

$$\begin{aligned} P(Y_0 \mid \mathbf{X} = \mathbf{x}) &= (1 - e(\mathbf{x}))P(Y_0 \mid T = 0, \mathbf{X} = \mathbf{x}) + e(\mathbf{x})P(Y_0 \mid T = 1, \mathbf{X} = \mathbf{x}), \\ P(Y_1 \mid \mathbf{X} = \mathbf{x}) &= (1 - e(\mathbf{x}))P(Y_1 \mid T = 1, \mathbf{X} = \mathbf{x}) + e(\mathbf{x})P(Y_1 \mid T = 0, \mathbf{X} = \mathbf{x}). \end{aligned} \quad (16)$$

The conditional distributions of each potential outcome given the observed treatment,  $P(Y_0 \mid T = 0, \mathbf{X} = \mathbf{x})$  and  $P(Y_1 \mid T = 1, \mathbf{X} = \mathbf{x})$ , are identifiable from observational data, whereas the conditional distributions of each potential outcome given the counterfactual treatment,  $P(Y_0 \mid T = 1, \mathbf{X} = \mathbf{x})$  and  $P(Y_1 \mid T = 0, \mathbf{X} = \mathbf{x})$  are not. Under ignorability,  $\{Y_0, Y_1\} \perp\!\!\!\perp T \mid \mathbf{X} = \mathbf{x}$ ,  $P(Y_0 \mid T = 0, \mathbf{X} = \mathbf{x}) = P(Y_0 \mid T = 1, \mathbf{X} = \mathbf{x})$  and  $P(Y_1 \mid T = 1, \mathbf{X} = \mathbf{x}) = P(Y_1 \mid T = 0, \mathbf{X} = \mathbf{x})$ . Therefore, any deviation from these equalities will be indicative of hidden confounding. However, because the distributions  $P(Y_0 \mid T = 1, \mathbf{X} = \mathbf{x})$  and  $P(Y_1 \mid T = 0, \mathbf{X} = \mathbf{x})$  are unidentifiable, the MSM postulates a relationship between each pair of identifiable and unidentifiable components.

The MSM assumes that  $P(Y_t \mid T = 1 - t, \mathbf{X} = \mathbf{x})$  is absolutely continuous with respect to  $P(Y_t \mid T = t, \mathbf{X} = \mathbf{x})$  for all  $t \in \mathcal{T}_B$ . Therefore, given that  $P(Y_t \mid T = t, \mathbf{X} = \mathbf{x})$  and  $P(Y_t \mid T = 1 - t, \mathbf{X} = \mathbf{x})$  are  $\sigma$ -finite measures, by the Radon-Nikodym theorem, there exists a function  $\lambda_B(Y_t, \mathbf{x}; t) : \mathcal{Y} \rightarrow [0, \infty)$  such that,

$$P(Y_t \mid T = 1 - t, \mathbf{X} = \mathbf{x}) = \int_{\mathcal{Y}} \lambda_B(Y_t, \mathbf{x}; t) dP(Y_t \mid T = t, \mathbf{X} = \mathbf{x}). \quad (17)$$

Rearranging terms,  $\lambda_B(Y_t, \mathbf{x}; t)$  is expressed as the Radon-Nikodym derivative or ratio of densities,

$$\begin{aligned} \lambda_B(Y_t, \mathbf{x}; t) &= \frac{dP(Y_t \mid T = 1 - t, \mathbf{X} = \mathbf{x})}{dP(Y_t \mid T = t, \mathbf{X} = \mathbf{x})}, \\ &= \frac{p(y_t \mid T = 1 - t, \mathbf{X} = \mathbf{x})}{p(y_t \mid T = t, \mathbf{X} = \mathbf{x})}. \end{aligned} \quad (18)$$

By Bayes's rule,  $\lambda(Y_0, \mathbf{x}; 0)$  and  $\lambda(Y_1, \mathbf{x}; 1)$  are expressed as odds ratios,

$$\begin{aligned} \lambda_B(Y_0, \mathbf{x}; 0) &= \frac{1 - e(\mathbf{x})}{e(\mathbf{x})} \bigg/ \frac{1 - e(\mathbf{x}, y_0)}{e(\mathbf{x}, y_0)}, \\ \lambda_B(Y_1, \mathbf{x}; 1) &= \frac{e(\mathbf{x})}{1 - e(\mathbf{x})} \bigg/ \frac{e(\mathbf{x}, y_1)}{1 - e(\mathbf{x}, y_1)}, \end{aligned} \quad (19)$$

where  $e(\mathbf{x}, y_t) \equiv \Pr(T = 1 \mid \mathbf{X} = \mathbf{x}, Y_t = y_t)$  is the unidentifiable complete propensity for treatment.

Finally, the MSM further postulates that the odds of receiving the treatment  $T = 1$  for subjects with covariates  $\mathbf{X} = \mathbf{x}$  can only differ from  $e(\mathbf{x})/(1 - e(\mathbf{x}))$  by at most a factor of  $\Lambda$ ,

$$\Lambda^{-1} \leq \lambda_B(Y_t, \mathbf{x}; t) \leq \Lambda. \quad (20)$$

$$\alpha(e(\mathbf{x}, t), \Lambda) = \frac{1}{\Lambda e(\mathbf{x}, t)} + 1 - \frac{1}{\Lambda} \leq \frac{1}{e(\mathbf{x}, t, y_t)} \leq \frac{\Lambda}{e(\mathbf{x}, t)} + 1 - \Lambda = \beta(e(\mathbf{x}, t), \Lambda) \quad (21)$$

## A.2 Modifying the MSM for categorical treatment values

For categorical treatments,  $\mathcal{T}_C = \{t_i\}_{i=1}^{n_c}$ , the (nominal) generalized propensity score [HI04],  $r(\mathbf{x}, t) \equiv Pr(T = t \mid \mathbf{X} = \mathbf{x})$ , states how the treatment status,  $t$ , depends on the covariates,  $\mathbf{x}$ , and is identifiable from observational data. The potential outcomes,  $\{Y_t : t \in \mathcal{T}_C\}$ , conditioned on the covariates,  $\mathbf{x}$ , are distributed as  $\{P(Y_t \mid \mathbf{X} = \mathbf{x}) : t \in \mathcal{T}_C\}$ . Again, each of these conditional distributions can be written as mixtures with weights based on the propensity density, yielding the following set of mixture distributions:

$$\left\{ P(Y_t \mid \mathbf{X} = \mathbf{x}) = \sum_{t' \in \mathcal{T}_C} r(\mathbf{x}, t') P(Y_t \mid T = t', \mathbf{X} = \mathbf{x}) \right\}. \quad (22)$$

Each conditional distribution of the potential outcome given the observed treatment,  $P(Y_t \mid T = t, \mathbf{X} = \mathbf{x})$ , is identifiable from observational data, but each conditional distribution of the potential outcome given the counterfactual treatment,  $P(Y_t \mid T = t', \mathbf{X} = \mathbf{x})$ , and therefore each mixture  $P(Y_t \mid \mathbf{X} = \mathbf{x})$ , is not. Under the ignorability assumption,  $P(Y_t \mid T = t, \mathbf{X} = \mathbf{x}) = P(Y_t \mid T = t', \mathbf{X} = \mathbf{x})$  for all  $t' \in \mathcal{T}_C$ .

In order to recover the form of the binary treatment MSM, we can postulate a relationship between the unidentifiable  $P(Y_t \mid \mathbf{X} = \mathbf{x}) - r(\mathbf{x}, t)P(Y_t \mid T = t, \mathbf{X} = \mathbf{x})$  and the identifiable  $P(Y_t \mid T = t, \mathbf{X} = \mathbf{x}) - r(\mathbf{x}, t)P(Y_t \mid T = t, \mathbf{X} = \mathbf{x})$ . Under the assumption that  $P(Y_t \mid \mathbf{X} = \mathbf{x}) - r(\mathbf{x}, t)P(Y_t \mid T = t, \mathbf{X} = \mathbf{x})$  is absolutely continuous with respect to  $P(Y_t \mid T = t, \mathbf{X} = \mathbf{x}) - r(\mathbf{x}, t)P(Y_t \mid T = t, \mathbf{X} = \mathbf{x})$ , we define the Radon-Nikodym derivative

$$\begin{aligned} \lambda_C(Y_t, \mathbf{x}; t) &= \frac{d(P(Y_t \mid \mathbf{X} = \mathbf{x}) - r(\mathbf{x}, t)P(Y_t \mid T = t, \mathbf{X} = \mathbf{x}))}{d(1 - r(\mathbf{x}, t))P(Y_t \mid T = t, \mathbf{X} = \mathbf{x})}, \\ &= \frac{1}{1 - r(\mathbf{x}, t)} \left( \frac{dP(Y_t \mid \mathbf{X} = \mathbf{x})}{dP(Y_t \mid T = t, \mathbf{X} = \mathbf{x})} - \frac{r(\mathbf{x}, t)dP(Y_t \mid T = t, \mathbf{X} = \mathbf{x})}{dP(Y_t \mid T = t, \mathbf{X} = \mathbf{x})} \right), \\ &= \frac{1}{1 - r(\mathbf{x}, t)} \left( \frac{\sum_{t' \in \mathcal{T}_C} r(\mathbf{x}, t')dP(Y_t \mid T = t', \mathbf{X} = \mathbf{x})}{dP(Y_t \mid T = t, \mathbf{X} = \mathbf{x})} - \frac{r(\mathbf{x}, t)dP(Y_t \mid T = t, \mathbf{X} = \mathbf{x})}{dP(Y_t \mid T = t, \mathbf{X} = \mathbf{x})} \right), \\ &= \frac{1}{1 - r(\mathbf{x}, t)} \left( \frac{\sum_{t' \in \mathcal{T}_C} r(\mathbf{x}, t')p(y_t \mid T = t', \mathbf{X} = \mathbf{x})}{p(y_t \mid T = t, \mathbf{X} = \mathbf{x})} - \frac{r(\mathbf{x}, t)p(y_t \mid T = t, \mathbf{X} = \mathbf{x})}{p(y_t \mid T = t, \mathbf{X} = \mathbf{x})} \right), \\ &= \frac{1}{1 - r(\mathbf{x}, t)} \left( \frac{\sum_{t' \in \mathcal{T}_C} \frac{r(\mathbf{x}, t')p(T=t' \mid y_t, \mathbf{x})p(y_t)}{r(\mathbf{x}, t')} - \frac{r(\mathbf{x}, t)p(T=t \mid y_t, \mathbf{x})p(y_t)}{r(\mathbf{x}, t)} \right), \\ &= \frac{r(\mathbf{x}, t)}{1 - r(\mathbf{x}, t)} \frac{1 - p(T = t \mid y_t, \mathbf{x})}{p(T = t \mid y_t, \mathbf{x})}, \\ &= \frac{r(\mathbf{x}, t)}{1 - r(\mathbf{x}, t)} \bigg/ \frac{r(\mathbf{x}, t, y_t)}{1 - r(\mathbf{x}, t, y_t)}, \end{aligned} \quad (23)$$

where,  $r(\mathbf{x}, t, y_t) \equiv p(T = t \mid y_t, \mathbf{x})$  is the unidentifiable complete propensity density for treatment.

Finally, the categorical MSM further postulates that the odds of receiving the treatment  $T = t$  for subjects with covariates  $\mathbf{X} = \mathbf{x}$  can only differ from  $r(\mathbf{x}, t)/(1 - r(\mathbf{x}, t))$  by at most a factor of  $\Lambda$ ,

$$\Lambda^{-1} \leq \lambda_C(Y_t, \mathbf{x}; t) \leq \Lambda. \quad (24)$$

$$\alpha(r(\mathbf{x}, t), \Lambda) = \frac{1}{\Lambda r(\mathbf{x}, t)} + 1 - \frac{1}{\Lambda} \leq \frac{1}{r(\mathbf{x}, t, y_t)} \leq \frac{\Lambda}{r(\mathbf{x}, t)} + 1 - \Lambda = \beta(r(\mathbf{x}, t), \Lambda) \quad (25)$$

## A.3 Defining the Continuous MSM (CMSM) in terms of densities for continuous-valued interventions

The conditional distributions of the potential outcomes given the observed treatment assigned,

$$\{P(Y_t \mid T = t, \mathbf{X} = \mathbf{x}) : t \in \mathcal{T}\},$$

are identifiable from observational data. However, the marginal distributions of the potential outcomes over all possible treatments,

$$\left\{ \begin{aligned} &P(Y_t | \mathbf{X} = \mathbf{x}) = \\ &\int_{\mathcal{T}} p(t' | \mathbf{x}) P(Y_t | T = t', \mathbf{X} = \mathbf{x}) dt' \\ &: t \in \mathcal{T} \end{aligned} \right\} \quad (26)$$

are not. This is because the component distributions,  $P(Y_t | T = t', \mathbf{X} = \mathbf{x})$ , are not identifiable when  $t' \neq t$  as  $Y_t$  cannot be observed for units under treatment level  $T = t'$ . Under the ignorability assumption,  $P(Y_t | T = t, \mathbf{X} = \mathbf{x}) = P(Y_t | T = t', \mathbf{X} = \mathbf{x})$  for all  $t' \in \mathcal{T}$ , and so  $P(Y_t | \mathbf{X} = \mathbf{x})$  and  $P(Y_t | T = t, \mathbf{X} = \mathbf{x})$  are identical. Therefore, any divergence between  $P(Y_t | \mathbf{X} = \mathbf{x})$  and  $P(Y_t | T = t, \mathbf{X} = \mathbf{x})$  will be indicative of hidden confounding.

Where in the binary setting the MSM postulates a relationship between the unidentifiable  $P(Y_t | T = 1 - t, \mathbf{X} = \mathbf{x})$  and identifiable  $P(Y_t | T = t, \mathbf{X} = \mathbf{x})$ , our CMSM postulates a relationship between the unidentifiable  $P(Y_t | \mathbf{X} = \mathbf{x})$  and the identifiable  $P(Y_t | T = t, \mathbf{X} = \mathbf{x})$ .

The Radon-Nikodym theorem involves a measurable space  $(X, \Sigma)$  on which two  $\sigma$ -finite measures are defined,  $\mu$  and  $\nu$ ."

– Wikipedia

In our setting, the measurable space is  $(\mathbb{R}, \Sigma)$ , and our  $\sigma$ -finite measures are,  $\mu = P(Y_t | T = t, \mathbf{X} = \mathbf{x})$  and  $\nu = P(Y_t | \mathbf{X} = \mathbf{x})$ :  $Y_t \in \mathcal{Y} \subseteq \mathbb{R}$ .

If  $\nu$  is absolutely continuous with respect to  $\mu$  (written  $\nu \ll \mu$ ), then there exists a  $\Sigma$ -measurable function  $f : X \rightarrow [0, \infty)$ , such that  $\nu(A) = \int_A f d\mu$  for any measurable set  $A \subseteq X$ .

– Wikipedia

We then need to assume that  $P(Y_t | \mathbf{X} = \mathbf{x}) \ll P(Y_t | T = t, \mathbf{X} = \mathbf{x})$ , that is  $P(A | T = t, \mathbf{X} = \mathbf{x}) = 0$  implies  $P(A | \mathbf{X} = \mathbf{x}) = 0$  for any measurable set  $A$ .

This leads us to a proof for Proposition 1

*Proof.* Further, in our setting we have  $f = \lambda(y_t; \mathbf{x}, t)$ , therefore

$$P(Y_t | \mathbf{X} = \mathbf{x}) = \int_{\mathcal{Y}} \lambda(y_t; \mathbf{x}, t) dP(Y_t | T = t, \mathbf{X} = \mathbf{x}). \quad (27)$$

Let the range of  $Y_t$  be the measurable space  $(\mathcal{Y}, \mathcal{A})$ , and  $\nu(A)$  denote the Lebesgue measure for any measurable  $A \in \mathcal{A}$ . Then,

$$\lambda(y_t; \mathbf{x}, t) = \frac{dP(Y_t | \mathbf{X} = \mathbf{x})}{dP(Y_t | T = t, \mathbf{X} = \mathbf{x})} \quad (28a)$$

$$= \frac{dP(Y_t | \mathbf{X} = \mathbf{x})}{d\nu} \frac{d\nu}{dP(Y_t | T = t, \mathbf{X} = \mathbf{x})} \quad (28b)$$

$$= \frac{dP(Y_t | \mathbf{X} = \mathbf{x})}{d\nu} \left( \frac{dP(Y_t | T = t, \mathbf{X} = \mathbf{x})}{d\nu} \right)^{-1} \quad (28c)$$

$$= \frac{d}{d\nu} \int_A p(y_t | \mathbf{X} = \mathbf{x}) d\nu \left( \frac{d}{d\nu} \int_A p(y_t | T = t, \mathbf{X} = \mathbf{x}) d\nu \right)^{-1} \quad (28d)$$

$$= \frac{p(y_t | \mathbf{X} = \mathbf{x})}{p(y_t | T = t, \mathbf{X} = \mathbf{x})} \quad (28e)$$

$$= \frac{p(t | \mathbf{X} = \mathbf{x})}{p(t | Y_t = y_t, \mathbf{X} = \mathbf{x})} \quad (28f)$$

Equation (28a) by the Radon-Nikodym derivative. Equation (28a)-Equation (28c) hold  $\nu$ -almost everywhere under the assumption  $P(Y_t \in A \mid \mathbf{x}) \ll \nu(A) \sim P(Y_t \in A \mid T = t, \mathbf{X} = \mathbf{x})$ . Equation (28c)-Equation (28d) by the Radon-Nikodym theorem. Equation (28d)-Equation (28e) by the fundamental theorem of calculus under the assumption that  $p(y_t \mid \mathbf{x})$  and  $p(y_t \mid T = t, \mathbf{X} = \mathbf{x})$  be continuous for  $y_t \in \mathcal{Y}$ . Equation (28e)-Equation (28f) by Bayes's Rule.  $\square$

The sensitivity analysis parameter  $\Lambda$  then bounds the ratio, which leads to our bounds for the inverse complete propensity density:

$$\begin{aligned} \frac{1}{\Lambda} &\leq \frac{p(t \mid \mathbf{x})}{p(t \mid y_t, \mathbf{x})} \leq \Lambda, \\ \frac{1}{\Lambda p(t \mid \mathbf{x})} &\leq \frac{1}{p(t \mid y_t, \mathbf{x})} \leq \frac{\Lambda}{p(t \mid \mathbf{x})} \\ \alpha(p(t \mid \mathbf{x}), \Lambda) &\leq \frac{1}{p(t \mid y_t, \mathbf{x})} \leq \beta(p(t \mid \mathbf{x}), \Lambda) \end{aligned} \tag{29}$$

### A.3.1 KL Divergence

The bounds on the density ratio can also be expressed as bounds on the Kullback-Leibler divergence between  $P(Y_t \mid T = t, \mathbf{X} = \mathbf{x})$  and  $P(Y_t \mid \mathbf{X} = \mathbf{x})$ .

$$\Lambda^{-1} \leq \frac{p(t \mid \mathbf{x})}{p(t \mid y_t, \mathbf{x})} \leq \Lambda, \tag{30}$$

$$\log(\Lambda^{-1}) \leq \log\left(\frac{p(t \mid \mathbf{x})}{p(t \mid y_t, \mathbf{x})}\right) \leq \log(\Lambda) \tag{31}$$

$$\mathbb{E}_{p(y|t, \mathbf{x})} \log(\Lambda^{-1}) \leq \mathbb{E}_{p(y|t, \mathbf{x})} \log\left(\frac{p(t \mid \mathbf{x})}{p(t \mid y_t, \mathbf{x})}\right) \leq \mathbb{E}_{p(y|t, \mathbf{x})} \log(\Lambda) \tag{32}$$

$$\log(\Lambda^{-1}) \leq \mathbb{E}_{p(y|t, \mathbf{x})} \log\left(\frac{p(t \mid \mathbf{x})}{p(t \mid y_t, \mathbf{x})}\right) \leq \log(\Lambda) \tag{33}$$

$$\log(\Lambda^{-1}) \leq \int_{\mathcal{Y}} \log\left(\frac{dP(Y_t \mid \mathbf{X} = \mathbf{x})}{dP(Y_t \mid T = t, \mathbf{X} = \mathbf{x})}\right) dP(Y_t \mid T = t, \mathbf{X} = \mathbf{x}) \leq \log(\Lambda) \tag{34}$$

$$\log(\Lambda^{-1}) \leq -D_{\text{KL}}(P(Y_t \mid T = t, \mathbf{X} = \mathbf{x}) \parallel P(Y_t \mid \mathbf{X} = \mathbf{x})) \leq \log(\Lambda) \tag{35}$$

$$|\log(\Lambda)| \geq D_{\text{KL}}(P(Y_t \mid T = t, \mathbf{X} = \mathbf{x}) \parallel P(Y_t \mid \mathbf{X} = \mathbf{x})) \tag{36}$$

## B Derivation of Equation (7)

**Lemma 1.**

$$\mu(\mathbf{x}, t) = \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} w(y, \mathbf{x})(y - \tilde{\mu}(\mathbf{x}, t))p(y \mid t, \mathbf{x})dy}{(\Lambda^2 - 1)^{-1} + \int_{\mathcal{Y}} w(y, \mathbf{x})p(y \mid t, \mathbf{x})dy} \tag{37}$$

*Proof.* Recall that the conditional average potential outcome,  $\mu(\mathbf{x}, t) = \mathbb{E}[Y_t \mid \mathbf{X} = \mathbf{x}]$ , is unidentifiable without further assumptions. Following [KMZ19], we start from,

$$\begin{aligned}
\mu(\mathbf{x}, t) &= \mathbb{E}[Y_t \mid \mathbf{X} = \mathbf{x}], \\
&= \frac{\int_{\mathcal{Y}} y_t p(y_t \mid \mathbf{x}) dy_t}{\int_{\mathcal{Y}} p(y_t \mid \mathbf{x}) dy_t}, \\
&= \frac{\int_{\mathcal{Y}} y_t \frac{p(t, y_t \mid \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(t, y_t \mid \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}, \\
&= \frac{\int_{\mathcal{Y}} y_t \frac{p(y_t \mid t, \mathbf{x}) p(t \mid \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(y_t \mid t, \mathbf{x}) p(t \mid \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}, \\
&= \frac{\int_{\mathcal{Y}} y_t \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t},
\end{aligned}$$

which is convenient as it decomposes  $\mu(\mathbf{x}, t)$  into its identifiable,  $p(y_t \mid t, \mathbf{x})$ , and unidentifiable,  $p(t \mid y_t, \mathbf{x})$ , parts.

Now, following [JMGS21], we add and subtract the empirical conditional outcome  $\tilde{\mu}(\mathbf{x}, t) = \mathbb{E}[Y \mid T = t, \mathbf{X} = \mathbf{x}]$  from the right-hand-side above:

$$\mu(\mathbf{x}, t) = \frac{\int_{\mathcal{Y}} y_t \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}, \quad (39a)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} y_t \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t} - \tilde{\mu}(\mathbf{x}, t), \quad (39b)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} y_t \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t} - \tilde{\mu}(\mathbf{x}, t) \frac{\int_{\mathcal{Y}} \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}, \quad (39c)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} y_t \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t} - \frac{\int_{\mathcal{Y}} \tilde{\mu}(\mathbf{x}, t) \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}, \quad (39d)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} (y - \tilde{\mu}(\mathbf{x}, t)) \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(y_t \mid t, \mathbf{x})}{p(t \mid y_t, \mathbf{x})} dy_t}. \quad (39e)$$

Following [KMZ19] again, we reparameterize the inverse complete propensity density as,  $\frac{1}{p(t \mid y_t, \mathbf{x})} = \alpha(\mathbf{x}; t, \Lambda) + w(y, \mathbf{x})(\beta(\mathbf{x}; t, \Lambda) - \alpha(\mathbf{x}; t, \Lambda))$  with  $w : \mathcal{Y} \times \mathcal{X} \rightarrow [0, 1]$ . We will shorten this

expression to  $\frac{1}{p(t|y_t, \mathbf{x})} = \alpha + w(y, \mathbf{x})(\beta - \alpha)$  below. This gives,

$$\mu(\mathbf{x}, t) = \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} (y - \tilde{\mu}(\mathbf{x}, t)) \frac{p(y_t | t, \mathbf{x})}{p(t | y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(y_t | t, \mathbf{x})}{p(t | y_t, \mathbf{x})} dy_t}, \quad (40a)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} (\alpha + w(y, \mathbf{x})(\beta - \alpha))(y - \tilde{\mu}(\mathbf{x}, t)) p(y_t | t, \mathbf{x}) dy_t}{\int_{\mathcal{Y}} (\alpha + w(y, \mathbf{x})(\beta - \alpha)) p(y_t | t, \mathbf{x}) dy_t}, \quad (40b)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\alpha \int_{\mathcal{Y}} (y - \tilde{\mu}(\mathbf{x}, t)) p(y_t | t, \mathbf{x}) dy_t + (\beta - \alpha) \int_{\mathcal{Y}} (y - \tilde{\mu}(\mathbf{x}, t)) w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}{\alpha \int_{\mathcal{Y}} p(y_t | t, \mathbf{x}) dy_t + (\beta - \alpha) \int_{\mathcal{Y}} w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}, \quad (40c)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\alpha \int_{\mathcal{Y}} (y - \tilde{\mu}(\mathbf{x}, t)) p(y_t | t, \mathbf{x}) dy_t + (\beta - \alpha) \int_{\mathcal{Y}} (y - \tilde{\mu}(\mathbf{x}, t)) w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}{\alpha + (\beta - \alpha) \int_{\mathcal{Y}} w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}, \quad (40d)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{(\beta - \alpha) \int_{\mathcal{Y}} (y - \tilde{\mu}(\mathbf{x}, t)) w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}{\alpha + (\beta - \alpha) \int_{\mathcal{Y}} w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}, \quad (40e)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} (y - \tilde{\mu}(\mathbf{x}, t)) w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}{\frac{\alpha}{\beta - \alpha} + \int_{\mathcal{Y}} w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}, \quad (40f)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} (y - \tilde{\mu}(\mathbf{x}, t)) w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}{\frac{1/(\Lambda p(t | \mathbf{x}))}{\Lambda/p(t | \mathbf{x}) - 1/(\Lambda p(t | \mathbf{x}))} + \int_{\mathcal{Y}} w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}, \quad (40g)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} (y - \tilde{\mu}(\mathbf{x}, t)) w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}{\frac{1}{\Lambda^2 - 1} + \int_{\mathcal{Y}} w(y, \mathbf{x}) p(y_t | t, \mathbf{x}) dy_t}, \quad (40h)$$

which concludes the proof.  $\square$

## C Approximating integrals using Gauss-Hermite quadrature

Gauss-Hermite quadrature is a numerical method to approximate indefinite integrals of the following form:  $\int_{-\infty}^{\infty} \exp(-y^2) f(y) dy$ . In this case,

$$\int_{-\infty}^{\infty} \exp(-y^2) f(y) dy \approx \sum_{i=1}^m g_i f(y_i),$$

where  $m$  is the number of samples drawn. The  $y_i$  are the roots of the physicists Hermite polynomial  $H_m^*(y)$  ( $i = 1, 2, \dots, m$ ) and the weights are given by

$$g_i = \frac{2^{m-1} m! \sqrt{\pi}}{m^2 [H_{m-1}^*(y_k)]^2}$$

This method can be used to calculate the expectation of a function,  $h(y)$ , with respect to a Gaussian distributed outcome  $p(y) = \mathcal{N}(y | \mu, \sigma^2)$  through a change of variables, such that,

$$\begin{aligned} \mathbb{E}_{p(y)}[h(y)] &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{\pi}} \exp(-y^2) h(\sqrt{2}\sigma y + \mu) dy \\ &\approx \frac{1}{\sqrt{\pi}} \sum_{i=1}^m g_i h(\sqrt{2}\sigma y_i + \mu). \end{aligned} \quad (41)$$

**Definition 2.** Gauss-Hermite quadrature integral estimator when  $p(y | t, \mathbf{x}, \boldsymbol{\theta})$  is a parametric Gaussian density estimator,  $\mathcal{N}(y | \tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta}), \tilde{\sigma}^2(\mathbf{x}, t; \boldsymbol{\theta}))$ :

$$I_G(h(y)) := \frac{1}{\sqrt{\pi}} \sum_{i=1}^m g_i h(\sqrt{2}\tilde{\sigma}^2(\mathbf{x}, t; \boldsymbol{\theta}) y + \tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta}))$$

Alternatively, when the density of the outcome is modelled using a  $n_y$  component Gaussian mixture,  $p(y) = \sum_{j=1}^{n_y} \pi_j \mathcal{N}(y \mid \mu_j, \sigma_j^2)$

$$\begin{aligned} \mathbb{E}_{p(y)}[h(y)] &= \frac{1}{\sqrt{\pi}} \sum_{j=1}^{n_y} \pi_j \int_{-\infty}^{\infty} \exp(-y^2) h(\sqrt{2}\sigma_j y + \mu_j) dy, \\ &\approx \frac{1}{\sqrt{\pi}} \sum_{j=1}^{n_y} \pi_j \sum_{i=1}^m g_i h(\sqrt{2}\sigma_j y + \mu_j). \end{aligned}$$

**Definition 3.** Gauss-Hermite quadrature integral estimator for expectations when  $p(y|\mathbf{t}, \mathbf{x}, \boldsymbol{\theta})$  is a parametric Gaussian Mixture Density,  $\sum_{j=1}^{n_y} \tilde{\pi}_j(\mathbf{x}, \mathbf{t}; \boldsymbol{\theta}) \mathcal{N}(y \mid \tilde{\mu}_j(\mathbf{x}, \mathbf{t}; \boldsymbol{\theta}), \tilde{\sigma}_j^2(\mathbf{x}, \mathbf{t}; \boldsymbol{\theta}))$ :

$$I_{GM}(h(y)) := \frac{1}{\sqrt{\pi}} \sum_{j=1}^{n_y} \tilde{\pi}_j(\mathbf{x}, \mathbf{t}; \boldsymbol{\theta}) \sum_{i=1}^m g_i h(\sqrt{2}\tilde{\sigma}_j(\mathbf{x}, \mathbf{t}; \boldsymbol{\theta})y + \tilde{\mu}_j(\mathbf{x}, \mathbf{t}; \boldsymbol{\theta}))$$

## D Optimization over step functions

**Lemma 2.** The sensitivity bounds given in Equations (8) and (9) have the following equivalent expressions:

$$\begin{aligned} \bar{\mu}(\mathbf{x}, \mathbf{t}; \Lambda) &= \sup_{w(y) \in \mathcal{W}_{nd}^H} \tilde{\mu}(\mathbf{x}, \mathbf{t}) + \frac{\int_{\mathcal{Y}} w(y)(y - \tilde{\mu}(\mathbf{x}, \mathbf{t}))p(y \mid \mathbf{t}, \mathbf{x})dy}{(\Lambda^2 - 1)^{-1} + \int_{\mathcal{Y}} w(y)p(y \mid \mathbf{t}, \mathbf{x})dy}, \\ \underline{\mu}(\mathbf{x}, \mathbf{t}; \Lambda) &= \inf_{w(y) \in \mathcal{W}_{ni}^H} \tilde{\mu}(\mathbf{x}, \mathbf{t}) + \frac{\int_{\mathcal{Y}} w(y)(y - \tilde{\mu}(\mathbf{x}, \mathbf{t}))p(y \mid \mathbf{t}, \mathbf{x})dy}{(\Lambda^2 - 1)^{-1} + \int_{\mathcal{Y}} w(y)p(y \mid \mathbf{t}, \mathbf{x})dy}, \end{aligned}$$

where  $\tilde{\mu}(\mathbf{x}, \mathbf{t}) = \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}, \mathbf{T} = \mathbf{t}]$ ,  $\mathcal{W}_{nd}^H = \{w : H(y - y_H)\}_{y_H \in \mathcal{Y}}$ ,  $\mathcal{W}_{ni}^H = \{w : H(y_H - y)\}_{y_H \in \mathcal{Y}}$ , and

$$H(z) := \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases},$$

*Proof.* We follow the argument of [KMZ19] and show that our alternative formulations of  $\alpha(\cdot, \Lambda)$  and  $\beta(\cdot, \Lambda)$  do not change the conclusions of their linear program solution. Starting from  $\mu(\mathbf{x}, \mathbf{t}) = \frac{\int_{\mathcal{Y}} y_t \frac{p(\mathbf{t}, y_t | \mathbf{x})}{p(\mathbf{t} | \mathbf{y}_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(\mathbf{t}, y_t | \mathbf{x})}{p(\mathbf{t} | \mathbf{y}_t, \mathbf{x})} dy_t}$ , and applying a one-to-one change of variables,  $\frac{1}{p(\mathbf{t} | \mathbf{y}_t, \mathbf{x})} = \alpha(\mathbf{x}; \mathbf{t}, \Lambda) + w(y)(\beta(\mathbf{x}; \mathbf{t}, \Lambda) - \alpha(\mathbf{x}; \mathbf{t}, \Lambda))$  with  $w : \mathcal{Y} \rightarrow [0, 1]$ ,  $\alpha(\mathbf{x}; \mathbf{t}, \Lambda) = 1/\Lambda p(\mathbf{t} \mid \mathbf{x})$ ,  $\beta(\mathbf{x}; \mathbf{t}, \Lambda) = \Lambda/p(\mathbf{t} \mid \mathbf{x})$ , we arrive at:

$$\bar{\mu}(\mathbf{x}, \mathbf{t}; \Lambda) = \sup_{w: \mathcal{Y} \rightarrow [0, 1]} \frac{\int_{\mathcal{Y}} y p(y \mid \mathbf{t}, \mathbf{x}) dy + (\lambda^2 - 1) \int_{\mathcal{Y}} y w(y) p(y \mid \mathbf{t}, \mathbf{x}) dy}{1 + (\lambda^2 - 1) \int_{\mathcal{Y}} w(y) p(y \mid \mathbf{t}, \mathbf{x}) dy}, \quad (42)$$

and

$$\underline{\mu}(\mathbf{x}, \mathbf{t}; \Lambda) = \inf_{w: \mathcal{Y} \rightarrow [0, 1]} \frac{\int_{\mathcal{Y}} y p(y \mid \mathbf{t}, \mathbf{x}) dy + (\lambda^2 - 1) \int_{\mathcal{Y}} y w(y) p(y \mid \mathbf{t}, \mathbf{x}) dy}{1 + (\lambda^2 - 1) \int_{\mathcal{Y}} w(y) p(y \mid \mathbf{t}, \mathbf{x}) dy}, \quad (43)$$

after some cancellations. Duality can be used to prove that the  $w^*(y)$  which achieves the supremum in Equation (42) belongs to the set of step functions  $\mathcal{W}_{nd}^H$ . An analogous proof for Equation (43) would show that the  $w^*(y)$  which achieves the infimum in Equation (43) belongs to the set of step functions  $\mathcal{W}_{ni}^H$ .

The optimization problem in Equation (42) can be rewritten as a linear-fractional program:

$$\text{maximize} \quad \frac{a \langle y, w(y) \rangle_{p(y|\mathbf{t}, \mathbf{x})} + c}{b \langle 1, w(y) \rangle_{p(y|\mathbf{t}, \mathbf{x})} + d} \quad (44a)$$

$$\text{subject to} \quad 0 \leq w(y) \leq 1 : \forall y \in \mathcal{Y}, \quad (44b)$$

where  $\langle \cdot, \cdot \rangle_{p(y|\mathbf{t}, \mathbf{x})}$  is the inner product with respect to  $p(y \mid \mathbf{t}, \mathbf{x})$ ,  $a = b = \lambda^2 - 1$ ,  $c = \int_{\mathcal{Y}} y p(y \mid \mathbf{t}, \mathbf{x}) dy$ , and  $d = \int_{\mathcal{Y}} p(y \mid \mathbf{t}, \mathbf{x}) dy$ .

The linear-fractional program of Equation (44) is equivalent to the following linear program:

$$\text{maximize } a\langle y, \tilde{w}(y) \rangle_{p(y|t, \mathbf{x})} + c\tilde{v}(\mathbf{x}) \quad (45a)$$

$$\text{subject to } \tilde{w}(y) \leq \tilde{v}(\mathbf{x}) : \forall y \in \mathcal{Y} \quad (45b)$$

$$-\tilde{w}(y) \leq 0 : \forall y \in \mathcal{Y} \quad (45c)$$

$$b\langle 1, \tilde{w}(y) \rangle_{p(y|t, \mathbf{x})} + d\tilde{v}(\mathbf{x}) = 1 \quad (45d)$$

$$\tilde{v}(\mathbf{x}) \geq 0, \quad (45e)$$

where

$$\tilde{w}(y) = \frac{w(y)}{b\langle 1, w(y) \rangle_{p(y|t, \mathbf{x})} + d} \quad \text{and} \quad \tilde{v}(\mathbf{x}) = \frac{1}{b\langle 1, w(y) \rangle_{p(y|t, \mathbf{x})} + d}$$

by the Charnes-Cooper transformation.

Let the dual function  $\rho(y)$  be associated with the primal constraint eq. (45b), the dual function  $\eta(y)$  be associated with the primal constraint eq. (45c), and  $\gamma$  be the dual variable associated with the primal constraint eq. (45d). The dual program is then:

$$\text{minimize } \gamma \quad (46a)$$

$$\text{subject to } \rho(y) - \eta(y) + \gamma b p(y | t, \mathbf{x}) = a y p(y | t, \mathbf{x}) : \forall y \in \mathcal{Y} \quad (46b)$$

$$-\langle 1, \rho(y) \rangle + \gamma d \geq c \quad (46c)$$

$$\rho(y) \in \mathbb{R}_+, \eta(y) \in \mathbb{R}_+, \gamma \in \mathbb{R} \quad (46d)$$

At most one of  $\rho(y)$  or  $\eta(y)$  is non-zero by complementary slackness; therefore, condition eq. (46b) implies that

$$\rho(y) = (\lambda^2 - 1)p(y | t, \mathbf{x}) \max\{y - \gamma, 0\} \text{ when } \eta = 0,$$

$$\eta(y) = (\lambda^2 - 1)p(y | t, \mathbf{x}) \max\{\gamma - y, 0\} \text{ when } \rho = 0.$$

[KMZ19] argue that constraint eq. (46c) ought to be tight (an equivalence) at optimality, otherwise there would exist a smaller, feasible  $\gamma$  that satisfies the linear program. Therefore,

$$\begin{aligned} -\langle 1, \rho(y) \rangle + \gamma d &= c, \\ -\int_{\mathcal{Y}} (\lambda^2 - 1)p(y | t, \mathbf{x}) \max\{y - \gamma, 0\} dy + \gamma \int_{\mathcal{Y}} p(y | t, \mathbf{x}) dy &= \int_{\mathcal{Y}} y p(y | t, \mathbf{x}) dy, \\ (\lambda^2 - 1) \int_{\mathcal{Y}} \max\{y - \gamma, 0\} p(y | t, \mathbf{x}) dy &= \int_{\mathcal{Y}} (\gamma - y) p(y | t, \mathbf{x}) dy. \end{aligned} \quad (47)$$

Letting  $C_Y > 0$  such that  $|\mathcal{Y}| \leq C_Y$ , it is impossible that either  $\gamma > C_Y$  (the r.h.s. would be 0 and the l.h.s. would be  $> 0$ ) or  $\gamma < -C_Y$  (the r.h.s. would be  $> 0$  and the l.h.s. would be  $< 0$ ). Thus,  $\exists y^* \in [-C_Y, C_Y]$  such that when  $y < y^*$ ,  $\eta > 0$  so  $w = 0$  and when  $y \geq y^*$ ,  $\rho > 0$  so  $w = 1$ . Therefore, the optimal  $w^*(y)$  that achieves the supremum in Equation (42) is in  $\mathcal{W}_{\text{nd}}^H$ .

This result holds under

$$\mu(\mathbf{x}, t) = \frac{\int_{\mathcal{Y}} y p(y | t, \mathbf{x}) dy + (\lambda^2 - 1) \int_{\mathcal{Y}} y w(y) p(y | t, \mathbf{x}) dy}{1 + (\lambda^2 - 1) \int_{\mathcal{Y}} w(y) p(y | t, \mathbf{x}) dy}, \quad (48a)$$

$$= \frac{\int_{\mathcal{Y}} y_t \frac{p(t, y_t | \mathbf{x})}{p(t | y_t, \mathbf{x})} dy_t}{\int_{\mathcal{Y}} \frac{p(t, y_t | \mathbf{x})}{p(t | y_t, \mathbf{x})} dy_t}, \quad (48b)$$

$$= \tilde{\mu}(\mathbf{x}, t) + \frac{\int_{\mathcal{Y}} w(y) (y - \tilde{\mu}(\mathbf{x}, t)) p(y | t, \mathbf{x}) dy}{(\lambda^2 - 1)^{-1} + \int_{\mathcal{Y}} w(y) p(y | t, \mathbf{x}) dy}, \quad (48c)$$

thus concluding the proof (eq. (48b)-eq. (48c) by Lemma 1).  $\square$

### D.0.1 Discrete search approaches

Let  $\hat{\mathcal{Y}} = \{y_i \in \mathcal{Y}\}_{i=1}^k$  be a set of  $k$  values of  $y$ , then

$$\underline{\mu}_{\theta}^H(\mathbf{x}, t) = \min_{y^*} \left\{ \hat{\kappa}_{\theta}(\mathbf{x}, t; \Lambda, H(y^* - y)) : y^* \in \hat{\mathcal{Y}} \right\},$$

$$\bar{\mu}_{\theta}^H(\mathbf{x}, t) = \max_{y^*} \left\{ \hat{\kappa}_{\theta}(\mathbf{x}, t; \Lambda, H(y - y^*)) : y^* \in \hat{\mathcal{Y}} \right\}.$$



$$H(y) := \begin{cases} 1, & y > 0 \\ 0, & y \leq 0 \end{cases}$$

---

**Algorithm 2** Line Search Interval Optimizer

---

**Require:**  $\mathbf{x}^*$  is an instance of  $\mathbf{X}$ ,  $t^*$  is a treatment level to evaluate,  $\Lambda$  is a belief in the amount of hidden confounding,  $\theta$  are optimized model parameters,  $\hat{\mathcal{Y}}$  is a set of unique values  $y \in \mathcal{Y}$  sorted in ascending order.

```

1: function LINESEARCH( $\mathbf{x}^*, t^*, \Lambda, \theta, \hat{\mathcal{Y}}$ )
2:    $\bar{\mu} \leftarrow -\infty, \bar{\kappa} \leftarrow \infty$ 
3:    $\underline{\mu} \leftarrow \infty, \underline{\kappa} \leftarrow -\infty$ 
4:    $\underline{\delta} \leftarrow \text{True}, \bar{\delta} \leftarrow \text{True}$ 
5:   while  $\underline{\delta}$  do
6:      $y^* \leftarrow \text{POP}(\hat{\mathcal{Y}}_c)$   $\triangleright \hat{\mathcal{Y}}_c$  a copy of  $\hat{\mathcal{Y}}$ 
7:      $\underline{\kappa} \leftarrow \hat{\kappa}_{\theta}(\mathbf{x}, t; \Lambda, H(y^* - y))$ 
8:     if  $\underline{\kappa} < \underline{\mu}$  then
9:        $\underline{\mu} \leftarrow \underline{\kappa}$ 
10:    else
11:       $\underline{\delta} \leftarrow \text{False}$ 
12:    while  $\bar{\delta}$  do
13:       $y^* \leftarrow \text{POP}(\hat{\mathcal{Y}}_c)$   $\triangleright \hat{\mathcal{Y}}_c$  a copy of  $\hat{\mathcal{Y}}$ 
14:       $\bar{\kappa} \leftarrow \hat{\kappa}_{\theta}(\mathbf{x}, t; \Lambda, H(y - y^*))$ 
15:      if  $\bar{\kappa} > \bar{\mu}$  then
16:         $\bar{\mu} \leftarrow \bar{\kappa}$ 
17:      else
18:         $\bar{\delta} \leftarrow \text{False}$ 
19:  return  $\underline{\mu}, \bar{\mu}$ 

```

---

## E Theorem 1

Assume that

1.  $m \rightarrow \infty$ ,
2.  $n \rightarrow \infty$ ,
3.  $(\mathbf{X} = \mathbf{x}, T = t) \in \mathcal{D}_n$ ,
4.  $p(y \mid t, \mathbf{x}, \theta)$  converges in measure to  $p(y \mid t, \mathbf{x})$ ,
5.  $\tilde{\mu}(\mathbf{x}, t; \theta)$  is a consistent estimator of  $\tilde{\mu}(\mathbf{x}, t)$ ,
6.  $p(t \mid y_t, \mathbf{x}) > 0, \forall y_t \in \mathcal{Y}$ .

Then,  $\underline{\mu}(\mathbf{x}, t; \Lambda, \theta) \xrightarrow{P} \underline{\mu}(\mathbf{x}, t; \Lambda)$  and  $\bar{\mu}(\mathbf{x}, t; \Lambda, \theta) \xrightarrow{P} \bar{\mu}(\mathbf{x}, t; \Lambda)$ .

*Proof.* We prove that  $\underline{\mu}(\mathbf{x}, t; \Lambda, \theta) \xrightarrow{P} \underline{\mu}(\mathbf{x}, t; \Lambda)$ , from which  $\bar{\mu}(\mathbf{x}, t; \Lambda, \theta) \xrightarrow{P} \bar{\mu}(\mathbf{x}, t; \Lambda)$  can be proved analogously. Note that  $\xrightarrow{P}$  denotes “convergence in probability”. We need to show that  $\lim_n P(|\underline{\mu}(\mathbf{x}, t; \Lambda, \theta_n) - \underline{\mu}(\mathbf{x}, t; \Lambda)| \geq \epsilon) = 0$ , for all  $\epsilon > 0$ . Where  $\theta_n$  are the model parameters corresponding to a dataset  $\mathcal{D}_n$  of  $n$  observations. Recall that,

$$\underline{\mu}(\mathbf{x}, t; \Lambda) := \tilde{\mu}(\mathbf{x}, t) + \inf_{w \in \mathcal{W}_{\text{in}}^H} \frac{\int_{\mathcal{Y}} w(y)(y - \tilde{\mu}(\mathbf{x}, t))p(y \mid t, \mathbf{x})dy}{(\Lambda^2 - 1)^{-1} + \int_{\mathcal{Y}} w(y)p(y \mid t, \mathbf{x})dy},$$

and

$$\underline{\mu}(\mathbf{x}, t; \Lambda, \theta_n) := \tilde{\mu}(\mathbf{x}, t; \theta_n) + \inf_{w \in \mathcal{W}_{\text{in}}^H} \frac{I_m(w(y)(y - \tilde{\mu}(\mathbf{x}, t; \theta_n)))}{(\Lambda^2 - 1)^{-1} + I_m(w(y))},$$

where

$$I_m(w(y)(y - \tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta}_n))) = \frac{1}{m} \sum_{i=1}^m w(y_i)(y_i - \tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta}_n)),$$

and

$$I_m(w(y)) = \frac{1}{m} \sum_{i=1}^m w(y_i),$$

with  $y_i \sim p(y \mid t, \mathbf{x}, \boldsymbol{\theta}_n)$ .

First, by Item 1 and the law of large numbers, both

$$\lim_{m \rightarrow \infty} I_m(w(y)(y - \tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta}_n))) = \int_{\mathcal{Y}} w(y)(y - \tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta}_n))p(y \mid t, \mathbf{x}; \boldsymbol{\theta}_n)dy,$$

and

$$\lim_{m \rightarrow \infty} I_m(w(y)) = \int_{\mathcal{Y}} w(y)p(y \mid t, \mathbf{x}; \boldsymbol{\theta}_n)dy.$$

Therefore,

$$\lim_{m \rightarrow \infty} \underline{\mu}(\mathbf{x}, t; \Lambda, \boldsymbol{\theta}_n) = \tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta}_n) + \inf_{w \in \mathcal{W}_m^H} \frac{\int_{\mathcal{Y}} w(y)(y - \tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta}_n))p(y \mid t, \mathbf{x}; \boldsymbol{\theta}_n)dy}{(\Lambda^2 - 1)^{-1} + \int_{\mathcal{Y}} w(y)p(y \mid t, \mathbf{x}; \boldsymbol{\theta}_n)dy}.$$

Note that this step was missed by [JMGS21].

From here, the proof for Theorem 1 from [JMGS21] can be followed, substituting in  $(\Lambda^2 - 1)^{-1}$  where they write  $\alpha'_{\omega}$  and  $\alpha'$ .

□

## F Optimization over continuous functions

Second, we need a functional estimator for  $w(y, \mathbf{x})$ . We use a neural network,  $w(y, \mathbf{x}; \boldsymbol{\omega})$ , parameterized by  $\boldsymbol{\omega}$  with sigmoid non-linearity on the output layer to satisfy the  $w : \mathcal{Y} \times \mathcal{X} \rightarrow [0, 1]$  constraint.

For each  $(\Lambda, t)$  pair, we then need to solve the following optimization problems:

$$\underline{\omega} = \arg \min_{\omega} \frac{1}{n} \sum_{i=1}^n \mu(w(y, \cdot; \omega); \mathbf{x}_i, t, \Lambda, \boldsymbol{\theta}), \quad \mathbf{x}_i \in \mathcal{D},$$

and

$$\overline{\omega} = \arg \min_{\omega} \frac{1}{n} \sum_{i=1}^n -\mu(w(y, \cdot; \omega); \mathbf{x}_i, t, \Lambda, \boldsymbol{\theta}), \quad \mathbf{x}_i \in \mathcal{D},$$

where

$$\begin{aligned} & \mu(w(y, \cdot; \omega); \mathbf{x}, t, \Lambda, \boldsymbol{\theta}) \\ &:= \tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta}) + \frac{I(w(y, \mathbf{x}; \omega)(y - \tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta})))}{(\Lambda^2 - 1)^{-1} + I(w(y, \mathbf{x}; \omega))}. \end{aligned}$$

Each of these problems can then be optimized using stochastic gradient descent [Rud16] and error back-propagation [RHW86]. Since the optimization over  $\boldsymbol{\omega}$  is non-convex, guarantees on this strategy finding the optimal solution have yet to be established. As an alternative, the line-search algorithm presented in [JMGS21] can also be used with small modifications. Under the assumptions of Theorem 1 in [JMGS21], with the additional assumption that  $T$  is a bounded random variable, we inherit their guarantees on the bound of the conditional average potential outcome.

The upper and lower bounds for the CAPO function under treatment  $T = t$  and sensitivity parameter  $\Lambda$  can be estimated for any observed covariate value,  $\mathbf{X} = \mathbf{x}$ , as

$$\underline{\mu}(\mathbf{x}, t; \Lambda, \boldsymbol{\theta}) = \mu(w(y, \cdot; \underline{\omega}); \mathbf{x}, t, \Lambda, \boldsymbol{\theta}),$$

and

$$\bar{\mu}(\mathbf{x}, t; \Lambda, \boldsymbol{\theta}) = \mu(w(y, \cdot; \bar{\omega}); \mathbf{x}, t, \Lambda, \boldsymbol{\theta}).$$

The upper and lower bounds for the APO (dose-response) function under treatment  $T = t$  and sensitivity parameter  $\Lambda$  can be estimated over any set of observed covariates  $\mathcal{D}_{\mathbf{x}} = \{\mathbf{x}_i\}_{i=1}^n$ , as

$$\begin{aligned}\underline{\mu}(t; \Lambda, \boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n \underline{\mu}(\mathbf{x}_i, t; \Lambda, \boldsymbol{\theta}), \quad \mathbf{x}_i \in \mathcal{D}_{\mathbf{x}}, \\ \bar{\mu}(t; \Lambda, \boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n \bar{\mu}(\mathbf{x}_i, t; \Lambda, \boldsymbol{\theta}), \quad \mathbf{x}_i \in \mathcal{D}_{\mathbf{x}}.\end{aligned}$$

## G Datasets

### G.1 Synthetic

$$\begin{aligned}u &:= N_u, \\ x &:= N_x, \\ t &:= N_t,\end{aligned}\tag{49}$$

$$y_t := t + \mathbf{x} \exp(-tx) - \gamma_y(u - 0.5) * (0.5 * x + 1) + N_y,$$

where,  $N_u \sim p(u) := \text{Bern}(u \mid 0.5)$ ,  $N_x \sim p(x) := \text{Unif}[x \mid 0.1, 2.0]$ ,  $N_t \sim p(t \mid x, u) := \text{Beta-Binomial}(t \mid n = 100, \alpha = x + \gamma_t u, \beta = 1)$ , and  $N_y \sim \mathcal{N}(0, 0.04)$ . For the results in this paper  $\gamma_t = 0.3$  and  $\gamma_y = 0.5$ .

The ground truth ratio,  $\lambda = \frac{p(t|x)}{p(t|x, u)}$ , is then given by,

$$\begin{aligned}\lambda^*(t, x, u) &= \frac{\mathbb{E}_{p(u)}[p(t \mid x, u)]}{p(t \mid x, u)} \\ &= \frac{\sum_{u'=0}^1 0.5 * \text{Beta-Binomial}(t \mid n = 100, \alpha = x + \gamma_t u', \beta = 1)}{\text{Beta-Binomial}(t \mid n = 100, \alpha = x + \gamma_t u, \beta = 1)}\end{aligned}\tag{50}$$

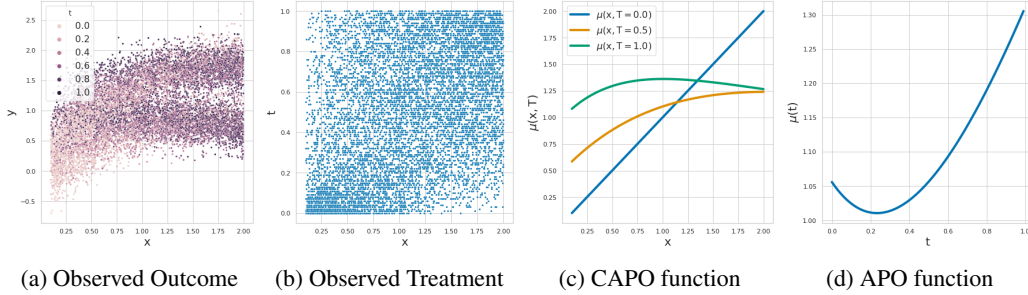


Figure 5: Synthetic data with hidden confounding

### G.2 Observations of clouds and aerosol

The Moderate Resolution Imaging Spectroradiometer (MODIS) instrument aboard the Aqua satellite observes the Earth twice daily at  $\sim 1 \text{ km} \times 1 \text{ km}$  resolution native resolution (Level 1) [BP06]. We used the daily mean,  $1^\circ \times 1^\circ$  gridded version (Level 2) in order to somewhat homogenize our observations of clouds and the atmosphere confined to a region off the coast of South America in the Pacific basin. MODIS observations are fed into the Modern-Era Retrospective analysis for Research and Applications version 2 (MERRA-2) real-time model in order to emulate the atmosphere and its components, such as aerosol [GMS<sup>+</sup>17]. Aerosol optical depth at 550nm from MERRA-2 is derived from MODIS observations of aerosol from multiple satellites (Terra, Aqua, Suomi-NPP), with corrections for sun glint and near-cloud optical effects [BAC<sup>+</sup>15]. We collocated all gridded observations of clouds and reanalysis aerosol with our meteorological proxies of the environment (EIS, SST, w500, RH700, RH850), then normalized our features before feeding them into the model.

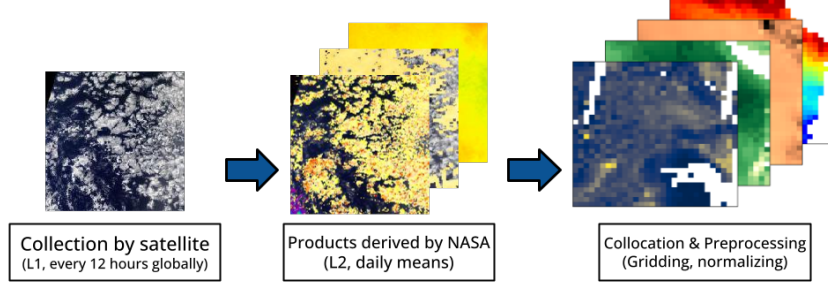


Figure 6: Workflow of observed clouds from satellite to ingestion by model.

Table 1: Sources of satellite observations.

Product name	Description
Cloud optical depth $\tau$	MODIS (1.6, 2.1, 3.7 $\mu\text{m}$ )
Precipitation	NOAA CMORPH
Sea Surface Temperature	NOAA WHOI
Vertical Motion	MERRA-2
Estimated Inversion Strength	MERRA-2
Relative Humidity	MERRA-2
Aerosol Optical Depth	MERRA-2

## H Implementation Details

Experiments 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. Code is written in python. Packages used include PyTorch [PGM<sup>+</sup>19], scikit-learn [PVG<sup>+</sup>11], Ray [MNW<sup>+</sup>18], NumPy, SciPy, and Matplotlib. We use ray tune [LLN<sup>+</sup>18] with HyperBand Bayesian Optimization [FKH17] search algorithm to optimize our network hyper-parameters. The hyper-parameters we consider are accounted for in Table 2. The final hyper-parameters used are given in Table 3. The hyper-parameter optimization objective is the batch-wise Pearson correlation averaged across all outcomes of the validation data for a single dataset realization with random seed 1331. All experiments reported can be completed in 30 hours using this setup.

Hyper-parameter	Search Space
hidden units	tune.qlograndint(32, 512, 32)
network depth	tune.randint(2, 5)
gmm components	tune.randint(1, 32)
attention heads	tune.randint(1, 8)
negative slope	tune.quniform(0.0, 0.5, 0.01)
dropout rate	tune.quniform(0.0, 0.5, 0.01)
layer norm	tune.choice([True, False])
batch size	tune.qlograndint(32, 256, 32)
learning rate	tune.quniform(1e-4, 1e-3, 1e-4)

Table 2: Hyper-parameter search space

### H.1 Model Architecture

The general model architecture is shown in Figure 7. The models are neural-network architectures with two basic components: a feature extractor,  $\phi(\mathbf{x}; \boldsymbol{\theta})$  ( $\phi$ , for short), and a conditional outcome prediction block  $f(\phi, \mathbf{t}; \boldsymbol{\theta})$ , or density estimator. The covariates  $\mathbf{x}$  (represented in blue) are given as input to the feature extractor, whose output is concatenated with the treatment  $\mathbf{t}$  (represented in purple)

Hyper-parameter	Synthetic	ACCE NN	ACCE Transformer
hidden units	96	256	256
network depth	4	3	3
gmm components	24	24	24
attention heads	NA	NA	4
negative slope	0.05	0.04	0.01
dropout rate	0.04	0.2	0.5
layer norm	False	False	False
batch size	32	2048	32
learning rate	0.0015	1e-4	2e-4

Table 3: Final hyper-parameters for each dataset/model

and given as input to the density estimator which outputs a Gaussian mixture density  $p(y | t, \mathbf{x}, \theta)$  from which we can sample to obtain samples of the outcomes (represented in **red**). Models are optimized by maximizing the log-likelihood,  $\log p(y | t, \mathbf{x}, \theta)$ , using mini-batch stochastic gradient descent.

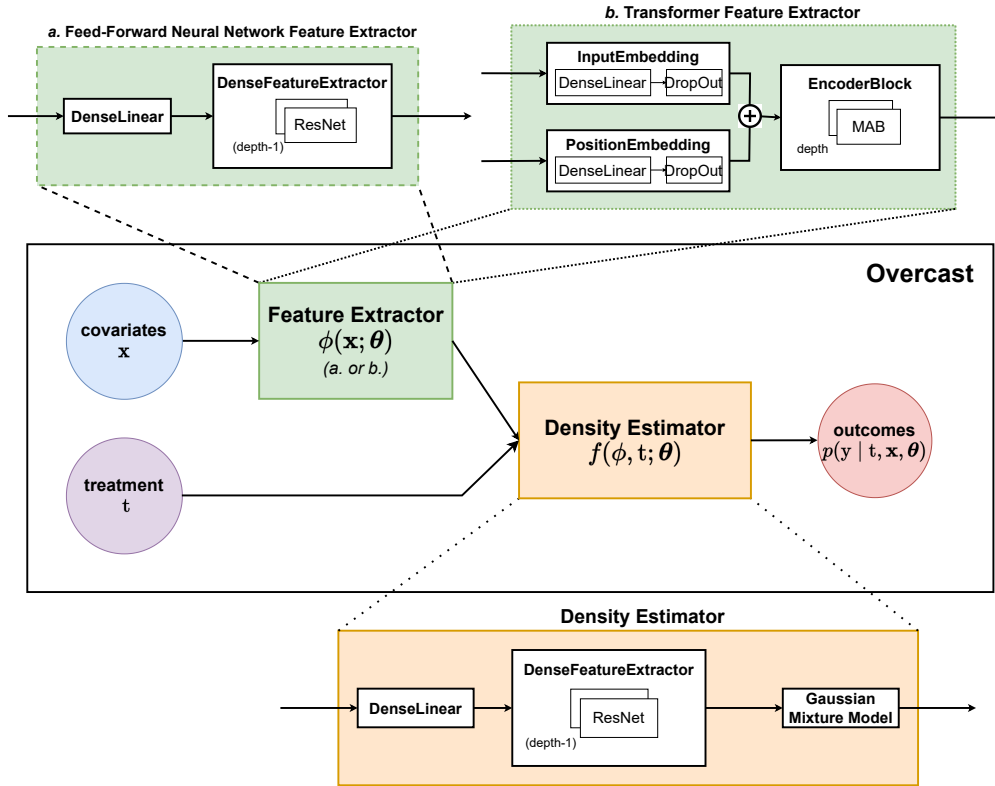


Figure 7: **Overcast model architecture.** The inputs are represented by circles, in **blue** the covariates, and in **purple** the treatment. In the **red** circle is the output of the model, the outcomes distribution. The model has different feature extractors (in **green**) for the feed-forward neural network and the transformer. It has a single density estimator (in **orange**).

### H.1.1 Feature extractor

The feature extractor design is problem and data specific. In our case, we look at using both a simple feed-forward neural network and also a transformer. The transformer has the advantage of allowing us to model the spatio-temporal correlations between the covariates on a given day using the geographical coordinates of the observations as positional encoding. This is interesting when

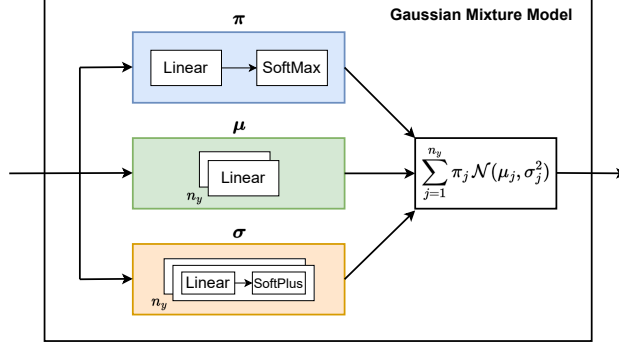


Figure 8: **Overcast Gaussian mixture model.** The mixing coefficients  $\pi$  are estimated with a linear layer and a SoftMax layer, to obtain  $\tilde{\pi}$ , represented in **blue** in the figure. The vector of means of the Gaussian kernels  $\tilde{\mu}$  is obtained by  $n_y$  linear layers (in **green** in the diagram), whilst the vector of variances  $\tilde{\sigma}$  is obtained by  $n_y$  blocks of linear layers and SoftPlus layers (in **orange** in the diagram).

studying ACI because confounding may be latent in the relationships between neighboring variables. Typically, environmental processes (which is one source of confounding) are dependent upon the spatial distribution of clouds, humidity and aerosol, and this feature extractor may capture these confounding effects better.

## H.2 Density Estimator

The conditional outcome prediction block, relies on a  $n_y$  component Gaussian mixture density represented in Figure 8. It outputs:

$$p(y \mid t, \mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^{n_y} \tilde{\pi}_j(\phi, t; \boldsymbol{\theta}) \mathcal{N}(y \mid \tilde{\mu}_j(\phi, t; \boldsymbol{\theta}), \tilde{\sigma}_j^2(\phi, t; \boldsymbol{\theta})),$$

and

$$\tilde{\mu}(\mathbf{x}, t; \boldsymbol{\theta}) = \sum_{j=1}^{n_y} \tilde{\pi}_j(\phi, t; \boldsymbol{\theta}) \tilde{\mu}_j(\phi, t; \boldsymbol{\theta}),$$

where  $\mathcal{N}(\cdot \mid \mu, \sigma^2)$  denotes a normal distribution with mean  $\mu$  and variance  $\sigma^2$ .

## I Additional Results

### I.1 Synthetic

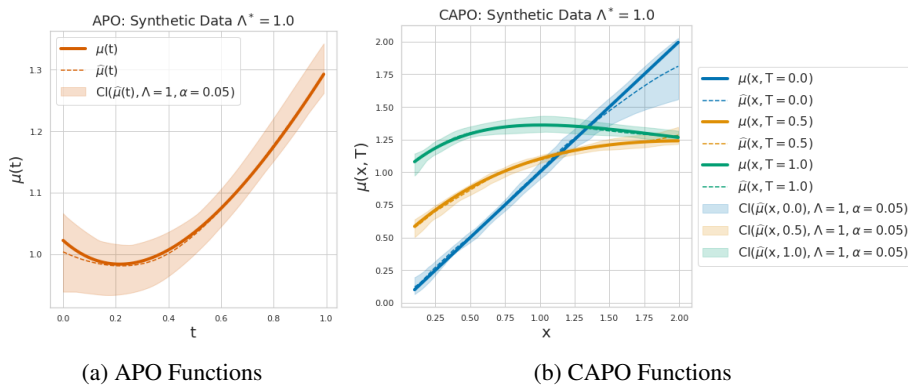


Figure 9: Investigating statistical uncertainty using unconfounded synthetic data.

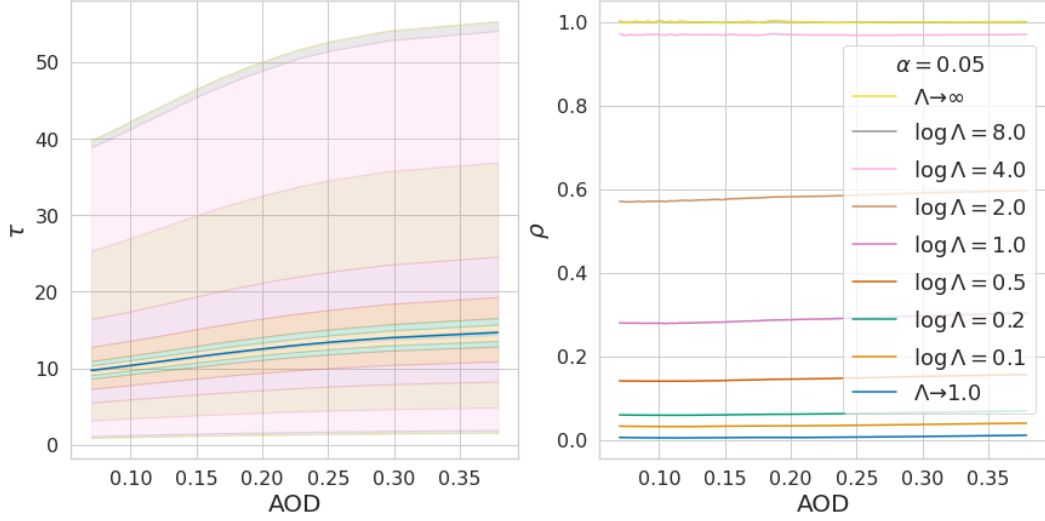


Figure 10: Interpreting  $\Lambda$  as a proportion ( $\rho$ ) of the unexplained range of  $Y_t$  attributed to unobserved confounding variables.

## I.2 Aerosol-Cloud-Climate Effects

In Figure 10 we show how  $\Lambda$  can be interpreted as the proportion,  $\rho$ , of the unexplained range of  $Y_t$  attributed to unobserved confounding variables. In the left figure, we plot the corresponding bounds for increasing values of  $\Lambda$  of the predicted AOD- $\tau$  dose-response curves. In the right figure we plot the  $\rho$  value for each  $\Lambda$  at each value of  $t$ . For the curves reported in Section 5.2: we find that  $\Lambda = 1.1$  leads to  $\rho \approx 0.04$ ,  $\Lambda = 1.2$  leads to  $\rho \approx 0.07$ , and  $\Lambda = 1.6$  leads to  $\rho \approx 0.15$ . This shows that when we let just a small amount of the unexplained range of  $Y_t$  be attributed to unobserved confounding, the range of the predicted APO curves become quite wide. If we were to completely relax the no-hidden-confounding assumption, the entire range seen in Figure 10 Left would be plausible for the APO function. This range dwarfs the predicted APO curve. These results highlight the importance of reporting such sensitivity analyses.

In Figure 11 we show additional dose response curves for cloud optical thickness ( $\tau$ ), water droplet effective radius ( $r_e$ ), and liquid water path (LWP).

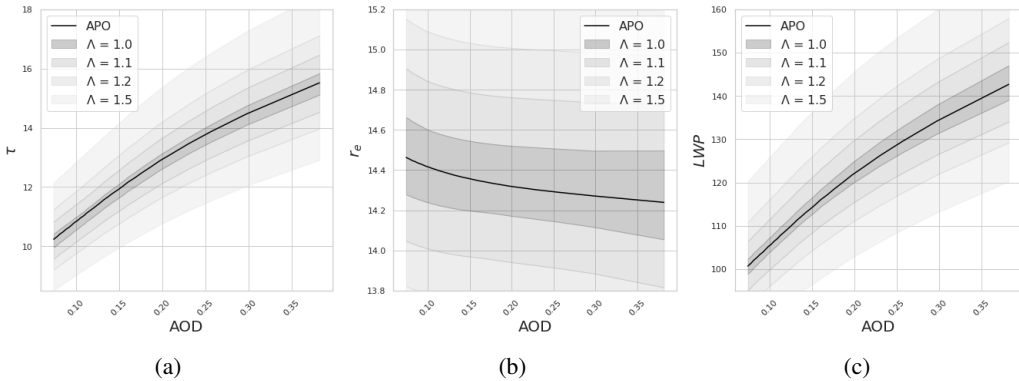


Figure 11: Average dose-response curves for other cloud properties. a) Cloud optical depth. b) Water droplet effective radius. c) Liquid water path.

In Figure 12 we show additional scatter plots comparing the neural network and transformer models for cloud optical thickness ( $\tau$ ), water droplet effective radius ( $r_e$ ), and liquid water path (LWP).

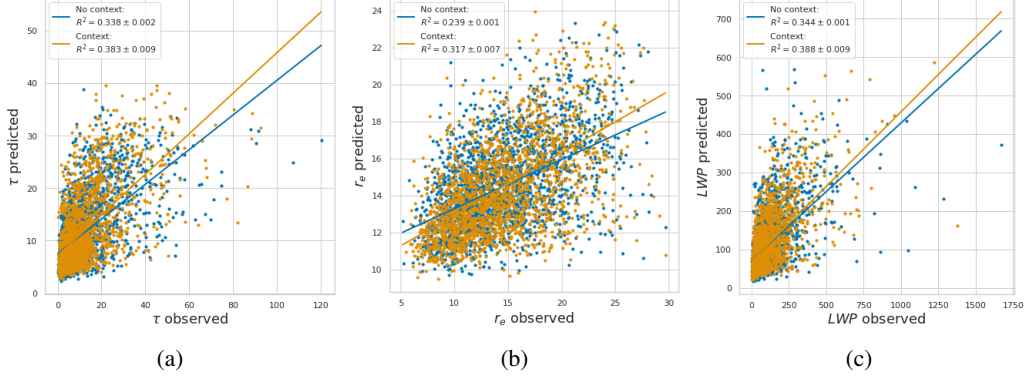


Figure 12: Comparing transformer to feed-forward feature extractor at predicting cloud properties given covariates and AOD. a) Cloud optical depth. b) Water droplet effective radius. c) Liquid water path. We see a significant improvement in pearson correlation ( $R^2$ ) in each case.

### I.3 $\omega_{500}$ experiment

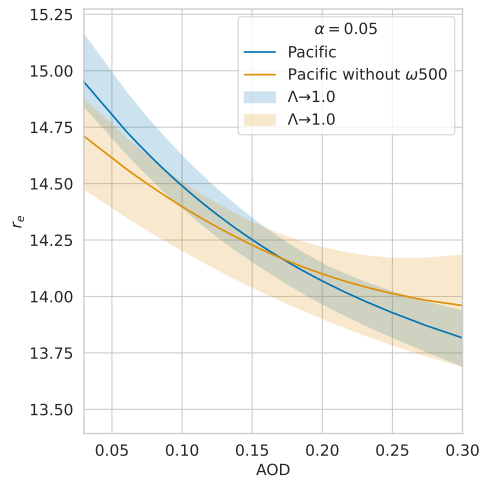
The Overcast models make use of expert knowledge about ACI to select the covariates. Ideally, they would include pressure profiles, temperature profiles and supersaturation since these are directly involved in cloud processes and impact the quality of AOD measurements as a proxy for aerosol concentration. Unfortunately, they are impossible to retrieve from satellite data, so we rely on meteorological proxies like relative humidity, sea surface temperature, inversion strengths, and vertical motion. Relying on these proxies however results in ignorability violations, which must be accounted for in the parameter  $\Lambda$  in order to derive appropriate plausible ranges of outcomes.

In the experiment that follows, we are removing a confounding variable from the model, therefore inducing hidden confounding. The covariate we remove is vertical motion at 500 mb, denoted by  $\omega_{500}$ . This experiment helps us gain some intuition about the influence of the parameter  $\Lambda$  and how it relates to the inclusion of confounding variables in the model.

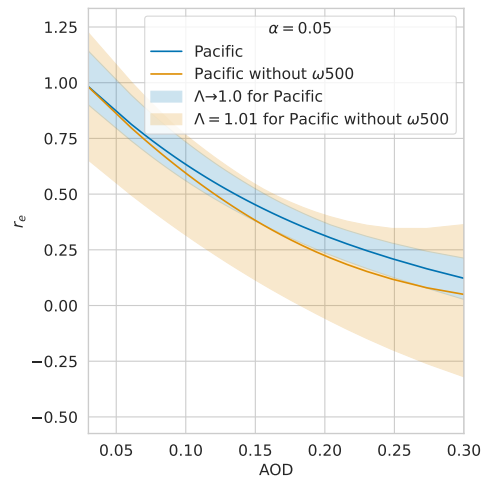
In Figure 13 we compare the same region with different covariates to identify an appropriate  $\Lambda$ . We fit one model on data from the Pacific (blue) and one model from the Pacific omitting  $\omega_{500}$  from the covariates (orange). The shaded bounds in blue are the ignorance region for  $\Lambda \rightarrow 1$  for the Pacific. We then find the  $\Lambda$  that results in an ignorance interval around the Pacific omitting  $\omega_{500}$  that covers the Pacific model prediction. From this, we can infer how the parameter  $\Lambda$  relates to the inclusion of covariates in the model. We show that we need to set  $\Lambda = 1.01$  to account for the fact that we are omitting  $\omega_{500}$  from our list of covariates. We also note that the slopes of the dose-response curves are slightly different, with worse predictions when omitting  $\omega_{500}$  from the covariates, as expected.

This work attempts to set a new methodology for setting  $\Lambda$  which can be summarised as followed. Working with two datasets, which vary in only aspect, we train two different models: (i), the control model, and (ii), the experimental model. After training both models, we plot the dose-response curves for (i) and (ii) on the same plot. We can compare the shape and slope of these curves as well as their uncertainty bounds under the unconfoundedness assumption by plotting the ignorance region for  $\Lambda \rightarrow 1$  for both models. Then, we are interested in setting  $\Lambda$  for model (ii) such that the uncertainty bounds cover the entire ignorance region of model (i) under the unconfoundedness assumption. For this, we are interested in comparing the slopes and thus min-max scale both curves.





(a) Unscaled



(b) Scaled by min-max with appropriate  $\Lambda$

Figure 13: Dose-response curves with or without vertical motion at 500 mb ( $\omega_{500}$ ) as a covariate.