# Supplementary Material: Alleviating Privacy Attacks via Causal Learning

## A  Generalization Properties of Causal Models

### A.1  Generalization over Different Distributions

We provide proofs for generalization properties of causal model over different distributions and a single datapoint.

**Theorem 1.** *Consider a structural causal graph* $G$ *that connects* $X$ *to* $Y$*, and causal features* $X_C \subseteq X$ *where* $X_C$ *represent the parents of* $Y$ *under* $G$*. Let* $P(X, Y)$ *and* $P^*(X, Y)$ *be two distributions with arbitrary* $P(X)$ *and* $P^*(X)$*, having overlap such that* $P(X = x) > 0$ *whenever* $P^*(X = x) > 0$*. In addition, the causal relationship between* $X_C$ *and* $Y$ *is preserved, which implies that* $P(Y|X_C) = P^*(Y|X_C)$*. Let* $L$ *be a symmetric loss function that obeys the triangle inequality (such as L2 loss), and let* $f : X_C \to Y$ *be the optimal predictor among all hypotheses using* $X_C$ *features under* $L$*, i.e.,* $f = \arg\min_h L_{x_c}(y, h(x_c))$ *for all* $x_c$*, and thus* $f$ *depends only on* $\Pr(Y|X_C)$ *(e.g.,* $f := \mathbb{E}[Y|X_C]$ *for L2 loss). Further, assume that* $\mathcal{H}_C$ *represents the set of causal models* $h_c : X_C \to Y$ *that may use all causal features and* $\mathcal{H}_A$ *represent the set of associational models* $h_a : X \to Y$ *that may use all available features, such that* $f \in \mathcal{H}_C$ *and* $\mathcal{H}_C \subseteq \mathcal{H}_A$*.*

1. *When generation of* $Y$ *is deterministic,* $y = f(X_C)$ *(e.g., when* $Y|X_C$ *is almost surely constant), the* ODE *loss for a causal model* $h_c \in \mathcal{H}_C$ *is bounded by:*

$$ODE_{P,P^*}(h_c, y) = \mathcal{L}_{P^*}(h_c, y) - \mathcal{L}_{S\sim P}(h_c, y)$$
$$\leq disc_{L,\mathcal{H}_C}(P, P^*) + IDE_P(h_c, y) \quad (4)$$

*Further, for any* $P$ *and* $P^*$*, the upper bound of* ODE *from a dataset* $S \sim P(X, Y)$ *to* $P^*$*(called* ODE−Bound*) for a causal model* $h_c \in \mathcal{H}_C$ *is less than or equal to the upper bound* ODE−Bound *of an associational model* $h_a \in \mathcal{H}_A$*, with probability at least* $(1 - \delta)^2$*.*

$$ODE{-}Bound_{P,P^*}(h_c, y; \delta) \leq ODE{-}Bound_{P,P^*}(h_a, y; \delta)$$

2. *When generation of* $Y$ *is probabilistic, the* ODE *error for a causal model* $h_c \in \mathcal{H}_C$ *includes additional terms for the loss between* $Y$ *and optimal causal models* $h_{c,P}^{OPT} = h_{c,P^*}^{OPT}$ *on* $P$ *and* $P^*$ *respectively.*

$$ODE_{P,P^*}(h_c, y) \leq disc_{L,\mathcal{H}_C}(P, P^*) + IDE_P(h_c, y) +$$
$$\mathcal{L}_{P^*}(h_{c,P^*}^{OPT}, y) + \mathcal{L}_P(h_{c,P}^{OPT}, y) \quad (5)$$

*While the loss of an associational model can be lower on* $P$*, there always exists a* $P^*$ *such that the worst case*

ODE−Bound *for an associational model is higher than the same for a causal model.*

$$\max_{P^*} ODE{-}Bound_{P,P^*}(h_c, y; \delta) \leq \max_{P^*} ODE{-}Bound_{P,P^*}(h_a, y; \delta)$$

*Proof.* The proof has three parts: General ODE Bound for a model, equivalence of loss-minimizing causal models on P and $P^*$, and finally the two claims from the Theorem.

#### I. GENERAL ODE BOUND

Consider a model $h : X \to Y$ belonging to a set of models $\mathcal{H}$, that was trained on $S \sim P(X, Y)$. From Def. 2 we write,

$$ODE_{P,P^*}(h, y) = \mathcal{L}_{P^*}(h, y) - \mathcal{L}_{S\sim P}(h, y)$$
$$= \mathcal{L}_{P^*}(h, y) - \mathcal{L}_P(h, y) +$$
$$\mathcal{L}_P(h, y) - \mathcal{L}_{S\sim P}(h, y) \quad (8)$$
$$= \mathcal{L}_{P^*}(h, y) - \mathcal{L}_P(h, y) + IDE_P(h, y)$$

where the last equation is to due to Def.1 of the in-distribution generalization error.

Let us denote the optimal loss-minimizing hypotheses over $\mathcal{H}$ for P and $P^*$ as $h_P^{OPT}$ and $h_{P^*}^{OPT}$.

$$h_P^{OPT} = \arg\min_{h\in\mathcal{H}} \mathcal{L}_P(h, y) \qquad h_{P^*}^{OPT} = \arg\min_{h\in\mathcal{H}} \mathcal{L}_{P^*}(h, y) \quad (9)$$

Using the triangle inequality of the loss function, we can write:

$$\mathcal{L}_{P^*}(h, y) \leq \mathcal{L}_{P^*}(h, h_P^{OPT}) + \mathcal{L}_{P^*}(h_P^{OPT}, y) \quad (10)$$

And,

$$\mathcal{L}_P(h, y) \geq \mathcal{L}_P(h, h_P^{OPT}) - \mathcal{L}_P(h_P^{OPT}, y)$$
$$\Rightarrow -\mathcal{L}_P(h, y) \leq -\mathcal{L}_P(h, h_P^{OPT}) + \mathcal{L}_P(h_P^{OPT}, y) \quad (11)$$

Thus, combining Eqns. 8, 10 and 11, we obtain,

$$ODE_{P,P^*}(h, y)$$
$$\leq IDE_P(h, y) + \mathcal{L}_{P^*}(h, h_P^{OPT}) +$$
$$\mathcal{L}_{P^*}(h_P^{OPT}, y) - \mathcal{L}_P(h, h_P^{OPT}) + \mathcal{L}_P(h_P^{OPT}, y)$$
$$= IDE_P(h, y) + (\mathcal{L}_{P^*}(h, h_P^{OPT}) - \mathcal{L}_P(h, h_P^{OPT})) +$$
$$\mathcal{L}_{P^*}(h_P^{OPT}, y) + \mathcal{L}_P(h_P^{OPT}, f)$$
$$\leq IDE_P(h, y) + disc_{L,\mathcal{H}}(P, P^*) +$$
$$\mathcal{L}_{P^*}(h_P^{OPT}, y) + \mathcal{L}_P(h_P^{OPT}, y)$$

$$(12)$$

where the last inequality is due to the definition of discrepancy distance (Definition 3).

Below we show that Eqn. 12 divides the out-of-distribution generalization error of a model $h$ in four parts. As defined in the Theorem statement, $\mathcal{H}_C$ refers to the class of models that uses all causal features ($X_C$), parents of $Y$ over the structural causal graph; and $\mathcal{H}_A$ refers to the class of associational models that may use all or a subset of all available features.

1. $\text{IDE}_P(h, y)$ denotes the in-distribution error of $h$. This can be bounded by typical generalization bounds, such as the uniform error bound that depends only on the VC dimension and sample size of $S$ (Shalev-Shwartz & Ben-David, 2014). Using a uniform error bound based on the VC dimension, we obtain, with probability at least $1 - \delta$,

$$\text{IDE} \leq \sqrt{8 \frac{\text{VCdim}(\mathcal{H})(\ln(2|S|) + 1) + \ln(4/\delta)}{|S|}} \quad (13)$$
$$= \text{IDE-Bound}(\mathcal{H}, S)$$

Since $\mathcal{H}_C \subseteq \mathcal{H}_A$, VC-dimension of causal models is not greater than that of associational models. Thus,

$$\text{VCDim}(\mathcal{H}_C) \leq \text{VCDim}(\mathcal{H}_A) \Rightarrow \text{IDE-Bound}(\mathcal{H}_C, \mathcal{S})$$
$$\leq \text{IDE-Bound}(\mathcal{H}_A, \mathcal{S}) \quad (14)$$

2. $\text{disc}_{L, \mathcal{H}}(P, P^*)$ denotes the distance between the two distributions. Given two distributions, the discrepancy distance does not depend on $h$, but only on the model class $\mathcal{H}$. From Definition 3, discrepancy distance is the maximum quantity over all pairs of models in a model class. Since $\mathcal{H}_C \subseteq \mathcal{H}_A$, we obtain that:

$$\text{disc}_{L, \mathcal{H}_C}(P, P^*) \leq \text{disc}_{L, \mathcal{H}_A}(P, P^*) \quad (15)$$

3. $\mathcal{L}_P(h_P^{\text{OPT}}, y)$ measures the error of the loss-minimizing model on P, when evaluated on P. While $h_P^{\text{OPT}}$ is optimal, there can still be error due to the true labeling function $f$ being outside the model class $\mathcal{H}$, or irreducible error due to probabilistic generation of $Y$.

4. $\mathcal{L}_{P^*}(h_P^{\text{OPT}}, y)$ measures the error of the loss-minimizing model on P, when evaluated on $P^*$. In addition to the reasons cited above, this error can be due to differences in both $\Pr(X)$ and $\Pr(Y|X)$ between P and $P^*$: change in the marginal distribution of inputs X, and/or change in the conditional distribution of $Y$ given X.

## II. SAME LOSS-MINIMIZING CAUSAL MODEL OVER P AND P*

Below we show that for a given distribution P and another distribution $P^*$ such that $P(Y|X_C) = P^*(Y|X_C)$, the loss minimizing model is the same for causal models ($h_{c,P}^{\text{OPT}} = h_{c,P^*}^{\text{OPT}}$), but not necessarily for associational models.

**Causal Model.** Given a structural causal network, let us construct a model using all parents of $X_C$ of $Y$. By property of the structural causal network, $X_C$ includes all parents of $Y$ and therefore there are no backdoor paths. Using Rule 2 of do-calculus from Pearl (2009):

$$\Pr(Y|\text{do}(X_c = x_c)) = P(Y|X_C = x_c) = P^*(Y|X_C = x_c) \quad (16)$$

where the last equality is assumed since data from $P^*$ also shares the same causal graph. Defining $h_{c,P}^{\text{OPT}} = \arg\min_{h_c \in \mathcal{H}_C} \mathcal{L}_P(h_c, y)$ and $h_{c,P^*}^{\text{OPT}} = \arg\min_{h_c \in \mathcal{H}_C} \mathcal{L}_{P^*}(h_c, y)$, we can write,

$$h_{c,P}^{\text{OPT}} = \arg\min_{h \in \mathcal{H}_C} \mathcal{L}_P(h, y)$$
$$= \arg\min_{h \in \mathcal{H}_C} \mathbb{E}_{P(x_c, y)} L(h(x_c), y) = f_{P(Y|X_C)} \quad (17)$$

since $f = \arg\min_h L_x(h(x_c), y)$ for all $x_c$ and thus does not depend on $\Pr(X_C)$, and $f \in \mathcal{H}_C$. Similarly, for $h_{c,P^*}^{\text{OPT}}$, we can write:

$$h_{c,P^*}^{\text{OPT}} = \arg\min_{h \in \mathcal{H}_C} \mathcal{L}_{P^*}(h, y)$$
$$= \arg\min_{h \in \mathcal{H}_C} \mathbb{E}_{P^*(x_c, y)} L(h(x_c), y) = f_{P^*(Y|X_C)} \quad (18)$$

Since $P(Y|X_C) = P^*(Y|X_C)$, we obtain,

$$f_{P(Y|X_C)} = f_{P^*(Y|X_C)} \Rightarrow h_{c,P}^{\text{OPT}} = h_{c,P^*}^{\text{OPT}} \quad (19)$$

**Associational Model.** In contrast, an associational model may use a subset $X_A \subseteq X$ that may not include all parents of $Y$, or may include parents but also include other extraneous variables. Following the derivation for causal models, let us define $h_{a,P}^{\text{OPT}} = \arg\min_{h_a \in \mathcal{H}_A} \mathcal{L}_P(h_a, y)$ and $h_{a,P^*}^{\text{OPT}} = \arg\min_{h_a \in \mathcal{H}_A} \mathcal{L}_{P^*}(h_a, y)$, we can write,

$$h_{a,P}^{\text{OPT}} = \arg\min_{h \in \mathcal{H}_A} \mathcal{L}_P(h, y)$$
$$= \arg\min_{h \in \mathcal{H}_A} \mathbb{E}_{P(x_a, y)} L(h(x_a), y) = f_{P(X_A, Y)} \quad (20)$$

where we define $f_A$ as, $f_A = \arg\min_h L_x(h(x_a), y)$ for any $x_a$. Similarly, for $h_{a,P^*}^{\text{OPT}}$, we can write:

$$h_{a,P^*}^{\text{OPT}} = \arg\min_{h \in \mathcal{H}_A} \mathcal{L}_{P^*}(h, y)$$
$$= \arg\min_{h \in \mathcal{H}_A} \mathbb{E}_{P^*(x_a, y)} L(h(x_a), y) = f_{P^*(X_A, Y)} \quad (21)$$

Now, in general,

$$P(X_A, Y) \neq P^*(X_A, Y) \Rightarrow f_{P(X_A, Y)} \neq f_{P^*(X_A, Y)}$$

Even if the optimal associational model $f_A \in \mathcal{H}_A$ (as we assumed for causal models), and thus $f_{P(X_A, Y)} = f_{P(Y|X_A)}$ and

$f_{P^*(X_A, Y)} = f_{P^*(Y|X_A)}$, they are not the same since $P(Y|X_A) \neq P^*(Y|X_A)$. Therefore we obtain,

$$f_{P(Y|X_A)} \neq f_{P^*(Y|X_A)} \Rightarrow h_{a,P}^{OPT} \neq h_{a,P^*}^{OPT} \tag{22}$$

That said, since $X_C \subset X$, it is possible that $X_A = X_C$ for some $X$ and $\mathcal{H}$, and thus the loss-minimizing associational model includes only the causal features of $Y$. Then $h_{a,P}^{OPT} = h_{a,P^*}^{OPT}$. In general, though, $h_{a,P}^{OPT} \neq h_{a,P^*}^{OPT}$.

## IIIa. CLAIM 1

As a warmup, consider the case when $Y$ is generated deterministically. That is, the optimal model $f$ has zero error. Then, both the loss-minimizing causal model and loss-minimizing associational model have zero error when evaluated on the same distribution that they were trained on. Thus, $\mathcal{L}_P(h_{c,P}^{OPT}, y) = \mathcal{L}_{P^*}(h_{c,P^*}^{OPT}, y) = 0$. Similarly, $\mathcal{L}_P(h_{a,P}^{OPT}, y) = 0$. (Note that here we consider only those cases where $f_{P(Y|X)} \in \mathcal{H}_A$ and $f_{P^*(Y|X)} \in \mathcal{H}_A$ for a fair comparison; otherwise, the error bound for $h_a \in \mathcal{H}_A$ is trivially larger than that for $h_c \in \mathcal{H}_C$).

Further, for a causal model, using Equation 19, we obtain:

$$\mathcal{L}_{P^*}(h_{c,P}^{OPT}, y) = \mathcal{L}_{P^*}(h_{c,P^*}^{OPT}, y) = 0 \tag{23}$$

However, the same does not hold for associational models: $\mathcal{L}_{P^*}(h_{a,P}^{OPT}, y)$ need not be zero.

We now present the loss bounds. Using Equations 19 and 23, we write Equation 12 for a causal model as:

$$\begin{aligned} \text{ODE}_{P,P^*}(h_c, y) &= \mathcal{L}_{P^*}(h_c, y) - \mathcal{L}_{S \sim P}(h_c, y) \\ &\leq \text{disc}_{L, \mathcal{H}_C}(P, P^*) + \text{IDE}_P(h_c, y) \end{aligned} \tag{24}$$

For an associational model, we obtain,

$$\begin{aligned} \text{ODE}_{P,P^*}(h_a, y) &= \mathcal{L}_{P^*}(h_a, y) - \mathcal{L}_{S \sim P}(h_a, y) \\ &\leq \text{disc}_{L, \mathcal{H}_A}(P, P^*) + \text{IDE}_P(h_a, y) \\ &\quad + \mathcal{L}_{P^*}(h_{a,P}^{OPT}, y) \end{aligned} \tag{25}$$

Using Eqn. 13 that bounds IDE with probability $1 - \delta$, and Eqns. 14 and 15 that compare IDE-Bound and discrepancy distance between causal and associational model classes, we can rewrite Eqn. 24. With probability at least $1 - \delta$:

$$\begin{aligned} \text{ODE}_{P,P^*}(h_c, y) &\leq \text{disc}_{L, \mathcal{H}_C}(P, P^*) + \text{IDE-Bound}_P(\mathcal{H}_C, S; \delta) \\ &= \text{ODE-Bound}_{P,P^*}(h_c, y; \delta) \\ &\leq \text{disc}_{L, \mathcal{H}_A}(P, P^*) + \text{IDE-Bound}_P(\mathcal{H}_A, S; \delta) \end{aligned} \tag{26}$$

Similarly, for the associational model,

$$\begin{aligned} \text{ODE}_{P,P^*}(h_a, y) &\leq \text{disc}_{L, \mathcal{H}_A}(P, P^*) + \text{IDE-Bound}_P(\mathcal{H}_A, S; \delta) \\ &\quad + \mathcal{L}_{P^*}(h_{a,P}^{OPT}, y) \\ &= \text{ODE-Bound}_{P,P^*}(h_a, y; \delta) \end{aligned} \tag{27}$$

Therefore, comparing Eqn. 26 and 27, we claim for any $P$ and $P^*$, with probability $(1 - \delta)^2$,

$$\text{ODE-Bound}_{P,P^*}(h_c, y; \delta) \leq \text{ODE-Bound}_{P,P^*}(h_a, y; \delta) \tag{28}$$

## IIIb. CLAIM 2

We now consider the general case when $Y$ is generated probabilistically. Thus, even though $f \in \mathcal{H}_C$ and $h_{c,P}^{OPT} = h_{c,P^*}^{OPT} = f$, $\mathcal{L}_P(h_{c,P}^{OPT}, y) \neq 0$ and $\mathcal{L}_{P^*}(h_{c,P^*}^{OPT}, y) \neq 0$.

Using the IDE bound from Eqn. 13, we write Eqn. 12 as,

$$\begin{aligned} \text{ODE}_{P,P^*}(h_c, y) &\leq \text{disc}_{L, \mathcal{H}_C}(P, P^*) + \text{IDE}_P(h_c, y) \\ &\quad + \mathcal{L}_{P^*}(h_{c,P}^{OPT}, y) + \mathcal{L}_P(h_{c,P}^{OPT}, y) \\ &\leq \text{disc}_{L, \mathcal{H}_C}(P, P^*) + \text{IDE-Bound}_P(\mathcal{H}_C, S; \delta) \\ &\quad + \mathcal{L}_{P^*}(h_{c,P}^{OPT}, y) + \mathcal{L}_P(h_{c,P}^{OPT}, y) \\ &= \text{ODE-Bound}_{P,P^*}(h_c, y; \delta) \\ &\leq \text{disc}_{L, \mathcal{H}_A}(P, P^*) + \text{IDE-Bound}_P(\mathcal{H}_A, S; \delta) \\ &\quad + \mathcal{L}_{P^*}(h_{c,P}^{OPT}, y) + \mathcal{L}_P(h_{c,P}^{OPT}, y) \end{aligned} \tag{29}$$
$$\tag{30}$$

where Eqn. 29 uses $h_{c,P}^{OPT} = h_{c,P^*}^{OPT}$ and Eqn. 30 uses inequalities comparing IDE and discrepancy distance from Eqns. 14 and 15.

Similarly, for associational model,

$$\begin{aligned} \text{ODE-Bound}_{P,P^*}(h_a, y) &= \text{disc}_{L, \mathcal{H}_A}(P, P^*) \\ &\quad + \text{IDE-Bound}_P(\mathcal{H}_A, S; \delta) + \mathcal{L}_{P^*}(h_{a,P}^{OPT}, y) + \mathcal{L}_P(h_{a,P}^{OPT}, y) \end{aligned} \tag{31}$$

Now, we compare the last two terms of Equations 30 and 31. Since $\mathcal{H}_C \subseteq \mathcal{H}_A$, loss of the loss-minimizing associational model can be lower than the loss of the causal model trained on the same distribution. Thus, $\mathcal{L}_P(h_{a,P}^{OPT}, y) \leq \mathcal{L}_P(h_{c,P}^{OPT}, y)$.

However, since $h_{a,P}^{OPT} \neq h_{a,P^*}^{OPT}$, loss of the loss-minimizing associational model trained on $P$ can be higher on $P^*$ than the loss of optimal causal model trained on $P^*$ and evaluated on $P^*$. Formally, let $\gamma_1 \geq 0$ be the loss reduction over $P$ due to use of associational model optimized on $P$, compared to the loss-minimizing causal model. Similarly, let $\gamma_2$ be the increase in loss over $P^*$ due to using the associational model optimized over $P$, compared to the loss-minimizing causal model.

$$\gamma_1 = \mathcal{L}_P(h_{c,P}^{OPT}, y) - \mathcal{L}_P(h_{a,P}^{OPT}, y) \tag{32}$$

$$\gamma_2 = \mathcal{L}_{P^*}(h_{a,P}^{OPT}, y) - \mathcal{L}_{P^*}(h_{c,P}^{OPT}, y) \tag{33}$$

Then, Eqn. 31 transforms to,

$$\begin{aligned} \text{ODE}_{P,P^*}(h_a, y) &\leq \text{disc}_{L, \mathcal{H}_A}(P, P^*) + \text{IDE-Bound}_P(\mathcal{H}_A, S; \delta) \\ &\quad + \mathcal{L}_{P^*}(h_{c,P^*}^{OPT}, y) + \mathcal{L}_P(h_{c,P}^{OPT}, y) + \gamma_2 - \gamma_1 \end{aligned} \tag{34}$$

Hence, as long as $\gamma_2 \geq \gamma_1$, we obtain,

$$\texttt{ODE-Bound}_{P,P^*}(h_c, y; \delta) \leq \texttt{ODE-Bound}_{P,P^*}(h_a, y; \delta) \tag{35}$$

Below we show that such a $P^*$ always exists, and further, the worst-case $\max_{P^*} \texttt{ODE-Bound}_{P,P^*}(h, y; \delta)$ is always lower for a causal model than an associational model.

**There exists $P^*$ such that $\gamma_2 \geq \gamma_1$.** The proof is by construction. As an example, consider L1 loss and a distribution P such that the optimal causal model $f$ for an input data point $x^{(i)}$ can be written as,

$$y^{(i)} = f_P(x_C^{(i)}) + \xi_i = f_{P^*}(x_C^{(i)}) + \xi_i \tag{36}$$

where $f(x_C) = h_{c,P}^{OPT} = h_{c,P^*}^{OPT}$ refers to the optimal causal model and is the same for P and $P^*$ (using Eqn. 19). Let $f_P(x_A) = h_{a,P}^{OPT}$ be the optimal associational model over P. We can rewrite $h_{a,P}^{OPT}$ as an arbitrary change from $h_{c,P}^{OPT}$, using $\lambda_{x_A}^{(i)}$ as a parameter that can be different for each data point $x^{(i)}$. That is,

$$h_{a,P}^{OPT}(x^{(i)}) = h_{c,P}^{OPT}(x_C^{(i)}) + \lambda_{x_A}^{(i)} \tag{37}$$

Based on Eqns. 36 and 37, $\gamma_1$ can be written as,

$$\mathcal{L}_P(h_{c,P}^{OPT}, y) = \mathbb{E}_P[|\xi|]$$
$$\mathcal{L}_P(h_{a,P}^{OPT}, y) = \mathbb{E}_P[|\xi - \lambda_{x_A}|] \tag{38}$$
$$\Rightarrow \gamma_1 = \mathbb{E}_P[|\lambda_{x_A}|]$$

Then, we can construct a $P^*(X, Y)$ such that (i) the relationship $(\Pr(Y|X_A))$ between $x_A$ and $y$ is reversed, and (ii) $\Pr(X)$ is chosen such that $\mathbb{E}_{P^*}[\lambda_{x_A}] \geq \mathbb{E}_P[\lambda_{x_A}]$ (e.g., by assigning higher probability weights to data points $i$ where $|\lambda_{x_A}^{(i)}|$ is high). That is, consider a $P^*$ such that we can write $h_{a,P^*}^{OPT}$ as,

$$h_{a,P^*}^{OPT}(x^{(i)}) = h_{c,P^*}^{OPT}(x_C^{(i)}) - \lambda_{x_A}^{(i)} \tag{39}$$

On such $P^*$, the loss-minimizing causal model remains the same. However, the loss of the associational model $h_{a,P}^{OPT}$ on such $P^*$ increases and can be written as:

$$\mathcal{L}_{P^*}(h_{c,P}^{OPT}, y) = \mathbb{E}_{P^*}[|\xi|]$$
$$\mathcal{L}_{P^*}(h_{a,P}^{OPT}, y) = \mathbb{E}_{P^*}[|\xi + \lambda_{x_A}|] \tag{40}$$
$$\Rightarrow \gamma_2 = \mathbb{E}_{P^*}[|\lambda_{x_A}|]$$

From condition (ii) above, $\mathbb{E}_{P^*}[\lambda_{x_A}] \geq \mathbb{E}_P[\lambda_{x_A}]$, thus $\gamma_2 \geq \gamma_1$.

Note that we did not use any special property of the L1 Loss above. In general, we can write the loss-minimizing function $h_{a,P}^{OPT}$ as adding some arbitrary value $\lambda_{x_A}^{(i)}$ to $h_{c,P}^{OPT}(x_c^{(i)})$; and then construct a $P^*$ such that the relationship $\Pr(Y|X_A)$ is reversed on $P^*$, and thus $h_{a,P^*}^{OPT}$ subtracts the same value. Further, the input data distribution $P^*(X)$ can be chosen

such that $\gamma_2 \geq \gamma_1$. That is, for a loss L, we can choose $\lambda$ such that $\mathcal{L}_{P^*}(h_{a,P}^{OPT}, y; \lambda) - \mathcal{L}_{P^*}(h_{c,P^*}^{OPT}, y) \geq \mathcal{L}_P(h_{c,P}^{OPT}, y) - \mathcal{L}_P(h_{a,P}^{OPT}, y; \lambda)$.

Hence, there exists a $P^*$ such that $\gamma_2 \geq \gamma_1$, and thus,

$$\texttt{ODE-Bound}_{P,P^*}(h_c, y; \delta) \leq \texttt{ODE-Bound}_{P,P^*}(h_a, y; \delta) \tag{41}$$

**Worst case ODE-bound for causal model is lower.** Finally, we show that the for a fixed P, the worst case $\texttt{ODE-Bound}$ also follows Eqn. 41. Looking at Eqns. 30 and 31, $\texttt{ODE-Bound}$ will be highest for a $P^*$ such that discrepancy between P and $P^*$ is highest and $\mathcal{L}_{P^*}(h_P^{OPT}, y)$ is highest. Below we show that discrepancy $\texttt{disc}_L(P, P^*)$ increases as $\mathcal{L}_{P^*}(h_P^{OPT}, y)$ increases.

$$\mathcal{L}_{P^*}(h_P^{OPT}, y) = \mathcal{L}_{P^*}(h_P^{OPT}, y) - \mathcal{L}_P(h_P^{OPT}, y) + \mathcal{L}_P(h_P^{OPT}, y)$$
$$\leq \texttt{disc}_L(P, P^*) + \mathcal{L}_P(h_P^{OPT}, y)$$
$$\Rightarrow \texttt{disc}_L(P, P^*) \geq \mathcal{L}_{P^*}(h_P^{OPT}, y) - \mathcal{L}_P(h_P^{OPT}, y) \tag{42}$$

where $\mathcal{L}_P(h_P^{OPT}, y)$ is fixed since P is fixed. Thus, the above equation shows that whenever $\mathcal{L}_{P^*}(h_P^{OPT}, y)$ is high, discrepancy is also high. Hence, for any $P_{max}^*$ that maximizes $\texttt{ODE-Bound}$, $P_{max}^* = \arg\max_{P^*} \texttt{ODE-Bound}_{P,P^*}(h, y; \delta)$, $\mathcal{L}_{P^*}(h_P^{OPT}, y)$ is also maximized.

Now, let us consider causal and associational models, and their respective worst case $P_{max}^*$. To complete the proof, we need to check whether $\gamma_2 \geq \gamma_1$ for such maximal $\mathcal{L}_{P^*}(h_{c,P}^{OPT}, y)$ and $\mathcal{L}_{P^*}(h_{a,P}^{OPT}, y)$. Since $\gamma_2$ increases monotonically with $\mathcal{L}_{P^*}(h_{a,P}^{OPT}, y)$ ( $\mathcal{L}_{P^*}(h_{c,P}^{OPT}, y)$ is bounded by $\max_x L_x(h_{c,P}^{OPT}, y)$), and there exists at least one $P^*$ such that $\gamma_2 \geq \gamma_1$, this implies that $\gamma_2 \geq \gamma_1$ for $P_{max}^*$ too. Therefore, using Equation 41,

$$\max_{P^*} \texttt{ODE-Bound}_{P,P^*}(h_c, y; \delta) \leq \max_{P^*} \texttt{ODE-Bound}_{P,P^*}(h_a, y; \delta) \tag{43}$$

$\square$

### A.2 Generalization over a Single Datapoint

**Theorem 2.** *Consider a causal model $h_{c,S}^{min} : X_C \rightarrow Y$ and an associational model $h_{a,S}^{min} : X \rightarrow Y$ trained on a dataset $S \sim P(X, Y)$ with loss L. Let $(x, y) \in S$ and $(x', y') \notin S$ be two input instances such that they share the same true labelling function on the causal features, $y \sim P(Y|X_C = x)$ and $y' \sim P(Y|X_C = x')$. Then, the worst-case generalization error for the causal model on such $x'$ is less than or equal to that for the associational model.*

$$\max_{x \in S, x'} L_{x'}(h_{c,S}^{min}, y) - L_x(h_{c,S}^{min}, y) \leq \max_{x \in S, x'} L_{x'}(h_{a,S}^{min}, y) - L_x(h_{a,S}^{min}, y)$$

*Proof.* For any model $h$, we can write,

$$\max_{x \in S, x'} L_{x'}(h, y) - L_x(h, y) = \max_{x'} L_{x'}(h, y) - \min_{x \in S} L_x(h, y) \tag{44}$$

since $x'$ and $x$ are independently selected. To prove the main result, we will show that the maximum loss on an

unseen $x'$, $\max_{x'} L_{x'}(h, y)$ is higher for a loss-minimizing associational model than a causal model, and that minimum loss on a training point $x \in S$, $\min_{x \in S} L_x(h, y)$ is lower for the associational model than a causal model.

**Loss on a training data point.**  First, consider loss on $x \in S$, $L_x(h, y)$.

$$h_{c,S}^{\min} = \arg \min_{h \in \mathcal{H}_C} \mathcal{L}_S(h_c, y) = \arg \min_h \frac{1}{N} \sum_{i=1}^{N} L_{x_i}(h, y)$$

$$h_{a,S}^{\min} = \arg \min_{h \in \mathcal{H}_A} \mathcal{L}_S(h_a, y) = \arg \min_h \frac{1}{N} \sum_{i=1}^{N} L_{x_i}(h, y)$$

Since $\mathcal{H}_C \subseteq \mathcal{H}_A$, the average training loss will be lower for the associational model.

$$\mathcal{L}_S(h_{c,S}^{\min}, y) \geq \mathcal{L}_S(h_{a,S}^{\min}, y) \tag{45}$$

Further, under a suitably complex $\mathcal{H}_A$ there exists a $h_{a,S}^{\min}$ such that the loss L is lower for any $x \in S$. Therefore,

$$\min_{x \in S} L_x(h_{c,S}^{\min}, y) \geq \min_{x \in S} L_x(h_{a,S}^{\min}, y) \tag{46}$$

**Loss on an unseen data point.**  Second, consider $L_{x'}(h, y)$. Without loss of generality, let us write the true function for some $(x', y') \sim P^*$ as,

$$y' = h_{c,P^*}^{OPT}(x_c') + \epsilon = h_{c,P}^{OPT}(x_c') + +\epsilon \tag{47}$$

where we use that $h_{c,P}^{OPT} = h_{c,P^*}^{OPT}$. Suppose there is a data point $(x_1', y_1')$ such that the loss L is maximum for $h_{c,S}^{\min}$.

$$\max_{x' \notin S} L_{x'}(h_{c,S}^{\min}, y) = L_{x_1'}(h_{c,S}^{\min}(x_1'), y_1')$$
$$= L_{x_1'}(h_{c,S}^{\min}(x_{c,1}'), h_{c,P}^{OPT}(x_{c,1}') + \epsilon_1) \tag{48}$$

Now for the associational model $h_{a,S}^{\min}$, the corresponding loss on $x_1'$ is,

$$L_{x_1'}(h_{a,S}^{\min}, y) = L_{x_1'}(h_{a,S}^{\min}, h_{c,P}^{OPT} + \epsilon_1) \tag{49}$$

Without loss of generality, we can write the output of the associational model $h_{a,S}^{\min}$ on a particular input $x'$ as,

$$h_{a,S}^{\min}(x') = h_{c,S}^{\min}(x_c') + h_a(x') \tag{50}$$

where $h_a$ is some associational function of x. Therefore the loss on $x_1'$ becomes,

$$L_{x_1'}(h_{a,S}^{\min}, y) = L_{x_1'}(h_{c,S}^{\min} + h_a, h_{c,P}^{OPT} + \epsilon_1) \tag{51}$$

Since $\Pr(Y|X_A)$ can change for different $x' \sim P^*$ (where $X_A = X \setminus X_C$ refers to the associational features), we will show that RHS of Eqn. 49 can always be greater than or

equal to the RHS of Eqn. 48. For ease of exposition, we consider L1 loss below. For a causal model, the loss can be written as,

$$L_{x_1'}(h_{c,S}^{\min}, h_{c,P}^{OPT} + +\epsilon)$$
$$= |h_{c,S}^{\min}(x_{c,1}') - h_{c,P}^{OPT}(x_{c,1}') - \epsilon_1| \tag{52}$$
$$= |h_{c,S}^{\min}(x_{c,1}') - h_{c,P}^{OPT}(x_{c,1}')| + |\epsilon_1|$$

where $x_1'$ (and thus $\epsilon_1$) is chosen such that $\epsilon_1(h_{c,P}^{OPT}(x_{c,1}') - h_{c,S}^{\min}(x_{c,1}')) \geq 0$ which leads to maximum loss. And for the associational model, the loss on the same $(x_1', y_1')$ can be written as,

$$L_{x_1'}(h_{a,S}^{\min}, h_{c,P}^{OPT} + \epsilon_1)$$
$$= |h_{c,S}^{\min}(x_{c,1}') + h_a(x_1') - h_{c,P}^{OPT}(x_{c,1}') - \epsilon_1| \tag{53}$$
$$= |(h_{c,S}^{\min}(x_{c,1}') - h_{c,P}^{OPT}(x_{c,1}')) + (h_a(x_1') - \epsilon_1)|$$

Comparing Eqns. 52 and 53, two cases arise. If $h_a(x_1')\epsilon_1 \leq 0$, then we obtain,

$$L_{x_1'}(h_{a,S}^{\min}, h_{c,P}^{OPT} + \epsilon_1)$$
$$= |h_{c,S}^{\min}(x_{c,1}') - h_{c,P}^{OPT}(x_{c,1}')| + |(h_a(x_1') - \epsilon_1)| \tag{54}$$
$$= |h_{c,S}^{\min}(x_{c,1}') - h_{c,P}^{OPT}(x_{c,1}')| + |h_a(x_1')| + |\epsilon_1|$$

which is greater than maximum loss on $x_1'$ using a causal model (Eqn. 52). Otherwise, if $h_a(x_1')\epsilon_1 > 0$, we can sample a new data point $(x_2', y_2')$ from some other $P^*$ such that its causal features are the same $(x_{c,1}' = x_{c,2}')$ and thus y is the same $(y_2' = y_1' = h_{c,P}^{OPT}(x_{c,1}') + \epsilon_1)$, but its associational features are different $(x_{a,1}' \neq x_{a,2}')$. Specifically, $x_{a,2}'$ is chosen such that $h_a(x_2')\epsilon_1 \leq 0$. Thus we again obtain,

$$L_{x_2'}(h_{a,S}^{\min}, h_{c,P}^{OPT} + \epsilon_1)$$
$$= |(h_{c,S}^{\min}(x_{c,2}') - h_{c,P}^{OPT}(x_{c,2}')) + (h_a(x_2') - \epsilon_1)|$$
$$= |(h_{c,S}^{\min}(x_{c,1}') - h_{c,P}^{OPT}(x_{c,1}')) + (h_a(x_2') - \epsilon_1)| \tag{55}$$
$$= |(h_{c,S}^{\min}(x_{c,1}') - h_{c,P}^{OPT}(x_{c,1}'))| + |(h_a(x_2')| + |\epsilon_1|$$

where the second equality uses $x_{c,2}' = x_{c,1}'$. Combining Eqns. 54 and 55 and comparing to Eqn. 52, we obtain,

$$\max_{x'} L_{x'}(h_{c,S}^{\min}, y) \leq \max_{x'} L_{x'}(h_{a,S}^{\min}, y) \tag{56}$$

Finally, using Eqns. 46 and 56 leads to the main result.

$$\max_{x'} L_{x'}(h_{c,S}^{\min}, y) - \min_{x \in S} L_x(h_{c,S}^{\min}, y)$$
$$\leq \max_{x'} L_{x'}(h_{a,S}^{\min}, y) - \min_{x \in S} L_x(h_{a,S}^{\min}, y)$$
$$\max_{x', x \in S} L_{x'}(h_{c,S}^{\min}, y) - L_x(h_{c,S}^{\min}, y) \leq \max_{x', x \in S} L_{x'}(h_{a,S}^{\min}, y) - L_x(h_{a,S}^{\min}, y)$$
$$\tag{57}$$

Using Eqns. 45 and 56 we also obtain an auxiliary result.

$$\max_{x'} L_{x'}(h_{c,S}^{\min}, y) - \mathcal{L}_S(h_{c,S}^{\min}, y) \leq \max_{x'} L_{x'}(h_{a,S}^{\min}, y) - \mathcal{L}_S(h_{a,S}^{\min}, y)$$
$$\tag{58}$$

$\square$

# B  Sensitivity of Causal and Associational Models

Before we prove Lemma 1, we prove Corollary 1 and restate a Lemma from Wu et al. (2015) for completeness.

**Corollary 1.** *Let* $S$ *be a dataset of* $n$ $(x, y)$ *values, such that* $y^{(i)} \sim P(Y|X_C = x^{(i)}) \forall (x^{(i)}, y^{(i)}) \in S$, *where* $P(Y|X_C)$ *is the invariant conditional distribution on the causal features* $X_C$. *Consider a neighboring dataset* $S'$ *such that* $S' = S \backslash (x, y) + (x', y')$ *where* $(x, y) \in S$, $(x', y') \notin S$, *and* $(x', y')$ *shares the same conditional distribution* $y' \sim P(Y|X_C = x'_c)$. *Then the maximum generalization error from* $S$ *to* $S'$ *for a causal model trained on* $S$ *is lower than or equal to that of an associational model.*

$$\max_{S,S'} \mathcal{L}_{S'}(h_{c,S}^{\min}, y) - \mathcal{L}_{S}(h_{c,S}^{\min}, y) \leq \max_{S,S'} \mathcal{L}_{S'}(h_{a,S}^{\min}, y) - \mathcal{L}_{S}(h_{a,S}^{\min}, y)$$

*Proof.* Let $S_{n-1} = S \backslash (x, y))$ and similarly $S'_{n-1} = S' \backslash (x', y')$. Since $S$ and $S'$ differ in only one data point, $S_{n-1} = S'_{n-1}$. We will add and subtract sum of losses on data points in $S_{n-1}$, $(n-1)L_{S_{n-1}}$ to Theorem 2 statement.

Considering the LHS of Theorem 2,

$$\max_{x \in S, x'} L_{x'}(h_{c,S}^{\min}, y) - L_{x}(h_{c,S}^{\min}, y)$$
$$= \max_{x \in S, x'} L_{x'}(h_{c,S}^{\min}, y) + (n-1)\mathcal{L}_{S_{n-1}}(h_{c,S}^{\min}, y)$$
$$- L_{x}(h_{c,S}^{\min}, y) - (n-1)\mathcal{L}_{S_{n-1}}(h_{c,S}^{\min}, y)$$
$$= \max_{x \in S, x'} L_{x'}(h_{c,S}^{\min}, y) + (n-1)\mathcal{L}_{S'_{n-1}}(h_{c,S}^{\min}, y)$$
$$- L_{x}(h_{c,S}^{\min}, y) - (n-1)\mathcal{L}_{S_{n-1}}(h_{c,S}^{\min}, y)$$
$$= \max_{S'} n\mathcal{L}_{S'}(h_{c,S}^{\min}, y) - n\mathcal{L}_{S}(h_{c,S}^{\min}, y) \quad (59)$$

Similarly, the RHS of Theorem 2 can be written as,

$$\max_{x \in S, x'} L_{x'}(h_{a,S}^{\min}, y) - L_{x}(h_{a,S}^{\min}, y)$$
$$= \max_{x \in S, x'} L_{x'}(h_{a,S}^{\min}, y) + (n-1)\mathcal{L}_{S_{n-1}}(h_{a,S}^{\min}, y)$$
$$- L_{x}(h_{a,S}^{\min}, y) - (n-1)\mathcal{L}_{S_{n-1}}(h_{a,S}^{\min}, y)$$
$$= \max_{S'} n\mathcal{L}_{S'}(h_{a,S}^{\min}, y) - n\mathcal{L}_{S}(h_{a,S}^{\min}, y) \quad (60)$$

Using Theorem 2 and dividing Eqns. 59 and 60 by $n$, we obtain,

$$\max_{S'} n\mathcal{L}_{S'}(h_{c,S}^{\min}, y) - n\mathcal{L}_{S}(h_{c,S}^{\min}, y)$$
$$\leq \max_{S'} n\mathcal{L}_{S'}(h_{a,S}^{\min}, y) - n\mathcal{L}_{S}(h_{a,S}^{\min}, y)$$
$$\Rightarrow \max_{S'} \mathcal{L}_{S'}(h_{c,S}^{\min}, y) - \mathcal{L}_{S}(h_{c,S}^{\min}, y)$$
$$\leq \max_{S'} \mathcal{L}_{S'}(h_{a,S}^{\min}, y) - \mathcal{L}_{S}(h_{a,S}^{\min}, y) \quad (61)$$

Finally, since the above holds for any $S \sim P$, it will also hold for the worst-case $S$. The result follows.

$$\max_{S,S'} \mathcal{L}_{S'}(h_{c,S}^{\min}, y) - \mathcal{L}_{S}(h_{c,S}^{\min}, y)$$
$$\leq \max_{S'} \mathcal{L}_{S'}(h_{a,S}^{\min}, y) - \mathcal{L}_{S}(h_{a,S}^{\min}, y) \quad (62)$$

□

**Lemma 3.** *[From Wu et al. (2015)] Let* $S$ *and* $S'$ *be two neighboring datasets as defined in Corollary 1 where* $S' = S \backslash (x, y) + (x', y')$. *Given a model class* $\mathcal{H}$, *Let* $h_S^{\min}$ *be the loss-minimizing model on* $S$ *and* $h_{S'}^{\min}$ *be the loss-minimizing model on* $S'$. *Then the difference in losses between the two models on the same dataset is bounded by,*

$$\mathcal{L}_{S}(h_{S'}^{\min}, y) - \mathcal{L}_{S}(h_{S}^{\min}, y)$$
$$\leq \frac{L_{x'}(h_{S}^{\min}, y) - L_{x'}(h_{S'}^{\min}, y)}{n} + \frac{L_{x}(h_{S'}^{\min}, y) - L_{x}(h_{S}^{\min}, y)}{n} \quad (63)$$

*Proof.* The proof follows from expanding loss over a dataset into individual terms for each data point and then using the fact that $h_{S'}^{\min}$ has the minimum loss on $S'$.

Using the definition of $\mathcal{L}_{S} = \frac{1}{n} \sum_{i=1}^{n} L_{x_i}(h, y)$, we can write the following for any two neighboring datasets $S$ and $S'$.

$$\mathcal{L}_{S}(h_{S'}^{\min}, y) - \mathcal{L}_{S}(h_{S}^{\min}, y)$$
$$= \mathcal{L}_{S'}(h_{S'}^{\min}, y) + \frac{L_{x}(h_{S'}^{\min}, y) - L_{x'}(h_{S'}^{\min}, y)}{n}$$
$$- (\mathcal{L}_{S'}(h_{S}^{\min}, y) + \frac{L_{x}(h_{S}^{\min}, y) - L_{x'}(h_{S}^{\min}, y)}{n})$$
$$= (\mathcal{L}_{S'}(h_{S'}^{\min}, y) - \mathcal{L}_{S'}(h_{S}^{\min}, y)) + \frac{L_{x}(h_{S'}^{\min}, y) - L_{x'}(h_{S'}^{\min}, y)}{n}$$
$$+ \frac{L_{x'}(h_{S}^{\min}, y) - L_{x}(h_{S}^{\min}, y)}{n})$$
$$\leq \frac{L_{x'}(h_{S}^{\min}, y) - L_{x'}(h_{S'}^{\min}, y)}{n} + \frac{L_{x}(h_{S'}^{\min}, y) - L_{x}(h_{S}^{\min}, y)}{n} \quad (64)$$

where the last inequality is since $h_{S'}^{\min}$ is the minimizer of $L_{S'}(h, y)$ and thus $\mathcal{L}_{S'}(h_{S'}^{\min}, y) - \mathcal{L}_{S'}(h_{S}^{\min}, y) \leq 0$.

□

**Lemma 1.** *Let* $S$ *and* $S'$ *be two datasets defined as in Corollary 1. Let a model* $h$ *be specified by a set of parameters* $\theta \in \Omega \subseteq \mathbb{R}^d$. *Let* $h_S^{\min}(x; \theta_S)$ *be a model learnt using* $S$ *as training data and* $h_{S'}^{\min}(x; \theta_{S'})$ *be the model learnt using* $S'$ *as training data, using a loss function* $L$ *that is* $\lambda$-*strongly convex over* $\Omega$, $\rho$-*Lipschitz, symmetric and obeys the triangle inequality. Then, under the conditions of Theorem 1 (optimal predictor* $f \in \mathcal{H}_C$*) and for a sufficiently large* $n$, *the sensitivity of a causal learning function* $\mathcal{F}_c$ *that outputs learnt empirical model* $h_{c,S}^{\min} \leftarrow \mathcal{F}_c(S)$ *and* $h_{c,S'}^{\min} \leftarrow \mathcal{F}_c(S')$ *is lower than or equal to the sensitivity of an associational learning function* $\mathcal{F}_a$ *that outputs* $h_{a,S}^{\min} \leftarrow \mathcal{F}_a(S)$ *and* $h_{a,S'}^{\min} \leftarrow \mathcal{F}_a(S')$,

$$\Delta\mathcal{F}_c = \max_{S,S'} ||h_{c,S}^{\min} - h_{c,S'}^{\min}||_1 \leq \max_{S,S'} ||h_{a,S}^{\min} - h_{a,S'}^{\min}||_1 = \Delta\mathcal{F}_a$$

*where the maximum is over all such datasets* $S$ *and* $S'$.

*Proof.* Since L is a strongly convex function over $\Omega$, we can write for the two models $h_{c,S}^{\min}$ and $h_{c,S'}^{\min}$ trained on S and S' respectively (Wu et al., 2015),

$$
\begin{aligned}
\mathcal{L}_S(h_{c,S}^{\min}, y) &\leq \mathcal{L}_S(\alpha h_{c,S}^{\min} + (1-\alpha)h_{c,S'}^{\min}, y) \\
&\leq \alpha \mathcal{L}_S(h_{c,S}^{\min}, y) + (1-\alpha)\mathcal{L}_S(h_{c,S'}^{\min}, y) \\
&\quad - \frac{\lambda}{2}\alpha(1-\alpha)\|h_{c,S'}^{\min} - h_{c,S}^{\min}\|^2
\end{aligned}
\tag{65}
$$

where $\alpha \in (0,1)$ and the first inequality is since $h_{c,S}^{\min}$ is the loss-minimizing model over S. Rearranging the terms and tending $\alpha$ to 1 leads to,

$$
\begin{aligned}
(1-\alpha)&(\mathcal{L}_S(h_{c,S}^{\min}, y) - \mathcal{L}_S(h_{c,S'}^{\min}, y)) \\
&\leq -\frac{\lambda}{2}(1-\alpha)\left\|h_{c,S'}^{\min} - h_{c,S}^{\min}\right\|^2 \\
\Rightarrow \frac{\lambda}{2}&\left\|h_{c,S'}^{\min} - h_{c,S}^{\min}\right\|^2 \leq \mathcal{L}_S(h_{c,S'}^{\min}, y) - \mathcal{L}_S(h_{c,S}^{\min}, y)
\end{aligned}
\tag{66}
$$

Now consider $\max_{S,S'}\left\|h_{c,S}^{\min} - h_{c,S'}^{\min}\right\|_1$. Without loss of generality, we can order the pair of datasets S, S' such that $\mathcal{L}_S(h_{c,S'}^{\min}, y) \leq \mathcal{L}_{S'}(h_{c,S}^{\min}, y)$. Then using Eqn. 66 and taking the maximum, we obtain,

$$
\begin{aligned}
\frac{\lambda}{2}\max_{S,S'}\left\|h_{c,S}^{\min} - h_{c,S'}^{\min}\right\|_1^2 &\leq \max_{S,S'}\mathcal{L}_S(h_{c,S'}^{\min}, y) - \mathcal{L}_S(h_{c,S}^{\min}, y) \\
&\leq \max_{S,S'}\mathcal{L}_{S'}(h_{c,S}^{\min}, y) - \mathcal{L}_S(h_{c,S}^{\min}, y) \\
&\leq \max_{S,S'}\mathcal{L}_{S'}(h_{a,S}^{\min}, y) - \mathcal{L}_S(h_{a,S}^{\min}, y)
\end{aligned}
\tag{67}
$$

where the last inequality is due to Theorem 2. Let $S_1$ and $S_1'$ denote the datasets that lead to the maximum in the RHS above. We know that $\mathcal{L}_{S_1}(h_{a,S_1'}^{\min}) \geq \mathcal{L}_{S_1'}(h_{a,S_1'}^{\min})$ since $h_{a,S_1'}^{\min}$ is the loss-minimizing model over $S_1'$. Therefore, we can rewrite,

$$
\begin{aligned}
\max_{S,S'}&\mathcal{L}_{S'}(h_{a,S}^{\min}, y) - \mathcal{L}_S(h_{a,S}^{\min}, y) \\
&= \mathcal{L}_{S_1'}(h_{a,S_1}^{\min}, y) - \mathcal{L}_{S_1}(h_{a,S_1}^{\min}, y) \\
&\leq \mathcal{L}_{S_1'}(h_{a,S_1}^{\min}, y) - \mathcal{L}_{S_1}(h_{a,S_1}^{\min}, y) + (\mathcal{L}_{S_1}(h_{a,S_1'}^{\min}, y) - \mathcal{L}_{S_1'}(h_{a,S_1'}^{\min}, y)) \\
&= (\mathcal{L}_{S_1'}(h_{a,S_1}^{\min}, y) - \mathcal{L}_{S_1'}(h_{a,S_1'}^{\min}, y)) + (\mathcal{L}_{S_1}(h_{a,S_1'}^{\min}, y) - \mathcal{L}_{S_1}(h_{a,S_1}^{\min}, y))
\end{aligned}
\tag{68}
$$

Now using Lemma 3, we obtain the following two bounds.

$$
\begin{aligned}
\mathcal{L}_{S_1}&(h_{a,S_1'}^{\min}, y) - \mathcal{L}_{S_1}(h_{a,S_1}^{\min}, y) \\
&\leq \frac{L_{x'}(h_{a,S_1}^{\min}, y) - L_{x'}(h_{a,S_1'}^{\min}, y)}{n} + \frac{L_x(h_{a,S_1'}^{\min}, y) - L_x(h_{a,S_1}^{\min}, y)}{n} \\
\mathcal{L}_{S_1'}&(h_{a,S_1}^{\min}, y) - \mathcal{L}_{S_1'}(h_{a,S_1'}^{\min}, y) \\
&\leq \frac{L_{x'}(h_{a,S_1}^{\min}, y) - L_{x'}(h_{a,S_1'}^{\min}, y)}{n} + \frac{L_x(h_{a,S_1'}^{\min}, y) - L_x(h_{a,S_1}^{\min}, y)}{n}
\end{aligned}
\tag{69}
$$

Since the loss function $L(.,y)$ is $\rho$-Lipschitz, we have $L_x(h_1, y) - L_x(h_2, y) \leq \rho\|h_1 - h_2\|_1$ for any data point x and any two models $h_1$ and $h_2$. Plugging Eqn. 69 and the $\rho$-Lipschitz property back in Eqn. 68,

$$
\begin{aligned}
\mathcal{L}_{S_1'}&(h_{a,S_1}^{\min}, y) - \mathcal{L}_{S_1'}(h_{a,S_1'}^{\min}, y) \leq \frac{2}{n}\rho\left\|h_{a,S_1}^{\min} - h_{a,S_1'}^{\min}\right\|_1 \\
\mathcal{L}_{S_1}&(h_{a,S_1'}^{\min}, y) - \mathcal{L}_{S_1}(h_{a,S_1}^{\min}, y) \leq \frac{2}{n}\rho\left\|h_{a,S_1}^{\min} - h_{a,S_1'}^{\min}\right\|_1 \\
\Rightarrow \max_{S,S'}&\mathcal{L}_{S'}(h_{a,S}^{\min}, y) - \mathcal{L}_S(h_{a,S}^{\min}, y) \\
&\leq (\mathcal{L}_{S_1'}(h_{a,S_1}^{\min}, y) - \mathcal{L}_{S_1'}(h_{a,S_1'}^{\min}, y)) + (\mathcal{L}_{S_1}(h_{a,S_1'}^{\min}, y) - \mathcal{L}_{S_1}(h_{a,S_1}^{\min}, y)) \\
&\leq \frac{2}{n}\rho\left\|h_{a,S_1}^{\min} - h_{a,S_1'}^{\min}\right\|_1 + \frac{2}{n}\rho\left\|h_{a,S_1}^{\min} - h_{a,S_1'}^{\min}\right\|_1 \\
&= \frac{4}{n}\rho\left\|h_{a,S_1}^{\min} - h_{a,S_1'}^{\min}\right\|_1
\end{aligned}
\tag{70}
$$

Finally, combining with Eqn. 67, we obtain,

$$
\begin{aligned}
\max_{S,S'}\left\|h_{c,S'}^{\min} - h_{c,S}^{\min}\right\|_1^2 &\leq \frac{8\rho}{\lambda n}\left\|h_{a,S_1}^{\min} - h_{a,S_1'}^{\min}\right\|_1 \\
&\leq \frac{8\rho}{\lambda n}\max_{S,S'}\left\|h_{a,S}^{\min} - h_{a,S'}^{\min}\right\|_1 \\
&\leq \max_{S,S'}\left\|h_{a,S}^{\min} - h_{a,S'}^{\min}\right\|_1
\end{aligned}
\tag{71}
$$

where the last inequality holds for a sufficiently large n such that $\frac{8\rho}{\lambda n} \leq 1$. If $\max_{S,S'}\left\|h_{c,S'}^{\min} - h_{c,S}^{\min}\right\|_1 \geq 1$, the result follows. Otherwise, we need a larger n such that $n \geq \frac{8\rho}{\lambda \max_{S,S'}\left\|h_{c,S'}^{\min} - h_{c,S}^{\min}\right\|_1}$. In both cases, we obtain,

$$
\max_{S,S'}\left\|h_{c,S}^{\min} - h_{c,S'}^{\min}\right\|_1 \leq \max_{S,S'}\left\|h_{a,S}^{\min} - h_{a,S'}^{\min}\right\|_1
\tag{72}
$$

Thus sensitivity of a causal learning function is lower than that of an associational learning function, that is, $\Delta\mathcal{F}_c \leq \Delta\mathcal{F}_a$. $\qquad\square$

## C Differential Privacy Guarantees with Tighter Data-dependent Bounds

In this section we present the differential privacy guarantee of a causal model based on a recent method (Papernot et al., 2017) that provides tighter data-dependent bounds.

Before analyzing the differential privacy guarantee, we provide a generalization result for a 0-1 classifier based on the results from Theorem 1: causal classification models trained on data from two different distributions $P(X)$ and $P^*(X)$ are more likely to output the same value for a new input than associational models.

**Lemma 4.** *Under the conditions of Theorem 1 and 0-1 loss, let* $h_{c,S}^{\min}$ *be the loss-minimizing causal classification model trained on a dataset* S *from distribution* P *and let* $h_{c,S^*}^{\min}$ *be the loss-minimizing model trained on a dataset* $S^*$

*from* $P^*$. *Similarly, let* $h_{a,S}^{\min}$ *and* $h_{a,S^*}^{\min}$ *be loss-minimizing associational classification models trained on* $S$ *and* $S^*$ *respectively. Then for any new data input* $x$,

$$\min_{S\sim P,S^*\sim P^*} \Pr\left(h_{c,S}^{\min}(x) = h_{c,S^*}^{\min}(x)\right)$$
$$\geq \min_{S\sim P,S^*\sim P^*} \Pr\left(h_{a,S}^{\min}(x) = h_{a,S^*}^{\min}(x)\right)$$

*As the size of the training sample* $|S| = |S^*| \to \infty$, *the LHS*$\to 1$.

*Proof.* Let $h_{a,P}^{\min} = \arg\min_{h\in\mathcal{H}_A} \mathcal{L}_S(h,y)$ and $h_{a,P^*}^{\min} = \arg\min_{h\in\mathcal{H}_A} \mathcal{L}_{S^*}(h,y)$ be the loss-minimizing associational hypotheses under the two datasets $S$ and $S^*$ respectively, where $\mathcal{H}_A$ is the set of hypotheses. We can analogously define $h_{c,P}^{\min}$ and $h_{c,P^*}^{\min}$. Likewise, let $h_{a,P}^{OPT} = \arg\min_{h\in\mathcal{H}_A} \mathcal{L}_P(h,y)$ and similarly let $h_{a,P^*}^{OPT} = \arg\min_{h\in\mathcal{H}_A} \mathcal{L}_{P^*}(h,y)$ be the loss-minimizing hypotheses over the two distributions. We can analogously define $h_{c,P}^{OPT}$ and $h_{c,P^*}^{OPT}$. For ease of exposition, let us consider a binary classification task where all associational models map $X \to \{0,1\}$ and causal models map $X_C \to \{0,1\}$.

**Infinite sample result.** As $|S| = |S^*| \to \infty$, each of the models on $S$ and $S^*$ approach their loss-minimizing functions on the distributions $P$ and $P^*$ respectively. Then, for any input $x$,

$$\lim_{|S|\to\infty} h_{a,S}^{\min} = h_{a,P}^{OPT} \qquad \lim_{|S^*|\to\infty} h_{a,S^*}^{\min} = h_{a,P^*}^{OPT} \qquad (73)$$

$$\lim_{|S|\to\infty} h_{c,S}^{\min} = h_{c,P}^{OPT} \qquad \lim_{|S^*|\to\infty} h_{c,S^*}^{\min} = h_{c,P^*}^{OPT} \qquad (74)$$

From Theorem 1 (Equation 19), we know that $h_{c,P}^{OPT} = h_{c,P^*}^{OPT}$. Therefore, for any new input $x$ for a causal model, we obtain $\Pr\left(h_{c,P}^{OPT}(x) = h_{c,P^*}^{OPT}(x)\right) = 1$, but not necessarily for associational models. This leads to,

$$\lim_{|S|\to\infty,|S^*|\to\infty} \Pr\left(h_{c,S}^{\min}(x) = h_{c,S^*}^{\min}(x)\right) = 1 \qquad (75)$$

$$\geq \lim_{|S|\to\infty,|S^*|\to\infty} \Pr\left(h_{a,S}^{\min}(x) = h_{a,S^*}^{\min}(x)\right) \qquad (76)$$

**Finite sample result.** Under finite samples, let $S_1$ and $S_2^*$ be the two datasets from $P$ and $P^*$ respectively that lead to the minimum probability of agreement between the two causal models $h_{c,S_1}^{\min}$ and $h_{c,S_1^*}^{\min}$.

$$\min_{S\sim P,S^*\sim P^*} \Pr\left(h_{c,S}^{\min}(x) = h_{c,S^*}^{\min}(x)\right) = \Pr\left(h_{c,S_1}^{\min}(x) = h_{c,S_1^*}^{\min}(x)\right) \qquad (77)$$

Now consider the probability of agreement for the two associational models trained on the same datasets, $h_{a,S_1}^{\min}$ and $h_{a,S_1^*}^{\min}$. Without loss of generality, we can write the associational models as,

$$h_{a,S_1}^{\min}(x) = |h_{c,S_1}^{\min}(x) - h_{a,S_1}(x)|$$
$$h_{a,S_1^*}^{\min}(x) = |h_{c,S_1^*}^{\min}(x) - h_{a,S_1^*}(x)| \qquad (78)$$

where $h_{a,S_1} : X \to \{0,1\}$ and $h_{a,S_1^*} : X \to \{0,1\}$ are any two functions. Effectively, when $h_{a,S_1}$ is 1, it flips the output of the loss-minimizing associational model compared to the loss-minimizing causal model on $S_1$ (and similarly for $h_{a,S_1^*}$ on $S_1^*$).

Now we can select a different $S_2^* \sim P_2^*$ where $y$ and the causal features remain the same as $S_1^*$ but associational features are changed for each input $x \in S_2^*$. Therefore we construct $S_2^*$ such that $h_{c,S_1^*}^{\min} = h_{c,S_2^*}^{\min}$ but the term in the loss-minimizing associational model $h_{a,S_2^*}^{\min}$ has the following property: $h_{a,S_2^*} \neq h_{a,S_1}(x)$ if $h_{c,S_1^*}^{\min} = h_{c,S_1}^{\min}$, and $h_{a,S_2^*} = h_{a,S_1}(x)$ if $h_{c,S_1^*}^{\min} \neq h_{c,S_1}^{\min}$. Thus, under $S_2^*$, for all $x$,

$$\left|h_{a,S_1}^{\min}(x) - h_{a,S_2^*}^{\min}(x)\right| \geq \left|h_{c,S_1}^{\min}(x) - h_{c,S_2^*}^{\min}(x)\right|$$
$$= \left|h_{c,S_1}^{\min}(x) - h_{c,S_1^*}^{\min}(x)\right| = \max_{S,S^*}\left|h_{c,S}^{\min}(x) - h_{c,S^*}^{\min}(x)\right| \qquad (79)$$

Therefore, the disagreement between two associational models trained on two datasets is greater than or equal to the disagreement between causal models on the worst-case $S_1$ and $S_1^*$. Since the loss is 0-1 loss, the worst-case probability of agreement is lower.

$$\max_{S,S^*}\left|h_{c,S}^{\min}(x) - h_{c,S^*}^{\min}(x)\right| \leq \max_{S,S^*}\left|h_{a,S}^{\min}(x) - h_{a,S^*}^{\min}(x)\right|$$
$$\Rightarrow \min_{S\sim P,S^*\sim P^*} \Pr\left(h_{c,S}^{\min}(x) = h_{c,S^*}^{\min}(x)\right)$$
$$\geq \min_{S\sim P,S^*\sim P^*} \Pr\left(h_{a,S}^{\min}(x) = h_{a,S^*}^{\min}(x)\right)$$

$\square$

Based on the above generalization property, we now show that causal models provide stronger differential privacy guarantees than corresponding associational models. We utilize the *subsample and aggregate* technique (Dwork et al., 2014) that was extended for machine learning in Hamm et al. (2016) and Papernot et al. (2017), for constructing a differentiably private model. The framework considers M arbitrary teacher models that are trained on a separate subsample of the dataset without replacement. Then, a student model is trained on some auxiliary unlabeled data with the (pseudo) labels generated from a majority vote of the teachers. Differential privacy can be achieved by either perturbing the number of votes for each class (Papernot et al., 2017), or perturbing the learnt parameters of the student model (Hamm et al., 2016). For any new input, the output of the model is a majority vote on the predicted labels from the M models. The privacy guarantees are better if a larger number of teacher models agree on each input, since by definition the majority decision could not have been changed by modifying a single data point (or a single teacher's vote). Since causal models generalize to new distributions, intuitively we expect causal models trained on separate samples to agree more. Below

we show that for a fixed amount of noise, a causal model is $\epsilon_c$-DP compared to $\epsilon$-DP for a associational model, where $\epsilon_c \leq \epsilon$.

**Theorem 4.** *Let* D *be any dataset generated from possibly a mixture of different distributions* $\Pr(\mathtt{X},\mathtt{Y})$ *such that* $\Pr(\mathtt{Y}|\mathtt{X_C})$ *remains the same. Let* $\mathtt{n_k}$ *be the votes for the* k*th class from* M *teacher models. Let* $\mathcal{M}$ *be the mechanism that produces a noisy max,* $\arg\max_{\mathtt{k}}\{\mathtt{n_k} + \mathtt{Lap}(2/\gamma)\}$. *Then the privacy budget* $\epsilon_c$ *for a causal model trained on* D *is lower than that for an associational model.*

*Proof.* Consider a change in a single input example $(x, y)$, leading to a new $D'$ dataset. Since sub-datasets are sampled without replacement, only a single teacher model can change in $D'$. Let $n'_k$ be the vote counts for a class $k$ under $D'$. Because the change in a single input can only affect one model's vote, $|n_k - n'_k| \leq 1$.

Let the noise added to each class be $r_k \sim Lap(2/\gamma)$. Let the majority class (class with the highest votes) using data from $D$ be $i$ and the class with the second largest votes be $j$. Let us consider the minimum noise $r^*$ required for class $i$ to be the majority output under $\mathcal{M}$ over $D$. Then,

$$n_i + r^* > n_j + r_j$$

For $i$ to have the maximum votes using $\mathcal{M}$ over $D'$ too, we need,

$$n'_i + r_i > n'_j + r_j$$

In the worst case, $n'_i = n_i - 1$ and $n'_j = n_j + 1$ for some $j$. Thus, we need,

$$n_i - 1 + r_i > n_j + 1 + r_j \Rightarrow n_i + r_i > n_j + 2 + r_j$$

$$(80)$$

which shows that $r_i > r^* + 2$. Note that $r^* > r_j - (n_i - n_j)$. We have two cases:

**CASE I:** The noise $r_j < n_i - n_j$, and therefore $r^* < 0$. Writing $\Pr(i|D')$ to denote the probability that class $i$ is chosen as the majority class under $D'$,

$$P(\mathtt{i}|D') = P(\mathtt{r_i} \geq \mathtt{r^*} + 2) = 1 - 0.5\exp(\gamma)\exp\left(\frac{1}{2}\gamma\mathtt{r^*}\right)$$
$$= 1 - \exp(\gamma)(1 - P(\mathtt{r_i} \geq \mathtt{r^*}))$$
$$= 1 - \exp(\gamma)(1 - P(\mathtt{i}|D))$$

$$(81)$$

where the equations on the right are due to Laplace c.d.f. Using the above equation, we can write:

$$\frac{P(\mathtt{i}|D')}{P(\mathtt{i}|D)} = \exp(\gamma) + \frac{1 - \exp(\gamma)}{P(\mathtt{i}|D)}$$
$$= \exp(\gamma) + \frac{1 - \exp(\gamma)}{P(\mathtt{r_i} \geq \mathtt{r^*})} \leq \exp(\epsilon)$$

$$(82)$$

for some $\epsilon > 0$. As $P(i|D) = P(r_i \geq r^*)$ increases, the ratio decreases and thus the effective privacy budget ($\epsilon$) decreases. Thus, a DP-mechanism based on teacher models with lower $r^*$ (effectively higher $|r^*|$) will exhibit the lowest $\epsilon$.

Below we show that the worst-case $|r^*|$ across any two datasets $\mathtt{D_1} \sim \mathtt{P}$, $\mathtt{D_2} \sim \mathtt{P^*}$ such that $\mathtt{P}(\mathtt{Y}|\mathtt{X_C}) = \mathtt{P^*}(\mathtt{Y}|\mathtt{X_C})$ is higher for a causal model, and thus $\max_\mathtt{D} \mathtt{P}(\mathtt{r_i} \geq \mathtt{r^*})$ is higher. Intuitively, $|r^*|$ is higher when there is more consensus between the M teacher models since $|r^*|$ is the difference between the votes for the highest voted class with the votes for the second-highest class. For two sub-datasets $\mathtt{D_1} \subset \mathtt{D}$ and $\mathtt{D_2} \subset \mathtt{D}$, let the two causal teacher models be $\mathtt{h_{c,D_1}}$ and $\mathtt{h_{c,D_2}}$, and the two associational teacher models be $\mathtt{h_{a,D_1}}$ and $\mathtt{h_{a,D_2}}$. From Lemma 4, for any new x, there is more consensus among causal models. For any dataset D,

$$\min_{\mathtt{D_1},\mathtt{D_2}} \Pr(\mathtt{h_{c,D_1}}(\mathtt{x}) = \mathtt{h_{c,D_2}}(\mathtt{x})) \geq \min_{\mathtt{D_1},\mathtt{D_2}} \Pr(\mathtt{h_{a,D_1}}(\mathtt{x}) = \mathtt{h_{a,D_2}}(\mathtt{x}))$$

Hence worst-case $r_c^* \leq r^*$. From Equation 82, $\epsilon_c \leq \epsilon$.

**CASE II:** The noise $r_j >= n_i - n_j$, and therefore $r^* >= 0$. Following the steps above, we obtain:

$$P(\mathtt{i}|D') = P(\mathtt{r_i} \geq \mathtt{r^*} + 2) = 0.5\exp(-\gamma)\exp\left(-\frac{1}{2}\gamma\mathtt{r^*}\right)$$
$$= \exp(-\gamma)(P(\mathtt{r_i} \geq \mathtt{r^*}))$$
$$= \exp(-\gamma)(P(\mathtt{i}|D))$$

$$(83)$$

Thus, the ratio does not depend on $r^*$.

$$\frac{P(\mathtt{i}|D')}{P(\mathtt{i}|D)} = \exp(-\gamma)$$

$$(84)$$

Under CASE II when the noise is higher to the differences in votes between the highest and second highest voted class, causal models provide the same privacy budget as associational models.

Hence, overall, $\epsilon_c \leq \epsilon$. $\qquad \square$

# D Maximum Advantage of a Differentially Private algorithm

**Theorem 6.** *Under the conditions of Theorem 1, let* $\mathtt{S} \sim \mathtt{P}(\mathtt{X},\mathtt{Y})$ *be a dataset sampled from* P. *Let* $\hat{\mathcal{F}}_{\mathtt{c,S}}$ *and* $\hat{\mathcal{F}}_{\mathtt{a,S}}$ *be the differentially private mechanisms trained on* S *by adding identical Laplacian noise to the causal and associational learning functions from Lemma 1 respectively. Assume that a membership inference adversary is provided inputs sampled from either* P *or* P^*, *where* P^* *is any distribution such that* $\mathtt{P}(\mathtt{Y}|\mathtt{X_C}) = \mathtt{P^*}(\mathtt{Y}|\mathtt{X_C})$. *Then, across all adversaries* $\mathcal{A}$

*that predict membership in* $S \sim P$, *the worst-case membership advantage of* $\hat{\mathcal{F}}_{c,S}$ *is not greater than that of* $\hat{\mathcal{F}}_{a,S}$.

$$\max_{\mathcal{A},P^*} \mathrm{Adv}(\mathcal{A}, \hat{\mathcal{F}}_{c,S}, n, P, P^*) \leq \max_{\mathcal{A},P^*} \mathrm{Adv}(\mathcal{A}, \hat{\mathcal{F}}_{a,S}, n, P, P^*)$$

*Proof.* Consider a neighboring dataset $S'$ to $S \sim P$ such that $S'$ replaces data point $x \in S$ with a different point $x'$. Following Theorem 1 proof from Yeom et al. (2018), the membership advantage of an adversary $\mathcal{A}$ on a differentially private algorithm $\hat{\mathcal{F}}$ can be written as:

$$\mathrm{Adv}(\mathcal{A}, \hat{\mathcal{F}}, n, P, P^*) = \Pr(\mathcal{A} = 1 | b = 1) - \Pr(\mathcal{A} = 1 | b = 0)$$
$$= \Pr(\mathcal{A}(x, \hat{\mathcal{F}}_S) = 1 | x \in S) - \Pr(\mathcal{A}(x, \hat{\mathcal{F}}_{S'}) = 1 | x \in S)$$
(85)

where $\mathcal{A}(x, \hat{\mathcal{F}}_S)$ denotes a membership adversary for an algorithm $\hat{\mathcal{F}}_S$ trained on a dataset $S$, and $\mathcal{A}(x, \hat{\mathcal{F}}_{S'})$ denotes an adversary for algorithm $\hat{\mathcal{F}}_{S'}$ trained on $S'$. Without loss of generality for the case where there are an infinite number of models $h$, assume that models are sampled from a discrete set of $K$ models: $\{h_1, h_2, ..., h_K\}$. Then using the law of total probability over the models yielded by the algorithms $\hat{\mathcal{F}}_S$ and $\hat{\mathcal{F}}_{S'}$,

$$\mathrm{Adv}(\mathcal{A}, \hat{\mathcal{F}}, n, P, P^*) = \sum_{j=1}^{K} \Pr(\mathcal{A}(x, h_j) = 1) \Pr(\hat{\mathcal{F}}_S = h_j)$$
$$- \sum_{j=1}^{K} \Pr(\mathcal{A}(x, h_j) = 1) \Pr(\hat{\mathcal{F}}_{S'} = h_j)$$
$$= \sum_{j=1}^{K} \Pr(\mathcal{A}(x, h_j) = 1)[\Pr(\hat{\mathcal{F}}_S = h_j) - \Pr(\hat{\mathcal{F}}_{S'} = h_j)]$$
(86)

where $\Pr(\mathcal{A}(x, h_j) = 1)$ can be interpreted as the non-negative weights in a sum. Thus, the above is a weighted sum and will be maximum when positive values for $\Pr(\hat{\mathcal{F}}_S = h_j) - \Pr(\hat{\mathcal{F}}_{S'} = h_j)$ have the highest weight and negative values for $\Pr(\hat{\mathcal{F}}_S = h_j) - \Pr(\hat{\mathcal{F}}_{S'} = h_j)$ have zero weight. It follows that to obtain the maximum advantage, the adversary will choose $\Pr(\mathcal{A}(x, h_j) = 1) = 1$ if $\Pr(\hat{\mathcal{F}}_S = h_j) - \Pr(\hat{\mathcal{F}}_{S'} = h_j) > 0$, and $0$ otherwise. In other words, the adversary predicts membership in train set for an input $x \in S$ whenever probability of the given model $h_j$ being generated from $\hat{\mathcal{F}}_S$ is higher than it being generated from $\hat{\mathcal{F}}_{S'}$.

Let $H_+ \subset H$ be the set of models for which $\Pr(\hat{\mathcal{F}}_S = h_j) - \Pr(\hat{\mathcal{F}}_{S'} = h_j) > 0$. Similarly, let $H_- = H \setminus H_+$ be the set of models that are more probable to be generated from $\hat{\mathcal{F}}_{S'}$: $\Pr(\hat{\mathcal{F}}_{S'} = h_j) - \Pr(\hat{\mathcal{F}}_S = h_j) \geq 0$. The worst-case adversary selects datasets $S \sim P, S'$ such that the sum $\sum_{h_j \in H_+} \Pr(\hat{\mathcal{F}}_S = h_j) - \Pr(\hat{\mathcal{F}}_{S'} = h_j)$ is the highest. Therefore, for a given distribution $P$ and a differentially

private algorithm $\hat{\mathcal{F}}_S$ learnt on $S \sim P$, we can write the maximum membership advantage as,

$$\max_{\mathcal{A},P^*} \mathrm{Adv}(\mathcal{A}, \hat{\mathcal{F}}_S, n, P, P^*)$$
$$= \max_{S,S'} \sum_{h_j \in H_+} [\mathrm{P}(\hat{\mathcal{F}}_S = h_j) - \mathrm{P}(\hat{\mathcal{F}}'_S = h_j)]$$
$$= \max_{S,S'} \mathrm{P}(\hat{\mathcal{F}}_S \in H_+) - \mathrm{P}(\hat{\mathcal{F}}'_S \in H_+)$$
$$= \max_{S,S'} \mathrm{P}(\hat{\mathcal{F}}_S \in H_+) - (1 - \mathrm{P}(\hat{\mathcal{F}}'_S \in H_-))$$
$$= \max_{S,S'} 2 \Pr(\hat{\mathcal{F}}_S \in H_+) - 1$$
(87)

where the last equality is since $\hat{\mathcal{F}}_S$ and $\hat{\mathcal{F}}_{S'}$ have Laplace noise added from identical distributions and thus $\Pr(\hat{\mathcal{F}}_S \in H_+) = \Pr(\hat{\mathcal{F}}_{S'} \in H_-)$. Equation 87 provides the maximum membership advantage for any $\epsilon$-DP mechanism $\hat{\mathcal{F}}_S$ with Laplace noise.

We next show that Eqn. 87 for a causal differentially private mechanism $\hat{\mathcal{F}}_c$ is not greater than that for an associational mechanism $\hat{\mathcal{F}}_a$. Let $\Pr(\hat{\mathcal{F}}_S)$ be a Laplace distribution with mean at $h_{S,\min}$ and $\Pr(\hat{\mathcal{F}}_{S'})$ be a Laplace distribution with mean $h_{S',\min}$ with identical scale/noise parameters. We would like to find the boundary model $h^\dagger$ of the set $H_+$ where $\mathrm{P}(\hat{\mathcal{F}}_S = h_j) = \mathrm{P}(\hat{\mathcal{F}}'_S = h_j)$, since $\Pr(\hat{\mathcal{F}}_S \in H_+)$ is the probability under the Laplace distribution, cut off at the point $h^\dagger$. Due to identical noise for $\hat{\mathcal{F}}_S$ and $\hat{\mathcal{F}}'_S$ and the symmetry of the Laplace distribution, the boundary $h^\dagger$ corresponds to the midpoint of $h_{S,\min}$ and $h_{S',\min}$: $0.5(h_{S,\min} + h_{S',\min})$. Alternatively, the $\ell_1$ distance of the boundary $h^\dagger$ from the means of the Laplace distributions can be written as (for a worst case $S, S'$),

$$\left\| h^\dagger - h_{S,\min} \right\|_1 = \frac{\left\| h_{S',\min} - h_{S,\min} \right\|_1}{2} = \frac{\Delta \mathcal{F}_S}{2} \quad (88)$$

where $\Delta \mathcal{F}_S$ is the sensitivity of the learning function $\mathcal{F}_S$ and the last equality is due to the choice of worst-case $S$ and $S'$.

From Lemma 1, we know that sensitivity of a causal learning function is lower than that of an associational learning function.

$$\Delta \mathcal{F}_{c,S} \leq \Delta \mathcal{F}_{a,S} \quad (89)$$

Thus, $\ell_1$ distance of $h_c^\dagger$ from the mean $h_{c,S}^{\min}$ is lower for a causal learning function, and consequently its probability $\Pr(\hat{\mathcal{F}}_S = h_c^\dagger)$ is higher. Now the set $H_+$ is a one-sided boundary on the values of $h$ and includes the mean of the Laplace distribution. Given symmetry of the Laplace distribution, probability of $\hat{\mathcal{F}}_S$ lying in $H_+$, $\Pr(\hat{\mathcal{F}}_S \in H_+)$ should be lower whenever the probability at the one-sided boundary is higher. Therefore, $\mathrm{P}(\hat{\mathcal{F}}_S \in H_+)$ is lower for a causal mechanism than the associational learning mechanism.

$$\Delta \mathcal{F}_{c,S} \leq \Delta \mathcal{F}_{a,S} \Rightarrow \Pr(\hat{\mathcal{F}}_{c,S} \in H_+) \leq \Pr(\hat{\mathcal{F}}_{a,S} \in H_+)$$
(90)

Finally, using the above equation in Eqn. 87 shows that the maximum membership advantage of a causal model is lower.

$$\max_{\mathcal{A}, P^*} \text{Adv}(\mathcal{A}, \hat{\mathcal{F}}_{c,S}, n, P, P^*) \leq \max_{\mathcal{A}, P^*} \text{Adv}(\mathcal{A}, \hat{\mathcal{F}}_{a,S}, n, P, P^*) \tag{91}$$

$\square$

## E  Infinite Sample Robustness to MI Attacks

**Corollary 2.** *Under the conditions of Theorem 1, let* $h_{c,S}^{min}$ *be a causal model trained using empirical risk minimization on a dataset* $S \sim P(X, Y)$ *with sample size* $n$. *As* $n \to \infty$, *membership advantage* $\text{Adv}(\mathcal{A}, h_{c,S}^{min}) \to 0$.

*Proof.* $h_{c,S}^{min}$ can be obtained by empirical risk minimization.

$$h_{c,S}^{min} = \arg\min_{h \in \mathcal{H}_c} \mathcal{L}_{S \sim P}(h, y) = \arg\min_{h \in \mathcal{H}_c} \frac{1}{n} \sum_{i=1}^{n} L_{x_i}(h, y) \tag{92}$$

As $|S| = n \to \infty$, $h_{c,S}^{min} \to h_{c,P}^{OPT}$. Suppose now that there exists another $S'$ of the same size such that $S' \sim P^*$. Then as $|S'| \to \infty$, $h_{c,S'}^{min} \to h_{c,P^*}^{OPT}$.

From Theorem 1, $h_{c,P}^{OPT} = h_{c,P^*}^{OPT}$. Thus,

$$\lim_{n \to \infty} h_{c,S}^{min} = \lim_{n \to \infty} h_{c,S'}^{min} \tag{93}$$

Equation 93 implies that as $n \to \infty$, the learnt $h_{c,S}^{min}$ does not depend on the training set, as long as the training set is sampled from any distribution $P^*$ such that $P(Y|X_c) = P^*(Y|X_c)$. That is, being the global minimizer over distributions, $h_{c,S}^{min} = h_{c,P}^{OPT}$ does not depend on its training set. Therefore, $h_{c,S}^{min}(x)$ is independent of whether $x$ is in the training set.

$$\lim_{n \to \infty} \text{Adv}(\mathcal{A}, h_{c,S}^{min}) = \Pr(\mathcal{A} = 1|b = 1) - \Pr(\mathcal{A} = 1|b = 0)$$
$$= \mathbb{E}[\mathcal{A}|b = 1] - \mathbb{E}[\mathcal{A}|b = 0]$$
$$= \mathbb{E}[\mathcal{A}(h_{c,S}^{min})|b = 1] - \mathbb{E}[\mathcal{A}(h_{c,S}^{min})|b = 0]$$
$$= \mathbb{E}[\mathcal{A}(h_{c,S}^{min})] - \mathbb{E}[\mathcal{A}(h_{c,S}^{min})] = 0$$

$$\tag{94}$$

where the second last equality follows since any function of $h_{c,S}^{min}$ is independent of the training dataset. $\square$

## F  Robustness to Attribute Inference Attacks

**Theorem 7.** *Given a dataset* $S(X, Y)$ *of size* $n$ *and a structural causal model that connects* $X$ *to* $Y$, *a causal model* $h_c$ *makes it impossible to infer non-causal features.*

*Proof.* The proof follows trivially from the definition of a causal model. $h_c$ includes only causal features during training. Thus, $h(x)$ is independent of all features not in $X_c$.

$$\text{Adv}(\mathcal{A}, h) = \Pr(\mathcal{A} = 1|x_s = 1) - \Pr(\mathcal{A} = 1|x_s = 0)$$
$$= \Pr(\mathcal{A}(h) = 1|x_s = 1) - \Pr(\mathcal{A}(h) = 1|x_s = 0)$$
$$= \Pr(\mathcal{A}(h) = 1) - \Pr(\mathcal{A}(h) = 1) = 0$$

$\square$

## G  Generating Train and Test Datasets

This section describes how the train and test data are generated for the target and attacker models. The target model is built using a dataset generated using the bnlearn library. We divide the total dataset into training and test datasets in a $60 : 40$ ratio.

To generate train and test sets for the attacker model, the output of the target model for each of the training and test dataset is again divided into 50:50 ratio. Note that the attacker model is trained on the confidence outputs of the target model. Therefore, the training set for the attacker model consists of output confidence values from the target model's training as well as the test dataset. The dataset creation process is summarized in Figure 5.
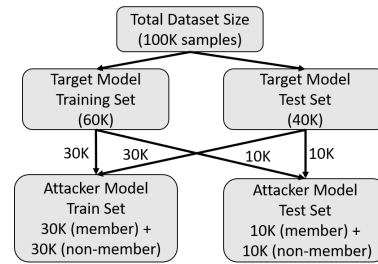


Figure 5: Dataset division for training target and attacker models.