# On the Importance of Modeling Unmeasured Confounders in Predictive Inference

Timothy Brathwaite

March 21st, 2020

## 1 Introduction

In a general prediction problem, the quantity we want to estimate is $P\left(Y^{*} \mid X^{*}, X, Y\right)$. Algebraically, this can be written using the rules of probability theory as follows.

$$
\begin{aligned}
P\left(Y^{*} \mid X^{*}, X, Y\right) &= \int P\left(Y^{*}, \lambda \mid X^{*}, X, Y\right) \partial \lambda \\
&= \int P\left(Y^{*} \mid X^{*}, X, Y, \lambda\right) P\left(\lambda \mid X^{*}, X, Y\right) \partial \lambda \\
\text{where } Y^{*} &= \text{The outcome value for a new observation.} \\
X^{*} &= \text{The explanatory variables for a new observation.} \\
Y &= \text{The historical outcome values.} \\
X &= \text{The historical explanatory values.} \\
\lambda &= \text{Everything we don't know that affects } X \text{ and/or } Y.
\end{aligned}
\tag{1}
$$

At this point, our expression contains two terms. The first term is a model for $Y^{*}$, $P\left(Y^{*} \mid X^{*}, X, Y, \lambda\right)$. The second is a probability for unknown variables $\lambda$, $P\left(\lambda \mid X^{*}, X, Y\right)$. As written, this expression is not useful because we don't know the form of $P\left(\lambda \mid X^{*}, X, Y\right)$. To re-express $P\left(\lambda \mid X^{*}, X, Y\right)$ in terms of quantities we know or can estimate, it is helpful to further specify what $\lambda$ is. In particular, we can decompose $\lambda$ as follows. Either:

1. $\lambda = (\tau, \theta)$ or

2. $\lambda = (\tau, \theta, \Omega)$

In these scenarios, $\tau$ represents unknown variables that only affect $X$, $\theta$ represents unknown variables that affect $Y$, and $\Omega$ represents variables that affect both $X$ and $Y$.

Likewise, it is useful to specify the type of relationships that can exist between $X$ and $Y$. We have four types of relationships:

1. $X$ causes $Y$

2. $Y$ causes $X$

3. $X$ does not cause $Y$ AND $Y$ does not cause $X$

4. $X$ causes $Y$ AND $Y$ causes $X$.

For the moment, we will ignore case four, because it is a mathematically hard case to deal with and because it is uncommon in most practical cases that I deal with. In most problems I face, $Y$ is an outcome of interest that occurs at time $t$ and $X$ is produced at time $t'$ where $t' < t$, so it is not sensible that $Y$ works backwards in time to cause $X$.

So, combining the three remaining types of relationships between $X$ and $Y$ with the two types of $\lambda$, we have the following 6 graphical models.

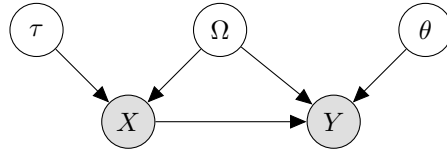| Case | $\lambda = (\tau, \theta)$ | $\lambda = (\tau, \theta, \Omega)$ |
|---|---|---|
| (i) | $\tau \rightarrow X \rightarrow Y \leftarrow \theta$ | $\tau \rightarrow X$, $\Omega \rightarrow X$, $\Omega \rightarrow Y$, $\theta \rightarrow Y$, $X \rightarrow Y$ |
| (ii) | $\tau \rightarrow X \leftarrow Y \leftarrow \theta$ | $\tau \rightarrow X$, $\Omega \rightarrow X$, $\Omega \rightarrow Y$, $\theta \rightarrow Y$, $Y \rightarrow X$ |
| (iii) | $\tau \rightarrow X$ $\quad$ $Y \leftarrow \theta$ | $\tau \rightarrow X$, $\Omega \rightarrow X$, $\Omega \rightarrow Y$, $\theta \rightarrow Y$ |

As will be shown in the Section 3, cases where $\lambda = (\tau, \theta)$ correspond to "typical" supervised learning models. Specifically, cases with $\lambda = (\tau, \theta)$ and:

- $X \rightarrow Y$ correspond to discriminative models (e.g. logistic regression).

- $X \leftarrow Y$ correspond to generative models (e.g. naive bayes).

- $X \quad Y$ (e.g. case iii) correspond to "constant" models, where no features are useful for predicting $Y$.

Scenarios where unobserved confounding is present, i.e. where $\lambda = (\tau, \theta, \Omega)$, are not so commonly discussed in machine learning texts.

## 2 The Current Scenario: A Discriminative Model with Unobserved Confounding

In particular, the case I want to focus on is depicted by the following graphical model.

In this scenario, there are unmeasured confounders $\Omega$ that cause both $X$ and $Y$, and we also believe that $X \rightarrow Y$. For example, if considering $X$ to be commute travel cost and $Y$ to be transportation mode (e.g. bike, walk, drive, transit), an example of an unobserved confounder might be travel distance if you don't know a user's exact home address. The travel distance will cause one's travel cost and may cause one's travel mode as well (e.g. if a person doesn't want to walk a long distance).

Let's see how we can calculate $P(Y^* \mid X^*, X, Y)$ in this scenario. Using Equation 1 and the factorization of $P(X^*, X, Y \mid \lambda)$ implied by the directed acyclic graph (DAG) above, we can continue writing

$$
\begin{aligned}
P(Y^* \mid X^*, X, Y) &= \int P(Y^* \mid X^*, X, Y, \lambda) P(\lambda \mid X^*, X, Y) \, \partial\lambda \\
&= \int P(Y^* \mid X^*, \lambda) P(\lambda \mid X^*, X, Y) \, \partial\lambda \\
&= \int P(Y^* \mid X^*, \Omega, \theta) P(\lambda \mid X^*, X, Y) \, \partial\lambda
\end{aligned}
\tag{2}
$$

Next, the term $P\left(\lambda \mid X^*, X, Y\right)$ can be calculated using Bayes' Rule and by (again) relying on the factorization from the DAG above.

$$
\begin{aligned}
P\left(\lambda \mid X^*, X, Y\right) &= \frac{P\left(X^*, X, Y \mid \lambda\right) P\left(\lambda\right)}{P\left(X^*, X, Y\right)} \\
&= \frac{P\left(X^* \mid \lambda\right) P\left(X, Y \mid \lambda\right) P\left(\lambda\right)}{P\left(X^*, X, Y\right)} \\
&= \frac{P\left(X^* \mid \lambda\right) P\left(Y \mid X, \lambda\right) P\left(X \mid \lambda\right) P\left(\lambda\right)}{P\left(X^*, X, Y\right)} \\
&= \frac{P\left(X^* \mid \tau, \Omega\right) P\left(Y \mid X, \theta, \Omega\right) P\left(X \mid \tau, \Omega\right) P\left(\lambda\right)}{P\left(X^*, X, Y\right)}
\end{aligned}
\tag{3}
$$

Substituting Equation 3 into Equation 2, we get the following expression.

$$
\begin{aligned}
P\left(Y^* \mid X^*, X, Y\right) &= \int \left[ \frac{P\left(Y^* \mid X^*, \theta, \Omega\right) P\left(X^* \mid \tau, \Omega\right) P\left(Y \mid X, \theta, \Omega\right) P\left(X \mid \tau, \Omega\right) P\left(\lambda\right)}{P\left(X^*, X, Y\right)} \right] \partial \lambda \\
&= \frac{\int P\left(Y^* \mid X^*, \theta, \Omega\right) P\left(X^* \mid \tau, \Omega\right) P\left(Y \mid X, \theta, \Omega\right) P\left(X \mid \tau, \Omega\right) P\left(\lambda\right) \partial \lambda}{P\left(X^*, X, Y\right)} \\
&= \frac{\int \left[ \frac{P\left(Y^* \mid X^*, \theta, \Omega\right) P\left(X^* \mid \tau, \Omega\right) P\left(Y \mid X, \theta, \Omega\right) P\left(X \mid \tau, \Omega\right) P\left(\lambda\right)}{P\left(X, Y\right)} \right] \partial \lambda}{\frac{P\left(X^*, X, Y\right)}{P\left(X, Y\right)}} \\
&= \frac{\int P\left(Y^* \mid X^*, \theta, \Omega\right) P\left(X^* \mid \tau, \Omega\right) P\left(\lambda \mid X, Y\right) \partial \lambda}{\frac{P\left(X^*, X, Y\right)}{P\left(X, Y\right)}} \\
&= \frac{\int P\left(Y^* \mid X^*, \theta, \Omega\right) P\left(X^* \mid \tau, \Omega\right) P\left(\lambda \mid X, Y\right) \partial \lambda}{P\left(X^* \mid X, Y\right)} \\
&\propto \int P\left(Y^* \mid X^*, \theta, \Omega\right) P\left(X^* \mid \tau, \Omega\right) P\left(\lambda \mid X, Y\right) \partial \lambda \\
&\approx \frac{1}{M} \sum_{i=1}^{M} P\left(Y^* \mid X^*, \theta_i, \Omega_i\right) P\left(X^* \mid \tau_i, \Omega_i\right)
\end{aligned}
\tag{4}
$$

$$
\begin{aligned}
\text{where } M &= \text{The number of samples from } P\left(\lambda \mid X, Y\right) \\
\lambda_i = \left(\tau_i, \theta_i, \Omega_i\right) &= \text{A singe sample from } P\left(\lambda \mid X, Y\right) \\
P\left(\lambda \mid X, Y\right) &= \frac{P\left(Y \mid X, \theta, \Omega\right) P\left(X \mid \tau, \Omega\right) P\left(\lambda\right)}{P\left(X, Y\right)}
\end{aligned}
$$

The reasoning above suggests that at prediction time, we should use a Monte Carlo estimate of the joint probability of $\left(X^*, Y^*\right)$. We can then normalize across $Y^*$ to get a probability in $[0, 1]$.
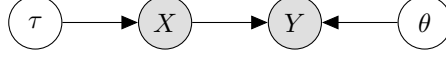
During training, we should estimate the posterior distribution $P\left(\lambda \mid X, Y\right)$. Given that $\Omega$ is expected to have at least one component per individual and we will only have one or at most a few observations per person in our dataset, we will likely not be in a world where the posterior distribution has concentrated on a single point. Accordingly, it would be best to estimate the entire posterior distribution since there may be wide uncertainty in our posterior beliefs about the value of the unmeasured confounders for some of the observations. However, if one persists in relying on point estimates, one can estimate the maximum a-posterior (MAP) point, which selects $\lambda$ to maximize the numerator: $P\left(Y \mid X, \theta, \Omega\right) P\left(X \mid \tau, \Omega\right) P\left(\lambda\right)$.

# 3 Deriving Known Models

At the end of Section 1, I stated that the situation where $\lambda = \left(\tau, \theta\right)$ corresponded to most typically used models such as discriminative and generative models. Below I'll provide a proof of these two points.

## 3.1 Discriminative Models

Discriminative models can be derived from the following graphical model where $\lambda = (\tau, \theta)$:



The derivation begins with Equation 1 and, using the DAG above, proceeds as follows.

$$
\begin{aligned}
P\left(Y^* \mid X^*, X, Y\right) &= \int P\left(Y^* \mid X^*, X, Y, \lambda\right) P\left(\lambda \mid X^*, X, Y\right) \partial \lambda \\
&= \int P\left(Y^* \mid X^*, \lambda\right) P\left(\lambda \mid X^*, X, Y\right) \partial \lambda \\
&= \int P\left(Y^* \mid X^*, \lambda\right) \frac{P\left(X^*, X, Y \mid \lambda\right) P(\lambda)}{P\left(X^*, X, Y\right)} \partial \lambda \\
&= \int \int P\left(Y^* \mid X^*, \tau, \theta\right) \frac{P\left(X^*, X, Y \mid \tau, \theta\right) P(\tau, \theta)}{P\left(X^*, X, Y\right)} \partial \tau \partial \theta \\
&= \int \int P\left(Y^* \mid X^*, \theta\right) \frac{P\left(X^* \mid \tau\right) P(X, Y \mid \tau, \theta) P(\tau, \theta)}{P\left(X^*, X, Y\right)} \partial \tau \partial \theta \\
&= \int \int P\left(Y^* \mid X^*, \theta\right) \frac{P\left(X^* \mid \tau\right) P(Y \mid X, \theta) P(X \mid \tau) P(\tau) P(\theta)}{P\left(X^*, X, Y\right)} \partial \tau \partial \theta \\
&= \int P\left(Y^* \mid X^*, \theta\right) P(Y \mid X, \theta) P(\theta) \partial \theta \int \frac{P\left(X^* \mid \tau\right) P(X \mid \tau) P(\tau)}{P\left(X^*, X, Y\right)} \partial \tau \\
&= \frac{1}{P\left(X^*, X, Y\right)} \int P\left(Y^* \mid X^*, \theta\right) P(Y \mid X, \theta) P(\theta) \partial \theta \int P\left(X^*, X, \tau\right) \partial \tau \\
&= \frac{P\left(X^*, X\right)}{P\left(X^*, X, Y\right)} \int P\left(Y^* \mid X^*, \theta\right) P(Y \mid X, \theta) P(\theta) \partial \theta \\
&= \frac{P\left(X^*, X\right)}{P\left(X^*, X, Y\right)} \int P\left(Y^* \mid X^*, \theta\right) \frac{P(Y \mid X, \theta) P(\theta)}{P(X, Y)} \partial \theta \\
&= \frac{\dfrac{P\left(X^*, X\right)}{P\left(X^*, X, Y\right)}}{P(X, Y)} \int P\left(Y^* \mid X^*, \theta\right) \frac{P(Y \mid X, \theta) P(\theta)}{P(X, Y)} \partial \theta \\
&= \frac{\dfrac{P\left(X^*, X\right)}{P(X)}}{\dfrac{P\left(X^*, X, Y\right)}{P(X, Y)}} \int P\left(Y^* \mid X^*, \theta\right) \frac{P(Y \mid X, \theta) P(X) P(\theta)}{P(X, Y)} \partial \theta \\
&= \frac{\dfrac{P\left(X^*, X\right)}{P(X)}}{\dfrac{P\left(X^*, X, Y\right)}{P(X, Y)}} \int P\left(Y^* \mid X^*, \theta\right) \frac{P(Y \mid X, \theta) P(X \mid \theta) P(\theta)}{P(X, Y)} \partial \theta \\
&= \frac{P\left(X^* \mid X\right)}{P\left(X^* \mid X, Y\right)} \int P\left(Y^* \mid X^*, \theta\right) \frac{P(Y, X, \theta)}{P(X, Y)} \partial \theta \\
&= \frac{P\left(X^* \mid X\right)}{P\left(X^* \mid X, Y\right)} \int P\left(Y^* \mid X^*, \theta\right) P(\theta \mid X, Y) \partial \theta \\
&= \int P\left(Y^* \mid X^*, \theta\right) P(\theta \mid X, Y) \partial \theta
\end{aligned}
\tag{5}
$$

The last line in Equation 5 follows from the following equality:

$$
\begin{aligned}
\frac{P\left(X^{*} \mid X\right)}{P\left(X^{*} \mid X, Y\right)} &= \frac{\int P\left(X^{*}, \tau \mid X\right) \partial \tau}{\int P\left(X^{*}, \tau \mid X, Y\right) \partial \tau} \\
&= \frac{\int P\left(X^{*} \mid \tau, X\right) P\left(\tau \mid X\right) \partial \tau}{\int P\left(X^{*} \mid \tau, X, Y\right) P\left(\tau \mid X, Y\right) \partial \tau} \\
&= \frac{\int P\left(X^{*} \mid \tau\right) P\left(\tau \mid X\right) \partial \tau}{\int P\left(X^{*} \mid \tau\right) P\left(\tau \mid X\right) \partial \tau} \\
&= 1
\end{aligned}
\tag{6}
$$

Moreover, the last line in Equation 5 justifies the discriminative learning technique of (1) building a model for $Y$ as a function of $X$ and $\theta$, and (2) taking the expectation of that model's outputs with respect to the posterior distribution of $\theta$. By Equation 5, no other models are required. In particular, no models of $X$ are needed. Note that this contrasts the equivalent scenario with unobserved confounding. There, one needs both a model for $Y$ AND a model for $X$.

## 3.2 Generative Models

In this section, I'll show how generative models come from the following graphical model:



First, it is important to note that this graphical model admits the following factorization of the joint probability $P\left(X, Y, \tau, \theta\right)$: $\{P\left(X \mid \tau, Y\right), P\left(Y \mid \theta\right), P\left(\tau\right), P\left(\theta\right)\}$. As a result, all of our probabilistic expressions for these variables should be written in terms of these basic elements. This requires that we re-express the following probability from Equation 1 as follows:

$$
\begin{aligned}
P\left(Y^{*} \mid X^{*}, X, Y, \lambda\right) &= P\left(Y^{*} \mid X^{*}, \lambda\right) \\
&= \frac{P\left(X^{*}, Y^{*} \mid \lambda\right)}{P\left(X^{*} \mid \lambda\right)} \\
&= \frac{P\left(X^{*} \mid Y^{*}, \lambda\right) P\left(Y^{*} \mid \lambda\right)}{P\left(X^{*} \mid \lambda\right)} \\
&= \frac{P\left(X^{*} \mid Y^{*}, \tau\right) P\left(Y^{*} \mid \theta\right)}{P\left(X^{*} \mid \tau, \theta\right)}
\end{aligned}
\tag{7}
$$

Similarly, we will need the following re-expression:

$$
\begin{aligned}
P\left(X^{*}, X, Y \mid \lambda\right) &= P\left(X^{*} \mid X, Y, \lambda\right) P\left(X, Y \mid \lambda\right) \\
&= P\left(X^{*} \mid \lambda\right) P\left(X, Y \mid \lambda\right) \\
&= P\left(X^{*} \mid \tau, \theta\right) P\left(X, Y \mid \tau, \theta\right) \\
&= P\left(X^{*} \mid \tau, \theta\right) P\left(X \mid Y, \tau, \theta\right) P\left(Y \mid \tau, \theta\right) \\
&= P\left(X^{*} \mid \tau, \theta\right) P\left(X \mid Y, \tau\right) P\left(Y \mid \theta\right)
\end{aligned}
\tag{8}
$$

With these expressions, we can now return to Equation 1 and proceed in the following manner:

$$
\begin{aligned}
P\left(Y^{*} \mid X^{*}, X, Y\right) &= \int P\left(Y^{*} \mid X^{*}, X, Y, \lambda\right) P\left(\lambda \mid X^{*}, X, Y\right) \partial \lambda \\
&= \int \frac{P\left(X^{*} \mid Y^{*}, \tau\right) P\left(Y^{*} \mid \theta\right)}{P\left(X^{*} \mid \tau, \theta\right)} P\left(\lambda \mid X^{*}, X, Y\right) \partial \lambda \\
&= \int \frac{P\left(X^{*} \mid Y^{*}, \tau\right) P\left(Y^{*} \mid \theta\right)}{P\left(X^{*} \mid \tau, \theta\right)} \frac{P\left(X^{*}, X, Y \mid \lambda\right) P\left(\lambda\right)}{P\left(X^{*}, X, Y\right)} \partial \lambda \\
&= \int \frac{P\left(X^{*} \mid Y^{*}, \tau\right) P\left(Y^{*} \mid \theta\right)}{P\left(X^{*} \mid \tau, \theta\right)} \frac{P\left(X^{*} \mid \tau, \theta\right) P\left(X, Y \mid \tau, \theta\right) P\left(\lambda\right)}{P\left(X^{*}, X, Y\right)} \partial \lambda \\
&= \int P\left(X^{*} \mid Y^{*}, \tau\right) P\left(Y^{*} \mid \theta\right) \frac{P\left(X, Y \mid \tau, \theta\right) P\left(\lambda\right)}{P\left(X^{*}, X, Y\right)} \partial \lambda \\
&= \int P\left(X^{*} \mid Y^{*}, \tau\right) P\left(Y^{*} \mid \theta\right) \frac{\dfrac{P\left(X, Y \mid \tau, \theta\right) P\left(\lambda\right)}{P\left(X, Y\right)}}{\dfrac{P\left(X^{*}, X, Y\right)}{P\left(X, Y\right)}} \partial \lambda \\
&= \int P\left(X^{*} \mid Y^{*}, \tau\right) P\left(Y^{*} \mid \theta\right) \frac{P\left(\tau, \theta \mid X, Y\right)}{P\left(X^{*} \mid X, Y\right)} \partial \lambda \\
&\propto \int P\left(X^{*} \mid Y^{*}, \tau\right) P\left(Y^{*} \mid \theta\right) P\left(\tau, \theta \mid X, Y\right) \partial \lambda
\end{aligned}
$$

$$
\begin{aligned}
\text{where } P\left(\tau, \theta \mid X, Y\right) &= \frac{P\left(X, Y \mid \tau, \theta\right)}{P\left(X, Y\right)} \\
&= \frac{P\left(X \mid Y, \tau\right) P\left(Y \mid \theta\right) P\left(\lambda\right)}{P\left(X, Y\right)}
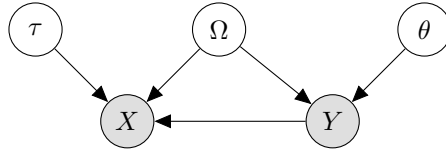\end{aligned}
$$

(9)

The last three lines of Equation 9 justify the typical generative model. First, models are built for $P\left(X \mid Y, \tau\right)$ and $P\left(Y \mid \theta\right)$. From these two models we can estimate the posterior distribution of $\lambda$: $P\left(\tau, \theta \mid X, Y\right)$. Secondly, given samples from $P\left(\tau, \theta \mid X, Y\right)$, we can use the two aforementioned models to construct a Monte Carlo estimate of $P\left(X^{*}, Y^{*} \mid X, Y\right)$. Finally, we can normalize across $Y^{*}$ to calculate $P\left(Y^{*} \mid X^{*}, X, Y\right)$.

Here, as in the case of a discriminative model with unobserved confounding, we need both a model of $X$ AND a model of $Y$.

# 4   A Generative Model with Unobserved Confounding

This case is conceptually very similar to the generative model that was just presented. As a result I will proceed by analogy and provide any missing derivations only as requested by others.

First, the graphical model that depicts this scenario is

Secondly, the desired predictive distribution is:

$$
\begin{aligned}
P\left(Y^* \mid X^*, X, Y\right) &= \int P\left(Y^* \mid X^*, X, Y, \lambda\right) P\left(\lambda \mid X^*, X, Y\right) \partial \lambda \\
&= \int \frac{P\left(X^* \mid Y^*, \tau, \Omega\right) P\left(Y^* \mid \theta, \Omega\right)}{P\left(X^* \mid \tau, \theta, \Omega\right)} P\left(\lambda \mid X^*, X, Y\right) \partial \lambda \\
&= \int \frac{P\left(X^* \mid Y^*, \tau, \Omega\right) P\left(Y^* \mid \theta, \Omega\right)}{P\left(X^* \mid \tau, \theta, \Omega\right)} \frac{P\left(X^*, X, Y \mid \lambda\right) P\left(\lambda\right)}{P\left(X^*, X, Y\right)} \partial \lambda \\
&= \int \frac{P\left(X^* \mid Y^*, \tau, \Omega\right) P\left(Y^* \mid \theta, \Omega\right)}{P\left(X^* \mid \tau, \theta, \Omega\right)} \frac{P\left(X^* \mid \tau, \theta, \Omega\right) P\left(X, Y \mid \tau, \theta, \Omega\right) P\left(\lambda\right)}{P\left(X^*, X, Y\right)} \partial \lambda \\
&= \int P\left(X^* \mid Y^*, \tau, \Omega\right) P\left(Y^* \mid \theta, \Omega\right) \frac{P\left(X, Y \mid \tau, \theta, \Omega\right) P\left(\lambda\right)}{P\left(X^*, X, Y\right)} \partial \lambda \\
&= \int P\left(X^* \mid Y^*, \tau, \Omega\right) P\left(Y^* \mid \theta, \Omega\right) \frac{\dfrac{P\left(X, Y \mid \tau, \theta, \Omega\right) P\left(\lambda\right)}{P\left(X, Y\right)}}{\dfrac{P\left(X^*, X, Y\right)}{P\left(X, Y\right)}} \partial \lambda & (10)\\
&= \int P\left(X^* \mid Y^*, \tau, \Omega\right) P\left(Y^* \mid \theta, \Omega\right) \frac{P\left(\tau, \theta, \Omega \mid X, Y\right)}{P\left(X^* \mid X, Y\right)} \partial \lambda \\
&\propto \int P\left(X^* \mid Y^*, \tau, \Omega\right) P\left(Y^* \mid \theta, \Omega\right) P\left(\tau, \theta, \Omega \mid X, Y\right) \partial \lambda
\end{aligned}
$$

$$
\begin{aligned}
\text{where } P\left(\tau, \theta, \Omega \mid X, Y\right) &= \frac{P\left(X, Y \mid \tau, \theta, \Omega\right)}{P\left(X, Y\right)} \\
&= \frac{P\left(X \mid Y, \tau, \Omega\right) P\left(Y \mid \theta, \Omega\right) P\left(\lambda\right)}{P\left(X, Y\right)}
\end{aligned}
$$

As shown by the last three lines of Equation 10 versus the last three lines of Equation 9, there is little mathematical difference between generative models with and without unobserved confounding. Moreover, the prediction equation for a generative model with confounding and a discriminative model with confounding are symmetric: switch $X$ and $Y$ and switch $\tau$ and $\theta$.