

Bounding Counterfactual Outcomes of Health Insurance Delay-and-Deny Practices (E-companion)

Martin B. Haugh and Raghav Singal

Appendix EC.1: Further Details on Sampling Hidden Paths

Lemma 1. *For $t > 1$, the forward probabilities obey the following recursion:*

$$\alpha(h_t) = e_{h_t x_t o_t} \sum_{h_{t-1} \in \mathbb{H}} q_{h_{t-1} o_{t-1} h_t} \times \alpha(h_{t-1}) \quad \forall h_t \in \mathbb{H}. \quad (4)$$

The recursion breaks at $t = 1$ with $\alpha(h_1) = p_{h_1} \times e_{h_1 x_1 o_1}$ for all $h_1 \in \mathbb{H}$.

Proof. We begin with $t = 1$ and note that for all $h_1 \in \mathbb{H}$, $\alpha(h_1) := \mathbb{P}(o_1 | h_1, x_1) \times \mathbb{P}(h_1 | x_1) \times \mathbb{P}(x_1) = \mathbb{P}(o_1 | h_1, x_1) \times \mathbb{P}(h_1) \times \mathbb{P}(x_1) = e_{h_1 x_1 o_1} \times p_{h_1}$, where we use the fact that the policy $x_{1:T}$ is deterministic. For $t > 1$, note that for all $h_t \in \mathbb{H}$,

$$\begin{aligned} \alpha(h_t) &= \sum_{h_{t-1}} \mathbb{P}(h_t, h_{t-1}, o_{1:t-1}, o_t, x_{1:t}) \\ &= \sum_{h_{t-1}} \mathbb{P}(o_t | h_t, h_{t-1}, o_{1:t-1}, x_{1:t}) \times \mathbb{P}(h_t | h_{t-1}, o_{1:t-1}, x_{1:t}) \\ &\quad \times \mathbb{P}(x_t | h_{t-1}, o_{1:t-1}, x_{1:t-1}) \times \mathbb{P}(h_{t-1}, o_{1:t-1}, x_{1:t-1}) \\ &= \sum_{h_{t-1}} \mathbb{P}(o_t | h_t, x_t) \times \mathbb{P}(h_t | h_{t-1}, o_{t-1}) \times \mathbb{P}(x_t) \times \mathbb{P}(h_{t-1}, o_{1:t-1}, x_{1:t-1}) \\ &= \mathbb{P}(x_t) \times \mathbb{P}(o_t | h_t, x_t) \sum_{h_{t-1}} \mathbb{P}(h_t | h_{t-1}, o_{t-1}) \alpha(h_{t-1}) = e_{h_t x_t o_t} \sum_{h_{t-1}} q_{h_{t-1} o_{t-1} h_t} \times \alpha(h_{t-1}), \end{aligned}$$

where we again use the fact that the policy $x_{1:T}$ is deterministic. This completes the proof. \square

Efficiently Computing the Forward Probabilities. See Algorithm 3. For numerical stability, we compute these probabilities on the log scale. In particular, the recursion in (4) (Lemma 1) can be expressed as follows (by taking log on both sides):

$$\log(\alpha(h_t)) = \log(e_{h_t x_t o_t}) + \log \sum_{h_{t-1}} \underbrace{q_{h_{t-1} o_{t-1} h_t} \times \alpha(h_{t-1})}_{=\exp\{\log(q_{h_{t-1} o_{t-1} h_t}) + \log(\alpha(h_{t-1}))\}},$$

where the “ $\log \sum \exp$ ” term can be evaluated using a standard logsumexp function.

Proposition 1. *Algorithm 1 outputs samples from the posterior distribution of $H_{1:T} | (o_{1:T}, x_{1:T})$.*

Proof. Observe that

$$\begin{aligned} \mathbb{P}(h_{1:T} | o_{1:T}, x_{1:T}) &= \mathbb{P}(h_1 | h_{2:T}, o_{1:T}, x_{1:T}) \times \mathbb{P}(h_{T-1} | h_T, o_{1:T}, x_{1:T}) \times \mathbb{P}(h_T | o_{1:T}, x_{1:T}) \\ &= \mathbb{P}(h_1 | h_2, o_{1:T}, x_{1:T}) \times \mathbb{P}(h_{T-1} | h_T, o_{1:T}, x_{1:T}) \times \mathbb{P}(h_T | o_{1:T}, x_{1:T}). \end{aligned}$$

We can therefore sample sequentially via the following two steps:

Algorithm 3 `compute_alpha(p, E, Q, o_{1:T}, x_{1:T})`

```

1:  $\log(\alpha(1, h)) = \log(p_h) + \log(e_{hx_1o_1})$  for all  $h \in \mathbb{H}$  % initialize period 1
2: for  $t = 2 : T$  do
3:   for  $h \in \mathbb{H}$  do
4:     t_minus_terms( $h'$ ) =  $\log(q_{h'o_{t-1}h}) + \log(\alpha(t-1, h'))$  for all  $h' \in \mathbb{H}$ 
5:      $\log(\alpha(t, h)) = \log(e_{hx_t o_t}) + \text{logsumexp}(\text{t\_minus\_terms})$ 
6:   end for
7: end for
8: return  $\exp(\log(\alpha))$ 

```

1. Draw h_T from $\mathbb{P}(h_T | o_{1:T}, x_{1:T})$.
2. For $t = T - 1, \dots, 1$, draw h_t from $\mathbb{P}(h_t | h_{t+1}, o_{1:T}, x_{1:T})$.

For Step 1, we can use $\alpha(h_T)$ since $\mathbb{P}(h_T | o_{1:T}, x_{1:T}) \propto \alpha(h_T)$. For Step 2, observe that for $t < T$,

$$\begin{aligned}
\mathbb{P}(h_t | h_{t+1}, o_{1:T}, x_{1:T}) &\propto \mathbb{P}(h_t, h_{t+1} | o_{1:T}, x_{1:T}) \\
&\propto \alpha(h_t) \times \mathbb{P}(h_{t+1}, x_{t+1} | h_t, o_t) && \text{[see (EC.1) below]} \\
&= \alpha(h_t) \times \mathbb{P}(h_{t+1} | h_t, o_t) \times \mathbb{P}(x_{t+1}) \\
&= \alpha(h_t) \times q_{h_t o_t h_{t+1}}.
\end{aligned}$$

For the second proportionality, note the *pairwise marginal* $\mathbb{P}(h_t, h_{t+1} | o_{1:T}, x_{1:T})$ is proportional to

$$\begin{aligned}
&\propto \mathbb{P}(o_{1:t}, o_{t+1}, o_{t+2:T}, x_{1:t}, x_{t+1}, x_{t+2:T}, h_{t+1}, h_t) \\
&= \mathbb{P}(o_{t+2:T}, x_{t+2:T} | o_{1:t}, o_{t+1}, x_{1:t}, x_{t+1}, h_{t+1}, h_t) \times \mathbb{P}(o_{1:t}, o_{t+1}, x_{1:t}, x_{t+1}, h_{t+1}, h_t) \\
&= \mathbb{P}(o_{t+2:T}, x_{t+2:T} | h_{t+1}, o_{t+1}) \times \mathbb{P}(o_{t+1} | o_{1:t}, h_{t+1}, h_t, x_{1:t}, x_{t+1}) \times \mathbb{P}(o_{1:T}, h_{t+1}, h_t, x_{1:t}, x_{t+1}) \\
&= \mathbb{P}(o_{t+2:T}, x_{t+2:T} | h_{t+1}, o_{t+1}) \times \mathbb{P}(o_{t+1} | h_{t+1}, x_{t+1}) \times \mathbb{P}(h_{t+1}, x_{t+1} | o_{1:T}, x_{1:t}, h_t) \times \mathbb{P}(o_{1:t}, x_{1:t}, h_t) \\
&= \mathbb{P}(o_{t+2:T}, x_{t+2:T} | h_{t+1}, o_{t+1}) \times \mathbb{P}(o_{t+1} | h_{t+1}, x_{t+1}) \times \mathbb{P}(h_{t+1}, x_{t+1} | h_t, o_t) \times \mathbb{P}(o_{1:t}, x_{1:t}, h_t) \\
&= \mathbb{P}(o_{t+2:T}, x_{t+2:T} | h_{t+1}, o_{t+1}) \times \mathbb{P}(o_{t+1} | h_{t+1}, x_{t+1}) \times \mathbb{P}(h_{t+1}, x_{t+1} | h_t, o_t) \times \alpha(h_t). \tag{EC.1}
\end{aligned}$$

This completes the proof. □

Appendix EC.2: Proofs of Lemmas 2 and 3

Lemma 2. *We have*

$$PN = 1 - \lim_{B \rightarrow \infty} \frac{1}{B} \sum_{b=1}^B \mathbb{P}_{\tilde{\mathbf{M}}(b)}(H_T^{\tilde{x}_{1:T}} = h^*). \tag{7}$$

Proof. Observe that

$$PN = 1 - \mathbb{P}_{\tilde{\mathbf{M}}}(H_T^{\tilde{x}_{1:T}} = h^*) \tag{by definition}$$

$$\begin{aligned}
&= 1 - \mathbb{E}_{\widetilde{\mathbf{M}}}[\mathbb{I}\{H_T^{\widetilde{x}_{1:T}} = h^*\}] && [\mathbb{P}(Y = y) = \mathbb{E}[\mathbb{I}\{Y = y\}]] \\
&= 1 - \mathbb{E}_{H_{1:T}}[\mathbb{E}_{\widetilde{\mathbf{M}}|H_{1:T}}[\mathbb{I}\{H_T^{\widetilde{x}_{1:T}} = h^*\}]] && [\text{law of total expectation}] \\
&= 1 - \lim_{B \rightarrow \infty} \frac{1}{B} \sum_{b=1}^B \mathbb{P}_{\widetilde{\mathbf{M}}(b)}(H_T^{\widetilde{x}_{1:T}} = h^*). && [\text{law of large numbers}]
\end{aligned}$$

The proof is now complete. \square

Lemma 3. $\mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{h}_t) := \mathbb{P}_{\widetilde{\mathbf{M}}(b)}(H_t^{\widetilde{x}_{1:T}} = \bar{h}_t)$ obeys the following recursion over $t \in \{T, \dots, 2\}$:

$$\mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{h}_t) = \sum_{\bar{h}_{t-1}, \bar{o}_{t-1}} \frac{\pi_{\bar{h}_{t-1}\bar{o}_{t-1}, h_{t-1}(b) o_{t-1}}(\bar{h}_t, h_t(b))}{q_{h_{t-1}(b) o_{t-1} h_t(b)}} \times \frac{\theta_{\bar{h}_{t-1}\bar{x}_{t-1}, h_{t-1}(b) x_{t-1}}(\bar{o}_{t-1}, o_{t-1})}{e_{h_{t-1}(b) x_{t-1} o_{t-1}}} \times \mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{h}_{t-1}).$$

The recursion breaks at $t = 1$ with $\mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{h}_1) = 1$ if $\bar{h}_1 = h_1(b)$ and 0 otherwise.

Proof. For $t \in \{T, T-1, \dots, 2\}$, observe that

$$\begin{aligned}
\mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{h}_t) &= \sum_{\bar{h}_{t-1} \in \mathbb{H}} \sum_{\bar{o}_{t-1} \in \mathbb{O}} \mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{h}_t, \bar{h}_{t-1}, \bar{o}_{t-1}) \\
&= \sum_{\bar{h}_{t-1} \in \mathbb{H}} \sum_{\bar{o}_{t-1} \in \mathbb{O}} \mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{h}_t | \bar{h}_{t-1}, \bar{o}_{t-1}) \mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{o}_{t-1} | \bar{h}_{t-1}) \mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{h}_{t-1}) \\
&= \sum_{\bar{h}_{t-1} \in \mathbb{H}} \sum_{\bar{o}_{t-1} \in \mathbb{O}} \tilde{q}_{t-1 \bar{h}_{t-1} \bar{o}_{t-1} \bar{h}_t}(b) \times \tilde{e}_{t-1 \bar{h}_{t-1} \bar{x}_{t-1} \bar{o}_{t-1}}(b) \times \mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{h}_{t-1}) \\
&= \sum_{\bar{h}_{t-1} \in \mathbb{H}} \sum_{\bar{o}_{t-1} \in \mathbb{O}} \frac{\pi_{\bar{h}_{t-1}\bar{o}_{t-1}, h_{t-1}(b) o_{t-1}}(\bar{h}_t, h_t(b))}{q_{h_{t-1}(b) o_{t-1} h_t(b)}} \times \frac{\theta_{\bar{h}_{t-1}\bar{x}_{t-1}, h_{t-1}(b) x_{t-1}}(\bar{o}_{t-1}, o_{t-1})}{e_{h_{t-1}(b) x_{t-1} o_{t-1}}} \times \mathbb{P}_{\widetilde{\mathbf{M}}(b)}(\bar{h}_{t-1}).
\end{aligned}$$

The base case ($t = 1$) holds since the counterfactual hidden state in period 1 equals the posterior sample $h_1(b)$ (as discussed in §4.2). The proof is now complete. \square

Appendix EC.3: The Structural Causal Model and Response Distributions

We define the SCM in §EC.3.1 and model it via response function distributions in §EC.3.2.

EC.3.1. The Structural Causal Model (SCM)

A generic SCM for the GHMM from Figure 1 is illustrated in Figure EC.1. Consider, for example, the observation O_t at any time t . It is modeled as a stochastic function of its parents (H_t, X_t) . To capture this stochasticity, we introduce an exogenous noise vector $\mathbf{V}_t := [V_{thx}]_{h,x}$, consisting of $|\mathbb{H}||\mathbb{X}| \mathbb{U}[0, 1]$ random variables. The vector representation of \mathbf{V}_t (rather than a scalar) allows us to account for the fact that each $O_{thx} := O_t | (H_t = h, X_t = x)$ defines a *distinct* random variable (or *potential outcome*) for all (h, x) . These random variables may be independent or exhibit positive or negative dependence. To model these scenarios, we associate each O_{thx} with a corresponding

noise variable V_{thx} . The *dependence structure* among $[V_{thx}]_{h,x}$ determines the dependence structure among $[O_{thx}]_{h,x}$. The structural equation for O_t is given by

$$O_t = f(H_t, X_t, \mathbf{V}_t) = \sum_{h,x} f_{hx}(V_{thx}) \mathbb{I}\{H_t = h, X_t = x\}, \quad (\text{EC.2a})$$

where $f_{hx}(\cdot)$ is the inverse CDF derived from the emission distribution $[e_{hxi}]_i$. This function generates $O_t | (H_t = h, X_t = x)$ using the noise variable $V_{thx} \sim \text{U}[0, 1]$.

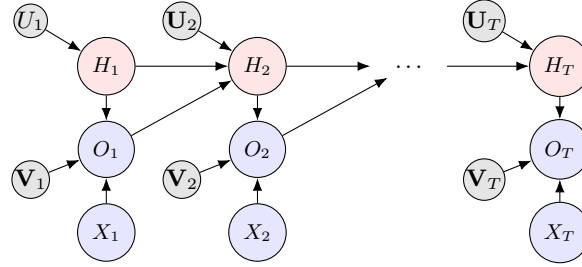


Figure EC.1 The SCM underlying the GHMM (Figure 1). The only difference between the SCM and the GHMM is the addition of exogenous noise nodes $[\mathbf{U}_t, \mathbf{V}_t]_t$.

Similarly, for $t > 1$, each $H_{thi} := H_t | (H_{t-1} = h, O_{t-1} = i)$ is a distinct random variable (or potential outcome) for all (h, i) . For time $t = 1$, there are no parents (H_0, O_0) for H_1 , so a single $\text{U}[0, 1]$ noise variable U_1 suffices to generate H_1 . For $t > 1$, we associate each H_{thi} with its own $\text{U}[0, 1]$ noise variable U_{thi} , such that

$$H_t = g(H_{t-1}, O_{t-1}, \mathbf{U}_t) = \sum_{h,i} g_{hi}(U_{thi}) \mathbb{I}\{H_{t-1} = h, O_{t-1} = i\} \quad (\text{EC.2b})$$

where $g_{hi}(\cdot)$ is the inverse CDF derived from the transition distribution $[q_{hih'}]_{h'}$. Here, $\mathbf{U}_t := [U_{thi}]_{h,i}$ consists of $|\mathbb{H}||\mathbb{O}|$ $\text{U}[0, 1]$ random variables.

The exogenous noise vectors $\{\mathbf{U}_t\}_{t=1}^T$ and $\{\mathbf{V}_t\}_{t=1}^T$ are assumed to be mutually independent across time and type. (Type = emissions vs transitions.) Specifically, (i) \mathbf{U}_t is independent of $\mathbf{U}_{t'}$ for all $t \neq t'$, (ii) \mathbf{V}_t is independent of $\mathbf{V}_{t'}$ for all $t \neq t'$, and (iii) \mathbf{U}_t is independent of $\mathbf{V}_{t'}$ for all t, t' . Such independence ensures that the SCM is consistent with the conditional-independence structure implied by the GHMM in Figure 1. This is exactly the restriction invoked in §4.2 to justify the local conditioning in (5). This assumption, akin to the commonly invoked *sequential ignorability* assumption in dynamic models, significantly reduces the space of feasible SCMs by prohibiting direct links between the \mathbf{U}_t 's and \mathbf{V}_t 's in Figure EC.1. Given the established validity of GHMMs for modeling disease progression (as discussed in earlier sections), this assumption seems very reasonable. Nonetheless, while this independence assumption defines a natural subclass of SCMs compatible with the underlying GHMM (an assumption on which our framework to bound

PN critically relies), there are other SCMs that would also be compatible with the underlying GHMM. For example, [Haugh and Singal \(2026\)](#) provide a simple example of such an SCM in the context of their dishonest casino model. The class of SCMs we do consider is still very rich, however, since we allow dependence among the components of $\mathbf{V}_t = [V_{thx}]_{h,x}$ and among the components of $\mathbf{U}_t = [U_{thi}]_{h,i}$.

The representation in (EC.2) enables us to model $[O_{thx}]_{h,x}$ ($[H_{thi}]_{h,i}$) and capture any dependence structure among these random variables by specifying the joint multivariate distribution of \mathbf{V}_t (\mathbf{U}_t). Since the univariate marginals of \mathbf{V}_t (\mathbf{U}_t) are known to be $U[0,1]$, specifying their multivariate distribution amounts to defining the dependence structure or the *copula* of \mathbf{V}_t (\mathbf{U}_t). For example, if the V_{thx} 's are mutually independent (the *independence copula*) and $H_t = h'$, then inferring the conditional distribution of $V_{th'x'}$ provides no information about the V_{thx} 's for $(h,x) \neq (h',x')$. Alternatively, if $V_{thx} = V_{th'x'}$ for all (h,x) and (h',x') , this models perfect positive dependence (the *comonotonic copula*), where inferring the conditional distribution of $V_{th'x'}$ simultaneously determines the conditional distribution of all V_{thx} 's.

We emphasize that the exogenous vectors \mathbf{U}_t and \mathbf{V}_t are critical for *counterfactual* analysis, as different joint distributions among the components of \mathbf{U}_t and among the components of \mathbf{V}_t can yield significantly different values of PN. However, if counterfactual analysis is not the goal and we only care about the joint distribution of a (subset of) $(O_{1:T}, H_{1:T})$, the analysis depends on \mathbf{U}_t and \mathbf{V}_t only through their known univariate marginals. In summary, the SCM is fully specified only when the dependence structure or copula is defined for each \mathbf{U}_t and \mathbf{V}_t . Finally, the emissions and the state transitions in the GHMM are time-homogeneous. It is, therefore, natural to also assume the copulas underlying \mathbf{V}_t and \mathbf{U}_t are time-homogeneous.

EC.3.2. Modeling the SCM via Response Distributions

While specifying the SCM in terms of the \mathbf{U}_t 's and \mathbf{V}_t 's via Figure EC.1 is conceptually useful, it is often more convenient to work with a direct but equivalent construction of the SCM. This alternative approach is well-established in the causal inference literature, where related concepts include *canonical partitions* and *principal strata* (see, e.g., [Duarte et al. 2024](#), [Zhang et al. 2022](#)).

Consider, for example, the relationship between $\mathbf{V}_t := [V_{thx}]_{h,x}$ and $[O_{thx}]_{h,x}$. From (EC.2), we know that the joint distribution of $[O_{thx}]_{h,x}$ is entirely determined by the joint distribution of \mathbf{V}_t . However, since our primary interest lies in the joint distribution of $[O_{thx}]_{h,x}$ itself, it is more natural to model this distribution *directly*, rather than indirectly via \mathbf{V}_t . Indeed, infinitely many joint distributions of \mathbf{V}_t can produce the same joint distribution of $[O_{thx}]_{h,x}$. This follows from Sklar's Theorem in copula theory ([Sklar 1959](#)).

Accordingly, in §4, we directly modeled the distributions of $[O_{hx}]_{h,x}$ and $[H_{hi}]_{h,i}$ using the $\boldsymbol{\theta}$ and $\boldsymbol{\pi}$ variables (response distributions). Specifically, we defined the pairwise marginals in (6) and the

full joint PMFs in (9). Nonetheless, the copula framework remains valuable in §7 where we leverage the independence and comonotonic copulas to evaluate the degree of sub-optimality introduced by our approximation for computational scalability.

Appendix EC.4: Using Copulas to Estimate the Probability of Necessity

Building on the SCM presented in §EC.3.1, we now describe the Monte Carlo simulation algorithms that we use in §6 and §7 for estimating the PN values for the independence (§EC.4.1) and comonotonic (§EC.4.2) SCMs. As stated in §EC.3.1, the SCM is fully specified only when the dependence structure or copula is defined for each \mathbf{U}_t and \mathbf{V}_t . When we refer to the independence (comonotonic) SCM below and indeed in the main body of the paper, what we have in mind is that the dependence structure of each \mathbf{U}_t and \mathbf{V}_t is given by the independence (comonotonic) copula. Of course, it would be easy to obtain other SCMs by allowing the \mathbf{U}_t 's to have one dependence structure, e.g., the independence copula, and the \mathbf{V}_t 's to have a different dependence structure, e.g., the monotonic copula. Indeed, each \mathbf{U}_t and \mathbf{V}_t could have their own different copulas and we would still obtain a legitimate SCM. There are countless combinations and so we will focus on just two SCMs here: the independence SCM where each \mathbf{U}_t and \mathbf{V}_t has the independence copula, and the comonotonic SCM where each \mathbf{U}_t and \mathbf{V}_t has the comonotonic copula.

EC.4.1. Counterfactual Simulations Under the Independence SCM

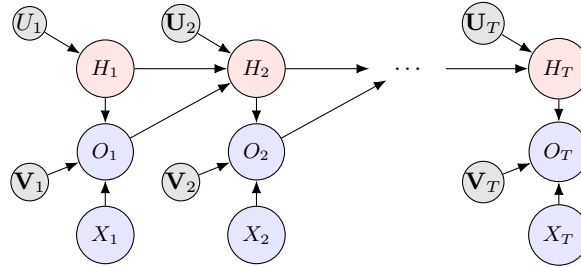


Figure EC.2 The same structural causal model as in Figure EC.1.

For convenience, we replicate Figure EC.1 from §EC.3.1, which now appears as Figure EC.2. Recall that $(o_{1:T}, x_{1:T})$ represents the observed data, and $\tilde{x}_{1:T}$ denotes the intervention policy. As in §4.3, we begin with the posterior samples $[h_{1:T}(b)]_{b=1}^B$, corresponding to the random path $H_{1:T} \mid (o_{1:T}, x_{1:T})$. These samples can be generated efficiently, as established in §4.1. For each b , our goal is to convert the sampled path $h_{1:T}(b)$ into a counterfactual path, denoted by $\tilde{h}_{1:T}(b)$. As noted in §4.2, the counterfactual hidden state in period 1 equals the posterior sample, i.e., $\tilde{h}_1(b) = h_1(b)$. To sample $\tilde{h}_2(b)$, we first need to sample the counterfactual emission $\tilde{o}_1(b)$ (cf. Figure EC.2). With the copula underlying \mathbf{V}_1 being the independence copula, it follows that

$$\tilde{o}_1(b) = \begin{cases} o_1 & \text{if } x_1 = \tilde{x}_1 \text{ and } h_1(b) = \tilde{h}_1(b) \\ \text{sample from the emission distribution } [e_{\tilde{h}_1(b)\tilde{x}_1}^i]_i & \text{otherwise.} \end{cases}$$

The counterfactual emission $\tilde{o}_1(b)$ allows us to sample the counterfactual state $\tilde{h}_2(b)$, which again leverages the fact that the copula underlying \mathbf{U}_2 is the independence copula:

$$\tilde{h}_2(b) = \begin{cases} h_2(b) & \text{if } h_1(b) = \tilde{h}_1(b) \text{ and } o_1 = \tilde{o}_1(b) \\ \text{sample from the transition distribution } [q_{\tilde{h}_1(b)\tilde{o}_1(b)h'}]_{h'} & \text{otherwise.} \end{cases}$$

We then generate the period 2 counterfactual emission $\tilde{o}_2(b)$ in a similar manner, and the process repeats iteratively until the end of the horizon T . The complete procedure is summarized in Algorithm 4.

Algorithm 4 Counterfactual simulations under the independence SCM

Require: (\mathbf{E}, \mathbf{Q}) , $(o_{1:T}, x_{1:T})$, $[h_{1:T}(b)]_{b=1}^B$, $\tilde{x}_{1:T}$

```

1: for  $b = 1$  to  $B$  do
2:    $\tilde{h}_1(b) = h_1(b)$ 
3:   for  $t = 1$  to  $T - 1$  do
4:     if  $x_t = \tilde{x}_t$  and  $h_t(b) = \tilde{h}_t(b)$  then
5:        $\tilde{o}_t(b) = o_t$ 
6:     else
7:        $\tilde{o}_t(b) \sim \text{Categorical}([e_{\tilde{h}_t(b)\tilde{x}_t i}]_i)$ 
8:     end if
9:     if  $h_t(b) = \tilde{h}_t(b)$  and  $o_t = \tilde{o}_t(b)$  then
10:       $\tilde{h}_{t+1}(b) = h_{t+1}(b)$ 
11:    else
12:       $\tilde{h}_{t+1}(b) \sim \text{Categorical}([q_{\tilde{h}_t(b)\tilde{o}_t(b)h'}]_{h'})$ 
13:    end if
14:  end for
15: end for
16: return  $[\tilde{h}_{1:T}(b)]_b$ 

```

EC.4.2. Counterfactual Simulations Under the Comonotonic SCM

Before formally describing how to estimate PN using the comonotonic SCM, we provide some intuition for our approach via a simple example consistent with the setup of Example 1. Consider two variables, X and Y , where $X \in \{0, 1\}$ represents the medical treatment, and $Y \in \{\text{bad}, \text{better}, \text{best}\}$ represents the patient outcome. The outcome $Y_x := Y \mid (X = x)$ obeys the following distribution: $Y_0 \sim \{\text{bad}, \text{better}, \text{best}\}$ w.p. $\{0.2, 0.3, 0.5\}$ and $Y_1 \sim \{\text{bad}, \text{better}, \text{best}\}$ w.p. $\{0.2, 0.2, 0.6\}$. The underlying SCM is shown in Figure EC.3. Note that X is deterministic and does not have an exogenous noise variable pointing to it.

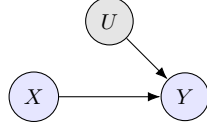


Figure EC.3 The comonotonic copula SCM for the “simple” example. The noise node is represented by the scalar $U \sim U[0, 1]$ rather than a vector \mathbf{U} . The structural equation is $Y = f(X, U)$, which we denote by $f_X(U)$, the *inverse transform* function corresponding to Y_x . That is, $f_0(u) = \text{bad, better, and best}$ if $u \in [0, 0.2]$, $u \in [0.2, 0.5]$, and $u \in [0.5, 1]$, resp. Similarly, $f_1(u) = \text{bad, better, and best}$ if $u \in [0, 0.2]$, $u \in [0.2, 0.4]$, and $u \in [0.4, 1]$, resp.

Consider a patient whose outcome Y was “better” under no treatment ($x = 0$). Given the prior $U \sim U[0, 1]$, the posterior distribution is $U \mid (Y_0 = \text{better}) \sim U[0.2, 0.5]$. Now, suppose we are interested in the counterfactual outcome under the intervention $\tilde{x} = 1$, i.e., the random variable $\tilde{Y} := Y_1 \mid (Y_0 = \text{better})$. Using the posterior $U[0.2, 0.5]$ for U and the functional form of $f_1(\cdot)$ (as defined in the caption of Figure EC.3), we observe that the interval $[0.2, 0.4]$ maps to “better” and $[0.4, 0.5]$ maps to “best”. Consequently, \tilde{Y} is “better” w.p. $2/3$ and “best” w.p. $1/3$.

We now generalize this approach to the GHMM. We first fix an ordering of the states (set \mathbb{H}) and the emissions (set \mathbb{O}), e.g., from “best” to “worst”. Denote by $r_H(h)$ the rank of state h with respect to this ordering and by $r_O(i)$ the rank of emission i . Furthermore, let $r_H^{-1}(r)$ and $r_O^{-1}(r)$ denote the inverse functions corresponding to $r_H(h)$ and $r_O(i)$, respectively. That is, $r_H^{-1}(r)$ returns the state with rank r and $r_O^{-1}(r)$ returns the emission with rank r . Also, for each (h, i) pair, observe that $[q_{hih'}]_{h'}$ denotes the transition distribution (which maps to the random variable H_{hi}). Corresponding to this distribution, we define the rank-ordered CDF as follows:

$$Q_{hih'} := \sum_{h'' : r_H(h'') \leq r_H(h')} q_{hih''} \quad \forall h'. \quad (\text{EC.3a})$$

Similarly, for each (h, x) pair, observe that $[e_{hxi}]_i$ denotes the emission distribution (which maps to the random variable O_{hx}). Corresponding to this distribution, we define the rank-ordered CDF as follows:

$$E_{hxi} := \sum_{j : r_O(j) \leq r_O(i)} e_{hxj} \quad \forall i. \quad (\text{EC.3b})$$

Also, define $Q_{hi0} = E_{hx0} = 0$ for all (h, i) and (h, x) . We discuss these orderings for the breast cancer application in §EC.5.4.

As in §EC.4.1, we begin with the samples $[h_{1:T}(b)]_{b=1}^B$, which we generate from the posterior distribution of $H_{1:T} \mid (o_{1:T}, x_{1:T})$. For each sample path $h_{1:T}(b)$, our goal is to generate a counterfactual path $\tilde{h}_{1:T}(b)$. As noted in §4.2, irrespective of the copula choice, the counterfactual hidden state in period 1 equals the posterior sample, i.e., $\tilde{h}_1(b) = h_1(b)$. To generate $\tilde{o}_1(b)$, consider the SCM in Figure EC.4. The exogenous noise nodes are now scalar $U[0, 1]$ random variables rather than

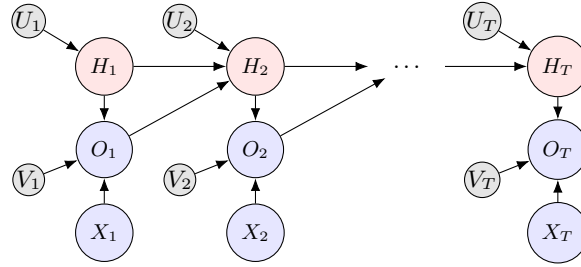


Figure EC.4 The comonotonic SCM underlying the GHMM. The key observation is that the noise nodes are now scalar $U[0, 1]$ random variables as opposed to random vectors, i.e., (U_t, V_t) as opposed to $(\mathbf{U}_t, \mathbf{V}_t)$.

random vectors with $U[0, 1]$ marginals. This is a direct implication of the comonotonic copula (as discussed in §EC.3.1). It follows from the structural equation (EC.2a) that $\tilde{o}_1(b)$ satisfies

$$\tilde{o}_1(b) = f(\tilde{h}_1(b), \tilde{x}_1, V_1) = f_{\tilde{h}_1(b)\tilde{x}_1}(V_1), \quad (\text{EC.4a})$$

where $f_{hx}(\cdot)$ is the inverse transform function corresponding to the rank-ordered CDF $[E_{hxi}]_i$ (recall (EC.3b)). Hence, all we need to sample $\tilde{o}_1(b)$ is the posterior distribution of V_1 , where the “posterior” corresponds to conditioning on $O_{th_1(b)x_1} = o_1$ where $O_{thx} := O_t | (H_t = h, X_t = x)$ was defined in §EC.3.1. (We must keep the t notation here since the posterior dynamics are time-dependent.) Given the prior $V_1 \sim U[0, 1]$, the posterior of V_1 satisfies

$$V_1 | (O_{1h_1(b)x_1} = o_1) \sim U[E_{h_1(b)x_1 o_1^-}, E_{h_1(b)x_1 o_1}], \quad (\text{EC.4b})$$

where $o^- := r_O^{-1}(r_O(o) - 1)$ is the emission ranked just below o . Hence, we can efficiently sample V_1 from its posterior, and this V_1 sample can be used to generate $\tilde{o}_1(b)$ (via (EC.4a)).

We can sample $\tilde{h}_2(b)$ similarly. By the structural equation (EC.2b), $\tilde{h}_2(b)$ satisfies

$$\tilde{h}_2(b) = g(\tilde{h}_1(b), \tilde{o}_1(b), U_2) = g_{\tilde{h}_1(b)\tilde{o}_1(b)}(U_2), \quad (\text{EC.5a})$$

where $g_{hi}(\cdot)$ is the inverse transform function corresponding to the rank-ordered CDF $[Q_{nih'}]_{h'}$ (recall (EC.3a)). Hence, all we need to sample $\tilde{h}_2(b)$ is the posterior distribution of U_2 , where the “posterior” corresponds to conditioning on $H_{2h_1(b)o_1} = h_2(b)$ where $H_{thi} := H_t | (H_{t-1} = h, O_{t-1} = i)$ was defined in §EC.3.1. (Once again, we must keep the t notation here since the posterior dynamics are time-dependent.) Given the prior $U_2 \sim U[0, 1]$, its posterior satisfies

$$U_2 | (H_{2h_1(b)o_1} = h_2(b)) \sim U[Q_{h_1(b)o_1 h_2(b)^-}, Q_{h_1(b)o_1 h_2(b)}], \quad (\text{EC.5b})$$

where $h^- := r_H^{-1}(r_H(h) - 1)$ is the state ranked just below h . Hence, we can efficiently sample U_2 from its posterior, and this U_2 sample can be used to generate $\tilde{h}_2(b)$ (via (EC.5a)).

We can now generate the period 2 counterfactual emission $\tilde{o}_2(b)$ in a similar manner and the process repeats until we reach the end of horizon T . We summarize the procedure in Algorithm 5.

Algorithm 5 Counterfactual simulations under the comonotonic SCM**Require:** (\mathbf{E}, \mathbf{Q}) , $(o_{1:T}, x_{1:T})$, $[h_{1:T}(b)]_{b=1}^B$, $\tilde{x}_{1:T}$, $r_H(\cdot)$, $r_O(\cdot)$

```

1: for  $b = 1$  to  $B$  do
2:    $\tilde{h}_1(b) = h_1(b)$ 
3:   for  $t = 1$  to  $T - 1$  do
4:      $V_t \sim \text{U}[E_{h_t(b)x_t o_t^-}, E_{h_t(b)x_t o_t}]$            % posterior sample of  $V_t$  (see (EC.4b))
5:      $\tilde{o}_t(b) = f_{\tilde{h}_t(b)\tilde{x}_t}(V_t)$                              % counterfactual emission (see (EC.4a))
6:      $u_{t+1} \sim \text{U}[Q_{h_t(b)o_t h_{t+1}(b)^-}, Q_{h_t(b)o_t h_{t+1}(b)}]$  % posterior sample of  $U_{t+1}$  (see (EC.5b))
7:      $\tilde{h}_{t+1}(b) = g_{\tilde{h}_t(b)\tilde{o}_t(b)}(u_{t+1})$                  % counterfactual state (see (EC.5a))
8:   end for
9: end for
10: return  $[\tilde{h}_{1:T}(b)]_b$ 

```

Remark EC.1. It should be intuitively clear that the comonotonic copula implies pathwise monotonicity. Our discussion of the simple example at the beginning of this sub-appendix further illustrates this point. In that example, a patient’s outcome Y under no treatment ($x = 0$) was “better”. The counterfactual outcome \tilde{Y} for this patient under treatment ($\tilde{x} = 1$) in the SCM with the comonotonic copula was either “better” (w.p. 2/3) or “best” (w.p. 1/3), thereby satisfying pathwise monotonicity.

Appendix EC.5: Further Details on the Case Study

We discuss further details on the breast cancer model primitives and their calibration in §EC.5.1, followed by showing how we exploit sparsity to reduce the number of decision variables (§EC.5.2). We then provide details on the pathwise monotonicity constraints and the comonotonic copula in §EC.5.3 and §EC.5.4, respectively.

EC.5.1. Calibration of Initial State Distribution \mathbf{p} and Emission Distribution \mathbf{E}

We have $\mathbf{p} := (p_1, \dots, p_7)$. Typically, breast cancer screening begins around age 40, with a prevalence among females aged 40-49 of 1.0183% (Table 4.24 of NIH (2020), all races, females): $p_2 + p_3 = 0.010183$. In-situ cancer accounts for 20% of new breast cancer diagnoses (Sprague and Trentham-Dietz 2009), giving $p_2 = 0.2 \times 0.010183$ and $p_3 = 0.8 \times 0.010183$. It is natural to set $p_4 = p_5 = p_6 = p_7 = 0$, leading to $p_1 = 1 - p_2 - p_3 = 1 - 0.010183$.

We have $\mathbf{E} := [e_{hxi}]_{h,x,i}$, where $e_{hxi} := \mathbb{P}(O_t = i \mid H_t = h, X_t = x)$. Before discussing the calibration, we highlight the sparse structure of \mathbf{E} . To illustrate this, we define the matrix $\mathbf{E}(x) := [e_{hxi}]_{hi}$ for each action x (with each row summing to 1), and observe the following structure:

State h	Policy x	Range of O_{hx}	Range cardinality
1	0	{2,3}	2
1	1	{1}	1
2	0	{2,4}	2
2	1	{1}	1
3	0	{2,5}	2
3	1	{1}	1
4	0	{4}	1
4	1	{4}	1
5	0	{5}	1
5	1	{5}	1
6	0	{6}	1
6	1	{6}	1
7	0	{7}	1
7	1	{7}	1

Table EC.1 Range of the 14 random variables $[O_{hx}]_{h,x}$ corresponding to $\theta_{1,\dots,K}$.

Consider the $\theta_{1,\dots,K}$ decision variables, which represent the joint PMF of $[O_k]_k$, where $k \equiv (h, x)$. To reduce the number of variables, we first determine the valid (h, x) pairs, rather than naively considering all $(h, x) \in \mathbb{H} \times \mathbb{X}$. Recall that the state h is encoded as follows: (1) healthy, (2) undiagnosed in-situ cancer, (3) undiagnosed invasive cancer, (4) diagnosed in-situ cancer, (5) diagnosed invasive cancer, (6) recovery, and (7) death. Furthermore, $x = 0$ indicates a mammogram is performed, while $x = 1$ indicates it is not performed. It is straightforward to verify that all 14 combinations of $(h, x) \in \mathbb{H} \times \mathbb{X}$ are valid, meaning none of the corresponding decision variables can be removed. Turning to the observations, they are encoded as follows: (1) no screening, (2) negative screening result, (3) positive mammogram followed by a negative biopsy, (4) diagnosed in-situ cancer, (5) diagnosed invasive cancer, (6) recovery, and (7) death. Table EC.1 summarizes the range of all $7 \times 2 = 14$ random variables $[O_{hx}]_{h,x}$. Multiplying the cardinalities from the last column of Table EC.1 reveals that only eight $\theta_{1,\dots,K}$ decision variables need to be considered, significantly fewer than the naive upper bound of $|\mathbb{O}|^{|\mathbb{H}||\mathbb{X}|} = 7^{14}$.

The same logic applies to the $\pi_{1,\dots,M}$ decision variables, which represent the joint PMF of $[H_m]_m$, where $m \equiv (h, i)$. For the $\pi_{1,\dots,M}$ decision variables, even the first step proves useful as some (h, i) pairs are invalid. For example, if $h = 1$ (healthy), then $i \notin \{4, 5, 6, 7\}$, as these emissions correspond to diagnosed cancer, recovery, or death. This initial step reduces the number of $[H_{hi}]_{h,i}$ random variables from $|\mathbb{H}||\mathbb{O}| = 49$ to 13 valid pairs. The second step trims the range of each of these 13 random variables. Table EC.2 documents these reductions, showing how the number of $\pi_{1,\dots,M}$ decision variables is reduced from the naive upper bound of $|\mathbb{H}|^{|\mathbb{H}||\mathbb{O}|} = 7^{49}$ to 15,552, which equals the product of the cardinalities presented in the last column of Table EC.2.

EC.5.3. Details on the Pathwise Monotonicity Constraints

Pathwise monotonicity can be enforced via linear constraints. We discussed this in §5 and now provide a comprehensive description of all pathwise monotonicity constraints embedded in our

State h	Observation i	Range of H_{hi}	Range cardinality
1	1	$\{1, 2, 3\}$	3
1	2	$\{1, 2, 3\}$	3
1	3	$\{1, 2, 3\}$	3
2	1	$\{2, 3\}$	2
2	2	$\{2, 3\}$	2
2	4	$\{4, 6\}$	2
3	1	$\{3, 7\}$	2
3	2	$\{3, 7\}$	2
3	5	$\{5, 6, 7\}$	3
4	4	$\{4, 6\}$	2
5	5	$\{5, 6, 7\}$	3
6	6	$\{6\}$	1
7	7	$\{7\}$	1

Table EC.2 Range of the 13 random variables $[H_{hi}]_{h,i}$ corresponding to $\pi_{1,\dots,M}$. Only 13 (h,i) pairs are shown as the other 36 are not valid.

breast cancer numerics. To recap from §5, suppose the patient has “low severity” cancer (state h) in period t , the cancer has not been detected (observation i), and the patient’s state remains in “low severity” in period $t + 1$. In the counterfactual world, if the cancer is detected in period t (observation \bar{i}), pathwise monotonicity requires that the cancer cannot transition to a worse state than “low severity” in period $t + 1$. Formally, this implies $\mathbb{P}(H_{h\bar{i}} = \bar{h}' \mid H_{hi} = h) = 0$ for all states \bar{h}' worse than h . There are multiple such cases to consider. We enforce all pathwise monotonicity constraints by setting the corresponding $\pi_{\bar{h}\bar{i},hi}(\bar{h}', h')$ variables to 0, as $\pi_{\bar{h}\bar{i},hi}(\bar{h}', h') = \mathbb{P}(H_{hi} = h')\mathbb{P}(H_{h\bar{i}} = \bar{h}' \mid H_{hi} = h')$.

Hence, to provide details on all pathwise monotonicity constraints we enforce, it suffices to enumerate the $(h, i, h', \bar{h}, \bar{i}, \bar{h}')$ combinations for which the $\pi_{\bar{h}\bar{i},hi}(\bar{h}', h')$ variables are set to 0. To achieve this, we iterate over each state $h \in \{1, \dots, 7\}$. (Note that for pathwise monotonicity, there are no $(h, x, \bar{i}, \bar{h}, \bar{x}, \bar{i})$ combinations for which we set the $\theta_{\bar{h}\bar{x},hx}(\bar{i}, i)$ variables equal to 0.)

Healthy ($h = 1$). Pathwise monotonicity is enforced for the following combinations:

- If (h, i, h') equals (healthy, whatever emission, healthy), then the counterfactual state \bar{h}' can not be in-situ, invasive, or death if \bar{h} is healthy. That is, $h = 1, i \in \mathbb{O}, h' = 1, \bar{h} = 1, \bar{i} \in \mathbb{O}$, and $\bar{h}' \in \{2, 3, 4, 5, 7\}$.
- If (h, i, h') equals (healthy, whatever emission, in-situ), then the counterfactual state \bar{h}' can not be healthy, invasive, or death if \bar{h} is healthy. That is, $h = 1, i \in \mathbb{O}, h' = 2, \bar{h} = 1, \bar{i} \in \mathbb{O}$, and $\bar{h}' \in \{1, 3, 5, 7\}$.
- If (h, i, h') equals (healthy, whatever emission, invasive), then the counterfactual state \bar{h}' can not be healthy, in-situ, or death if \bar{h} is healthy. That is, $h = 1, i \in \mathbb{O}, h' = 3, \bar{h} = 1, \bar{i} \in \mathbb{O}$, and $\bar{h}' \in \{1, 2, 4, 7\}$.

- If (h, i, h') equals (healthy, whatever emission, death), then the counterfactual state \bar{h}' can not be healthy, in-situ, or invasive if \bar{h} is healthy. That is, $h = 1$, $i \in \mathbb{O}$, $h' = 7$, $\bar{h} = 1$, $\bar{i} \in \mathbb{O}$, and $\bar{h}' \in \{1, 2, 3, 4, 5\}$.

Undiagnosed In-situ ($h = 2$). Pathwise monotonicity is enforced for the following combinations:

- If (h, i, h') equals (in-situ, undetected, in-situ), then the counterfactual state \bar{h}' can not be invasive or death if \bar{h} is healthy or in-situ. That is, $h = 2$, $i \in \{1, 2, 3\}$, $h' = 2$, $\bar{h} \in \{1, 2, 4\}$, $\bar{i} \in \mathbb{O}$, and $\bar{h}' \in \{3, 5, 7\}$.
- If (h, i, h') equals (in-situ, detected, in-situ), then the counterfactual state \bar{h}' can not be invasive or death if \bar{h} is in-situ and detected. That is, $h = 2$, $i = 4$, $h' = 4$, $\bar{h} \in \{2, 4\}$, $\bar{i} = 4$, and $\bar{h}' \in \{3, 5, 7\}$.
- If (h, i, h') equals (in-situ, undetected, invasive), then the counterfactual state \bar{h}' can not be death if \bar{h} is healthy or in-situ. That is, $h = 2$, $i \in \{1, 2, 3\}$, $h' = 3$, $\bar{h} \in \{1, 2, 4\}$, $\bar{i} \in \mathbb{O}$, and $\bar{h}' = 7$.
- If (h, i, h') equals (in-situ, detected, invasive), then the counterfactual state \bar{h}' can not be death if \bar{h} is in-situ and detected. That is, $h = 2$, $i = 4$, $h' = 5$, $\bar{h} \in \{2, 4\}$, $\bar{i} = 4$, and $\bar{h}' = 7$.
- If (h, i, h') equals (in-situ, detected, recovered), then the counterfactual state \bar{h}' can not be in-situ, invasive, or death if \bar{h} is in-situ and detected. That is, $h = 2$, $i = 4$, $h' = 6$, $\bar{h} \in \{2, 4\}$, $\bar{i} = 4$, and $\bar{h}' \in \{2, 3, 4, 5, 7\}$.
- If (h, i, h') equals (in-situ, undetected, death), then the counterfactual state \bar{h}' can not be in-situ, invasive, or recovered if \bar{h} is in-situ and undetected. That is, $h = 2$, $i \in \{1, 2, 3\}$, $h' = 7$, $\bar{h} = 2$, $\bar{i} \in \{1, 2, 3\}$, and $\bar{h}' \in \{2, 3, 4, 5, 6\}$.
- If (h, i, h') equals (in-situ, detected, death), then the counterfactual state \bar{h}' can not be in-situ, invasive, or recovered if \bar{h} is in-situ and detected. That is, $h = 2$, $i = 4$, $h' = 7$, $\bar{h} \in \{2, 4\}$, $\bar{i} = 4$, and $\bar{h}' \in \{2, 3, 4, 5, 6\}$.

Undiagnosed Invasive ($h = 3$). Pathwise monotonicity enforced for the following combinations:

- If (h, i, h') equals (invasive, undetected, invasive), then the counterfactual state \bar{h}' can not be death if \bar{h} is healthy, in-situ, or invasive. That is, $h = 3$, $i \in \{1, 2, 3\}$, $h' = 3$, $\bar{h} \in \{1, 2, 3, 4, 5\}$, $\bar{i} \in \mathbb{O}$, and $\bar{h}' = 7$.
- If (h, i, h') equals (invasive, detected, invasive), then the counterfactual state \bar{h}' can not be death if \bar{h} is invasive and detected. That is, $h = 3$, $i = 5$, $h' = 5$, $\bar{h} \in \{3, 5\}$, $\bar{i} = 5$, and $\bar{h}' = 7$.
- If (h, i, h') equals (invasive, detected, recovered), then the counterfactual state \bar{h}' can not be invasive or death if \bar{h} is invasive and detected. That is, $h = 3$, $i = 5$, $h' = 6$, $\bar{h} \in \{3, 5\}$, $\bar{i} = 5$, and $\bar{h}' \in \{3, 5, 7\}$.

- If (h, i, h') equals (invasive, undetected, death), then the counterfactual state \bar{h}' can not be invasive or recovered if \bar{h} is invasive and undetected. That is, $h = 3$, $i \in \{1, 2, 3\}$, $h' = 7$, $\bar{h} \in \{3, 5\}$, $\bar{i} \in \{1, 2, 3\}$, and $\bar{h}' \in \{3, 5, 6\}$.
- If (h, i, h') equals (invasive, detected, death), then the counterfactual state \bar{h}' can not be invasive or recovered if \bar{h} is invasive and detected. That is, $h = 3$, $i = 5$, $h' = 7$, $\bar{h} \in \{3, 5\}$, $\bar{i} = 5$, and $\bar{h}' \in \{3, 5, 6\}$.

Diagnosed In-situ ($h = 4$). Pathwise monotonicity is enforced for the following combinations:

- If (h, i, h') equals (in-situ, detected, in-situ), then the counterfactual state \bar{h}' can not be invasive, recovered, or death if \bar{h} is in-situ and detected. That is, $h = 4$, $i = 4$, $h' = 4$, $\bar{h} \in \{2, 4\}$, $\bar{i} = 4$, and $\bar{h}' \in \{5, 6, 7\}$.
- If (h, i, h') equals (in-situ, detected, invasive), then the counterfactual state \bar{h}' can not be in-situ, recovered, or death if \bar{h} is in-situ and detected. That is, $h = 4$, $i = 4$, $h' = 5$, $\bar{h} \in \{2, 4\}$, $\bar{i} = 4$, and $\bar{h}' \in \{4, 6, 7\}$.
- If (h, i, h') equals (in-situ, detected, recovery), then the counterfactual state \bar{h}' can not be in-situ, invasive, or death if \bar{h} is in-situ and detected. That is, $h = 4$, $i = 4$, $h' = 6$, $\bar{h} \in \{2, 4\}$, $\bar{i} = 4$, and $\bar{h}' \in \{4, 5, 7\}$.
- If (h, i, h') equals (in-situ, detected, death), then the counterfactual state \bar{h}' can not be in-situ, invasive, or recovered if \bar{h} is in-situ and detected. That is, $h = 4$, $i = 4$, $h' = 7$, $\bar{h} \in \{2, 4\}$, $\bar{i} = 4$, and $\bar{h}' \in \{4, 5, 6\}$.

Diagnosed Invasive ($h = 5$). Pathwise monotonicity is enforced for the following combinations:

- If (h, i, h') equals (invasive, detected, invasive), then the counterfactual state \bar{h}' can not be recovered or death if \bar{h} is invasive and detected. That is, $h = 5$, $i = 5$, $h' = 5$, $\bar{h} \in \{3, 5\}$, $\bar{i} = 5$, and $\bar{h}' \in \{6, 7\}$.
- If (h, i, h') equals (invasive, detected, recovery), then the counterfactual state \bar{h}' can not be invasive or death if \bar{h} is invasive and detected. That is, $h = 5$, $i = 5$, $h' = 6$, $\bar{h} \in \{3, 5\}$, $\bar{i} = 5$, and $\bar{h}' \in \{5, 7\}$.
- If (h, i, h') equals (invasive, detected, death), then the counterfactual state \bar{h}' can not be invasive or recovery if \bar{h} is invasive and detected. That is, $h = 5$, $i = 5$, $h' = 7$, $\bar{h} \in \{3, 5\}$, $\bar{i} = 5$, and $\bar{h}' \in \{5, 6\}$.

Recovery ($h = 6$) and Death ($h = 7$). Pathwise monotonicity enforced for no combinations.

EC.5.4. Details on the Comonotonic Copula

We discussed the counterfactual simulation under the comonotonic copula for a general GHMM in §EC.4.2. In this section, we connect that discussion to the breast cancer application. To do so, it

suffices to define the rank functions $r_H(\cdot)$ (for states) and $r_O(\cdot)$ (for emissions). For states, there are two possible orderings that seem “natural” (from “best” to “worst”):

- (1, 6, 4, **2, 5**, 3, 7)
- (1, 6, 4, **5, 2**, 3, 7).

Recalling the $\mathbf{Q}(i)$ notation from §6, note that columns 2 and 5 are never “active” simultaneously in any row of $\mathbf{Q}(i)$ (for any i). Thus, the choice of ordering between the two options above is inconsequential, and we can select either. Suppose we choose the first ordering. This defines the rank function, e.g., $r_H(6) = 2$, meaning the rank of state 6 is 2. For the inverse function, $r_H^{-1}(2) = 6$.

Defining $r_O(\cdot)$ for the breast cancer setting is unclear, but fortunately, it does not matter. To see why, consider the following path of interest, which generalizes both Paths 1 and 2 from §6:

$$\underbrace{O_1, \dots, O_{\tau_s-1}}_{\in\{2,3\}}, \underbrace{O_{\tau_s}, \dots, O_{\tau_e}}_{=1}, \underbrace{O_{\tau_e+1}, \dots, O_{T-1}}_{\in\{4,5\}}, \underbrace{O_T}_{=7}.$$

For the first $\tau_s - 1$ periods, observe that the counterfactual emission $\tilde{o}_t(b)$ equals the observed emission o_t for each b . This is because the intervention policy \tilde{x}_t equals the observed policy x_t for $t \leq \tau_s - 1$. Now, consider periods τ_s to τ_e , during which the screening was not done, i.e., $x_{\tau_s:\tau_e} = 1$. Hence, the corresponding emissions $o_t = 1$ w.p. 1 (see the matrix $\mathbf{E}(1)$ in §EC.5.1). This means that the emissions do not contain any information regarding the underlying noise variables $V_{\tau_s:\tau_e}$ (see Figure EC.4) and hence, their posterior equals their prior, which is $U[0, 1]$. As such, for $t \in \{\tau_s, \dots, \tau_e\}$, we can sample $\tilde{o}_t(b)$ using the categorical distribution over the probability vector $[e_{\tilde{h}_t(b)\tilde{x}_t i}]_i$. Note that we can use $\tilde{o}_{\tau_s-1}(b)$ to sample $\tilde{h}_{\tau_s}(b)$, which we can then use to sample $\tilde{o}_{\tau_s}(b)$, and so on (until we have sampled $\tilde{h}_{\tau_e+1}(b)$). Now, consider $t = \tau_e + 1$. We know $o_t \in \{4, 5\}$:

- If $o_t = 4$, then $h_t(b) = 2$ and $\tilde{h}_t(b) \in \{2, 4, 6\}$ (cf. pathwise monotonicity).
 - If $\tilde{h}_t(b) \in \{4, 6\}$, then $\tilde{o}_t(b) = \tilde{h}_t(b)$ (since rows 4 and 6 of $\mathbf{E}(0)$ have 1 on the diagonal).
 - Else, if $\tilde{h}_t(b) = 2$ ($= h_t(b)$), then $\tilde{o}_t(b) = o_t = 4$.
- Else, if $o_t = 5$, then $h_t(b) = 3$ and $\tilde{h}_t(b) \in \{3, 4, 5, 6\}$ (cf. pathwise monotonicity).
 - If $\tilde{h}_t(b) \in \{4, 5, 6\}$, then $\tilde{o}_t(b) = \tilde{h}_t(b)$ (since rows 4, 5, 6 of $\mathbf{E}(0)$ have 1 on the diagonal).
 - If $\tilde{h}_t(b) = 3$ ($= h_t(b)$), then $\tilde{o}_t(b) = o_t = 5$.

Finally, for $t \geq \tau_e + 2$, we know $h_t(b) \geq 4$ and that the corresponding rows in $\mathbf{E}(0)$ are 0-1. Hence, the posterior of V_t equals the prior and we can sample $\tilde{o}_t(b)$ using the categorical distribution over the probability vector $[e_{\tilde{h}_t(b)\tilde{x}_t i}]_i$. By construction, the comonotonic copula will obey pathwise monotonicity and hence, will ensure that in the counterfactual world, the patient does not die before period T .