

## Proof of the causal model

We need to calculate  $p(\mathbf{e}|\mathbf{x}, \mathbf{c})$ . Applying Bayes' theorem, we have that

$$p(\mathbf{e}|\mathbf{x}, \mathbf{c}) = \frac{p(\mathbf{x}|\mathbf{e}, \mathbf{c})p(\mathbf{e}|\mathbf{c})}{p(\mathbf{x}|\mathbf{c})}. \quad (1)$$

If the causal model is captured by the Directed Acyclic Graph (DAG)  $\mathbf{c} \rightarrow \mathbf{e} \rightarrow \mathbf{x}$ , we also have that

$$p(\mathbf{x}|\mathbf{e}, \mathbf{c}) = p(\mathbf{x}|\mathbf{e}). \quad (2)$$

Applying Bayes' theorem again, we have that

$$p(\mathbf{x}|\mathbf{e}) = \frac{p(\mathbf{e}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{e})}. \quad (3)$$

Thanks to equations 2 and 3, we can substitute  $p(\mathbf{x}|\mathbf{e}, \mathbf{c})$  in equation 1 with  $\frac{p(\mathbf{e}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{e})}$  and write

$$p(\mathbf{e}|\mathbf{x}, \mathbf{c}) = \frac{p(\mathbf{e}|\mathbf{x})p(\mathbf{e}|\mathbf{c})p(\mathbf{x})}{p(\mathbf{x}|\mathbf{c})p(\mathbf{e})}. \quad (4)$$

Now, because  $p(\mathbf{e}|\mathbf{x}, \mathbf{c})$  is a probability, we know that

$$\int_{\mathbf{e}} p(\mathbf{e}|\mathbf{x}, \mathbf{c}) d\mathbf{e} = 1. \quad (5)$$

But, because of equation 4, this also implies that

$$\int_{\mathbf{e}} \frac{p(\mathbf{e}|\mathbf{x})p(\mathbf{e}|\mathbf{c})p(\mathbf{x})}{p(\mathbf{x}|\mathbf{c})p(\mathbf{e})} d\mathbf{e} = 1. \quad (6)$$

We can observe that  $p(\mathbf{x})$  and  $p(\mathbf{x}|\mathbf{c})$  are constant with respect to  $\mathbf{e}$ , therefore they can be taken outside the integral:

$$\frac{p(\mathbf{x})}{p(\mathbf{x}|\mathbf{c})} \int_{\mathbf{e}} \frac{p(\mathbf{e}|\mathbf{x})p(\mathbf{e}|\mathbf{c})}{p(\mathbf{e})} d\mathbf{e} = 1. \quad (7)$$

In the next step, we use the assumption that the prior  $p(\mathbf{e})$  is uniform. If the  $p(\mathbf{e})$  is uniform, it does not depend on  $\mathbf{e}$ , therefore, we can take it outside the integral too:

$$\frac{p(\mathbf{x})}{p(\mathbf{x}|\mathbf{c})p(\mathbf{e})} \int_{\mathbf{e}} p(\mathbf{e}|\mathbf{x})p(\mathbf{e}|\mathbf{c})d\mathbf{e} = 1. \quad (8)$$

Multiplying on both sides by  $\frac{p(\mathbf{x}|\mathbf{c})p(\mathbf{e})}{p(\mathbf{x})}$ , we obtain that

$$\int_{\mathbf{e}} p(\mathbf{e}|\mathbf{x})p(\mathbf{e}|\mathbf{c})d\mathbf{e} = \frac{p(\mathbf{x}|\mathbf{c})p(\mathbf{e})}{p(\mathbf{x})}. \quad (9)$$

At this point, we can substitute  $\frac{p(\mathbf{x}|\mathbf{c})p(\mathbf{e})}{p(\mathbf{x})}$  in equation 4 with the integral on the left hand side of equation 9, obtaining

$$p(\mathbf{e}|\mathbf{x}, \mathbf{c}) = \frac{p(\mathbf{e}|\mathbf{x})p(\mathbf{e}|\mathbf{c})}{\int_{\mathbf{e}} p(\mathbf{e}|\mathbf{x})p(\mathbf{e}|\mathbf{c})d\mathbf{e}}. \quad (10)$$

Note that  $\int_{\mathbf{e}} p(\mathbf{e}|\mathbf{x})p(\mathbf{e}|\mathbf{c})d\mathbf{e}$  is integrated over all values of  $\mathbf{e}$ , and therefore it is constant with respect to changes in  $\mathbf{e}$ . Thus, in conclusion:

$$p(\mathbf{e}|\mathbf{x}, \mathbf{c}) \propto p(\mathbf{e}|\mathbf{x})p(\mathbf{e}|\mathbf{c}). \quad (11)$$