

Figure 4: The problem described in Example 3, shown in fig. 3 (right), but where we consider the chance of selecting the single best intervention ( $\hat{x} = 1$ ) as a function of prior  $\alpha$  and for different number of observations. The prior is necessary to obtain stable estimates of the expected information gain, however, it will impact the estimate of the expected information gain slightly differently and therefore a large prior may change the ordering. The method selects the best intervention with a much higher probability than chance even for very low counts.

ples have been included, in order to show that the expected gain eventually converges to the mean-field value.

We note that the expected information gain can be overestimated in the small-sample limit (see e.g.  $\hat{y} \neq 1$ ) where the prior term  $\alpha$  will be more important. To gain more insight into this, we consider the same example, but now show the probability that the expected gain will be the largest for the optimal intervention  $\hat{x} = 1$ , i.e. the chance that the agent will actually select the optimal intervention (see fig. 4). This probability is plotted as a function of the regularization parameter  $\alpha$  for three representative numbers of samples. Although the ability to select the optimal intervention is more impacted by the prior in the small-sample limit, we note that even for just 20 samples, it is much higher than chance ( $\frac{1}{8}$ ).

**Example 4: Active learning** This example will consider a concrete Active learning setting and show that Bayesian causal induction can learn quicker when actions are selected based on eq. (7), compared to random selection. We consider the ground-truth as fixed at  $H_0 = h_{X \rightarrow Y}^0$ , and generate larger random problems by setting  $K_X = K_Y = 8$  and setting  $p_{xy} = \frac{u_{xy}}{\sum_{x,y} u_{xy}}$ , where  $u_{xy}$  are i.i.d. uniform random variables in  $[0, 1]$ . Given  $P_{xy}$ , we sample  $n$  as in the previous examples, and compute the information gain in favor of  $h_{X \rightarrow Y}^0$  based on an intervention  $\hat{x}$  or  $\hat{y}$  and corresponding observations of  $y$  and  $x$  using the Bayesian update given in eq. (5). We then consider the case where interventions are selected randomly, as well as the case where they are selected using the maximal expected information gain computed using eq. (7). The results are averaged over  $10^5$  simulations using  $\alpha = 2$ . In both cases, Bayesian causal induction gains information about the true causal orientation, however, the informative action selections result in about twice as large gain in evidence on average.

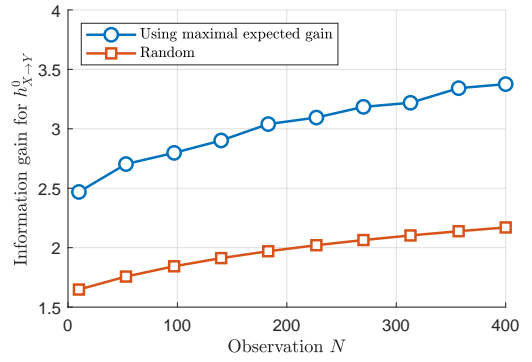


Figure 5: Evaluations of method for an actual intervention-selection problem. The truth is considered fixed as  $H^0 = h_{X \rightarrow Y}^0$ . From this,  $10^5$  random joint distributions  $P(X, Y)$  are generated (see text), and the information gains towards  $h_{X \rightarrow Y}^0$  are computed when interventions are either selected randomly, or when using the maximum anticipated gain eq. (7). As shown, the information gain is in both cases positive, but about twice as large when interventions are selected using our method.

## Discussion and conclusion

Probability trees are conceptually the simplest possible approach to causal inference: consider the causal orientation as an event, specify a prior, and compute the posterior. Although recent work has demonstrated how probability trees can represent both interventions and counterfactual, as well as context-dependent causal statements that cannot be expressed in a directed acyclic causal model (?), their practical use has remained limited.

In this work, we have highlighted another aspect of probability trees and Bayesian causal induction, namely the ability to make predictions about the information gain *prior* to performing interventions. We have illustrated this in the simplest possible situation, and shown how to express both the the expected gain *before* making an intervention, and the realized gain *after* making an intervention.

In making statements such as this, it is important to emphasize that a reduction in uncertainty is by itself meaningless, rather, what matters is the information gain in favor of the true hypothesis, a distinction we make using  $H^0$  and  $H^a$ .

In experiments, we have shown these measures can quantify the information gain and distinguish between different interventions. An active-learning example (fig. 5) shows an increased information gain when our method is used to select optimal interventions.

Many interesting avenues remain unexplored, such as the generalization to larger graphs, and concrete concentration bounds on the expected and realized information gains in terms of  $P$ .

Ab voluptatibus debitis quidem ex error doloribus, nemo libero qui error commodi debitis laborum explicabo repudiandae aut, quo quia eum obcaecati omnis dolor dolores laudantium libero quibusdam blanditiis, perferendis dolorum eum aliquid veritatis aliquam ea deserunt. Corrupti teneatur atque quis voluptatum pariatur a recusandae, cum ratione

molestiae commodi facilis porro minima repudiandae consequatur accusamus cupiditate.Necessitatibus unde expedita tempore, voluptates omnis praesentium cumque