

Generic Priors Yield Competition Between Independently-Occurring Preventive Causes

Derek Powell¹ (derekpowell@ucla.edu)

M. Alice Merrick¹ (m.a.merrick@gmail.com)

Hongjing Lu^{1,2} (hongjing@ucla.edu)

Keith J. Holyoak¹ (holyoak@lifesci.ucla.edu)

Departments of Psychology¹ and Statistics², University of California, Los Angeles
Los Angeles, CA, USA

Abstract

Recent work on causal learning has investigated the possible role of generic priors in guiding human judgments of causal strength. One proposal has been that people have a preference for causes that are *sparse and strong*—i.e., few in number and individually strong (Lu et al., 2008). Sparse-and-strong priors predict that competition can be observed between candidate causes of the same polarity (i.e., generative or else preventive) even if they occur independently. For instance, the strength of a moderately strong cause should be underestimated when a strong cause is also present, relative to when a weaker cause is present. In previous work (Powell et al., 2013) we found such competition effects for causal setups involving multiple generative causes. Here we investigate whether analogous competition is found for strength judgments about multiple *preventive* causes. An experiment revealed that a cue competition effect is indeed observed for preventive causes; moreover, the effect appears to be more persistent (as the number of observations increases) than the corresponding effect observed for generative causes. These findings, which are consistent with predictions of a Bayesian learning model with sparse-and-strong priors, provide further evidence that a preference for parsimony guides inferences about causal strength.

Keywords: causal learning; generic priors; causal strength; parsimony; Bayesian modeling

Introduction

Prior Beliefs in Causal Learning

Humans (and other intelligent organisms) are able to extract causal knowledge from patterns of covariation among cues and outcomes (for a review see Holyoak & Cheng, 2011). This knowledge enables us to predict the likely effects of interventions, and to better anticipate the outcomes of actions and events in the world. Causal knowledge goes beyond mere strength of association, instead depending on inferences about underlying causal links. Because causal relationships are not themselves observable, their existence must be posited *a priori*. Thus a certain degree of prior knowledge about causal relationships is required before causal learning can occur.

Viewed from a Bayesian perspective, causal inferences are expected to be a joint function of likelihoods (the probability of observing the data given potential causal links of various possible strengths) and priors (expectations about causal links that the learner brings to the task). An understanding of how causal links operate and act together is encoded in the likelihood term, whereas prior knowledge about what sorts of causal relationships are most likely is encoded in the prior term. *Generic* causal priors constitute preferences for certain types of causal explanations, based on abstract properties rather than domain-specific knowledge.

Although some Bayesian models have assumed uninformative priors (e.g., Griffiths & Tenenbaum, 2005), other models have incorporated substantive generic priors about the nature of causes. Lu et al. (2008) proposed that people have a preference for causes that are *sparse and strong*: i.e., a preference for causal models that include a relatively small number of strong causes (rather than a larger number of weak causes). Using an iterative-learning method, Yeung and Griffiths (2011) empirically derived a different (though non-uniform) prior that was suggestive of a preference for strong causes, but that lacked the competitive pattern associated with the sparse prior. However, in their study the iterative method did not fully converge for the background cause; hence their results are open to multiple interpretations. Sparse-and-strong priors can be viewed as a special case of a more general pressure to encourage parsimony (Chater & Vitanyi, 2003), which implies a combination of simplicity and explanatory power (see also Novick & Cheng, 2004; Lombrozo, 2007).

Generic Prior: Sparse-and-Strong (SS) Causes

Lu et al. (2008) formalized the “SS power” model with sparse-and-strong (SS) priors for simple causal models with a single candidate cue and a constantly-present background cause. By default (in accord with the power PC theory of causal learning; see Cheng, 1997), the background cause is assumed to be generative (making the effect happen). When

the candidate cause also generates (rather than prevents) the effect, SS priors create an expectation that the candidate cause will be strong (strength close to 1) and the background weak (strength close to 0), or vice versa. A single free parameter, α , controls the impact of the prior (when $\alpha = 0$, the prior distribution of strength is uniform). Lu et al. (2008) fit several causal learning models to parametric data for human strength judgments. They found the best fit was provided by a Bayesian implementation of the power PC theory that incorporated SS priors with an α value of 5 (not 0), implying a human preference for sparse-and-strong causes.

The original formulation only considered the simple case of competition between one generative candidate cause and the (generative) background cause. More recently, Powell et al. (2013) generalized the model to situations involving multiple generative candidate causes. The qualitative signature of SS priors is the preference for one strong cause of a given polarity (generative or else preventive) coupled with other weak causes; for example, a set of “ideal” causal strengths for three generative causes might be $w_A=1$, $w_B=0$ and $w_C=0$ (where each w indicates a value of causal strength, ranging from 0 to 1, for a particular candidate cause). Note that the SS prior does not express any preference about precisely *which* cause(s) are strong and which are weak.

This preference instantiated in SS priors implies a key empirical prediction: in judgments of causal strength, competition will be observed when multiple causes of the same polarity co-occur. That is, if a candidate cause appears along with one or more other causes of greater strength (as defined by likelihoods), then the strength of the weaker candidate cause will be underestimated.

Powell et al. (2013) tested the prediction of SS priors using a causal setup with three generative candidate causes. The objective power of one cause (A) was fixed at a moderate value of .50. The power of a second generative cause (B) was varied across conditions to be either low (.20) or high (.80). A third cause (C) served as a “visible” background cause, with power always set to a low value (.10). The contingencies were arranged so that the occurrences of cues A and B were uncorrelated with each other. Neither A nor B ever appeared alone (whereas C sometimes did), creating uncertainty in estimating the strengths of A and B. After participants had made a relatively small number of observations, cue A was judged to be weaker when the alternative (but uncorrelated) cause B was strong than when it was weak. After additional cases were presented, the two conditions converged. Powell et al. showed that a sparse-and-strong prior accurately characterized the expectations that would create the observed competition effect. This competition dynamic cannot be explained by naïve Bayesian models that assume uninformative priors (Busemeyer et al., 1993).

The Case of Preventive Causes

In the present paper we consider whether or not similar competition effects can be observed in situations involving multiple *preventive* causes. To the best of our knowledge, no previous study has investigated this possibility. The original formulation of SS priors (Lu et al., 2008) only considered the case of a single candidate cause coupled with the background cause. By default the background cause is assumed to be generative (since unless something generates the effect, it is impossible to assess whether other cues prevent it). Hence, in the case of a single candidate cause, SS priors predict competition effects for a generative candidate (which will compete with the generative background cause), but not for a preventive candidate (which will not compete with the generative background, given the basic assumption that only causes of the same polarity compete). In accord with these predictions, Lu et al. reported experiments showing competition effects for strength judgments about a generative candidate cause, but not a preventive cause.

Preventive causes are known to exhibit several asymmetries relative to generative causes (Carroll, Cheng & Lu, 2013; Cheng et al., 2007), so it is not obvious whether or not competition effects would be observed for preventive causes in more complex setups. However, the general principle of SS priors implies that multiple preventive causes will also compete. Testing this hypothesis requires creating a more complex, multi-cue causal setup, similar to the three-cause setup Powell et al. (2013) used to show competition among multiple generative causes. In the present paper we extend the Bayesian formulation of SS priors to the case of multiple preventive causes, and report an experiment that tests whether or not multiple preventive causes compete in causal strength judgments.

Generating Contingency Data

As reviewed above, the sparse-and-strong prior predicts competition between co-occurring causes of the same polarity. We constructed a set of contingency data D (summarized in Table 1) based on the presence of one generative cause (C) and the occurrence or non-occurrence of two preventive causes (A and B). Two conditions were created. The causal powers of A and C were held constant across the two conditions (preventive at -.50 and generative

Table 1. *Contingency learning data for one experimental block (40 trials) by trial type*

Conditions		C	AC	BC	ABC
Weak-B	E present	8	4	6	3
	E absent	2	6	4	7
Strong-B	E present	8	4	2	1
	E absent	2	6	8	9

at .80, respectively, coding preventive causes as negative strength values). The causal power of B varied from one condition (-.25, weak-B condition) to the other (-.75, strong-B condition). The occurrences of causes A and B were independent in both conditions.

In constructing these contingency data, we employed the noisy-OR likelihood function for generative causes, and the noisy-AND-NOT function for preventive causes (in accord with the power PC theory; Cheng, 1997), since binary causes and effects were used in the experiments. For example, if cues A and B are preventive causes, and cue C is a generative cause, the probability of the effect can be calculated as:

$$P(E = 1 | w_A, w_B, w_C) = w_C(1 - w_A)(1 - w_B) \quad (1)$$

where $E = 1$ indicates the presence of the effect.

Sparse-and-Strong Prior for Multiple Preventive Causes

The sparse-and-strong prior for multiple generative and preventive causes can be extended using mixture distributions as proposed in Lu et al. (2008). The key idea is that causes with the same polarity compete. For example, in the situation involving cue C as a generative cause and two other cues, A and B, as preventive causes, the sparse-and-strong prior can be defined as:

$$P(w_A, w_B, w_C) \propto e^{-\alpha(1-w_C)}(e^{-\alpha w_A - \alpha(1-w_B)} + e^{-\alpha(1-w_A) - \alpha w_B}). \quad (2)$$

Importantly, participants in our experiment were not informed of the polarity of the individual causes they observed. That is, prior to observing the contingency data, they did not know which cause(s) generated the effect and which prevented it. As discussed by Lu et al. (2008), the SS power model can be extended to cover situations in which the causal polarity of cues is not known in advance. Since human participants observed a C-only condition, we assume for simplicity that the polarity of cause C (generative) is known (as the likelihood that C is a preventer rapidly falls to zero in light of the contingency data). To predict participants' responses in this experimental context, we therefore derived sparse-and-strong prior distributions for four combinations of generative and preventive A and B cues. The prior probability was computed over a domain of strengths for both A and B in the range [-1,1]. Density in positive regions represents probability density for generative strengths, whereas density in negative regions represents probability density for preventive strengths. Equation 3 gives the sparse-and-strong prior functions for each quadrant of the domain. Quadrants are labeled with a three-letter code specifying the polarity of each cause ("P" for preventive and "G" for generative) in the order A, B, C. For instance, the region representing A and B as preventers and C as a generator is labeled "PPG".

Thus, we derived the following sparse-and-strong prior distribution over the four quadrants of the space depending on causal polarities for cues A and B:

$$P(w_A, w_B, w_C) \propto \begin{cases} e^{-\alpha(1-w_C)}(e^{-\alpha w_A - \alpha(1-w_B)} + e^{-\alpha(1-w_A) - \alpha w_B}), & PPG \\ (e^{-\alpha(1-w_B) - \alpha w_C} + e^{-\alpha w_B - \alpha(1-w_C)})(e^{-\alpha w_A} + e^{-\alpha(1-w_A)}), & PGG \\ (e^{-\alpha(1-w_A) - \alpha w_C} + e^{-\alpha w_A - \alpha(1-w_C)})(e^{-\alpha w_B} + e^{-\alpha(1-w_B)}), & GPG \\ e^{-\alpha(1-w_A) - \alpha w_B - \alpha w_C} + e^{-\alpha w_A - \alpha(1-w_B) - \alpha w_C} + e^{-\alpha w_A - \alpha w_B - \alpha(1-w_C)}, & GGG \end{cases} \quad (3)$$

Similarly, the likelihood function was specified over each of the four regions:

$$P(D | w_A, w_B, w_C) = \begin{cases} w_C(1 - w_A)(1 - w_B), & PPG \\ (1 - (1 - w_B)(1 - w_C))(1 - w_A), & PGG \\ (1 - (1 - w_A)(1 - w_C))(1 - w_B), & GPG \\ 1 - (1 - w_A)(1 - w_B)(1 - w_C), & GGG \end{cases} \quad (4)$$

Finally, the posterior was calculated based on Bayes rule as

$$P(w_A, w_B, w_C | D) = \frac{P(D | w_A, w_B, w_C)P(w_A, w_B, w_C)}{P(D)}. \quad (5)$$

To compare with human ratings, the SS power model computes the mean of estimated causal strength derived from the posterior distribution:

$$\bar{w}_A = \int_0^1 w_A P(w_A | D). \quad (6)$$

The posterior distribution $P(w_A | D)$ is obtained by marginalizing the posterior probability distribution calculated in Equation 5. The means of causal strengths estimated from the posterior distributions were used to compare with ratings from human participants, as shown in Figure 2 and Table 2.

The model makes two key predictions. First, the preventive cause A should be underestimated when in the presence of a stronger preventer, relative to when it is accompanied by a weaker preventer. Notably, this competition effect is predicted despite the fact that A and B occur independently. In contrast, a model assuming uniform priors predicts no such competition (see Figure 2). Second, this competition effect should diminish across learning trials. Priors tend to be most informative when data is scarce. As the model is exposed to more data, the likelihood term should eventually swamp out the competition driven by the prior term. However, as shown in Figure 2, the predicted impact of the priors persists even after three blocks of contingency data (the maximum data that participants received in our experiment). When compared to the corresponding predictions for three generative causes (Powell et al., 2013), the competition effect is predicted to be more persistent for the preventive case (potentially due to the greater uncertainty for causes of unknown polarities).

Experiment

Method

Participants Participants were 107 undergraduate students at the University of California, Los Angeles, who

participated for class credit (29 male, mean age = 21 years). Of these, 49 were assigned to the strong-B condition and 58 to the weak-B condition.

Procedure Participants read the following cover story: “Imagine that you are assisting a doctor at a new island resort. Many of the guests at this new resort have become ill, and your job is to determine the cause of the illnesses. The resort’s doctor suspects that the illnesses may be caused by the food the guests are eating. The resort provides a complimentary salad bar that offers a selection of exotic vegetables. The guests’ salads can be different mixes of these vegetables. One or more of the vegetables could be making the guests sick. However, not all of the guests who ate salad are getting sick. One or more of the vegetables may have medicinal properties that prevent the illness. Your job is to independently examine the doctor’s theory and determine what effects these different vegetables have on the guests. These are the three vegetables [pictures of vegetables are shown]. You will be reviewing a number of case files that describe what a guest ate and whether they became sick. Please pay attention to each case.... At several points throughout the study, you will be asked to give your assessment of the vegetables. You will be asked to assess whether each vegetable causes or prevents the illness. You will also be asked to estimate how likely each vegetable is to make someone sick or prevent the illness.”

These vegetables were labeled A, B, and C, and were

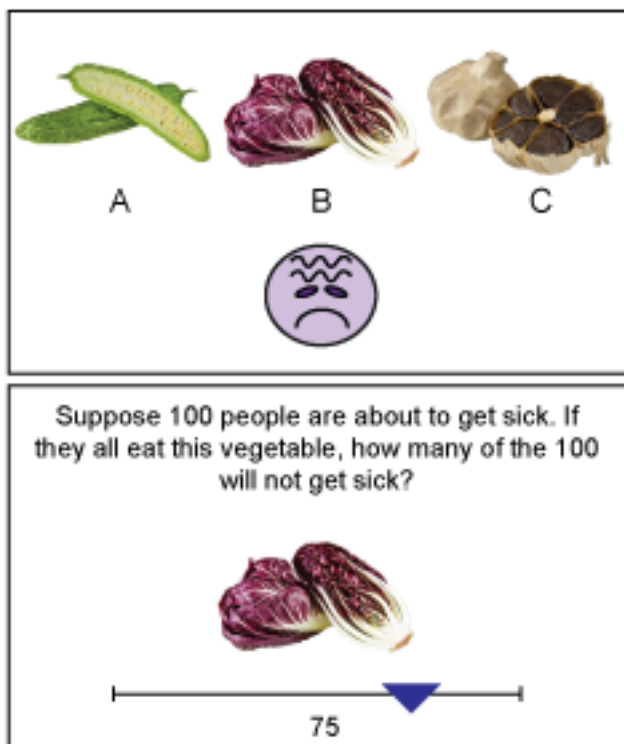


Figure 1: Example trial showing sick guest who ate A, B, and C vegetables (top). Example response trial for preventive cause (bottom).

represented by photographs of actual exotic vegetables: radicchio, bitter melon, and black garlic. The assignment of vegetables to the labels A, B and C was randomized across participants. During the learning phase, participants viewed “case files” for individual guests, showing which combination of vegetables they had eaten, and whether or not they had fallen ill (see Figure 1, top, for an example):

There were four possible combinations of vegetables: each guest had either eaten vegetable C alone, vegetables A and C, vegetables B and C, or all three vegetables A, B, and C. These four combinations were presented in equal number, such that A and B both occurred 50 percent of the time, and the correlation between the occurrence of A and B was 0 (see Table 1). Forty cases (10 of each type) were used per block, as this was the minimum number required to reflect the underlying causal powers in the presented distribution of cause combinations and their associated outcomes.

The numbers of guests who became sick after eating each combination were determined by the causal powers assigned to each vegetable. This number was calculated according to the noisy-OR and noisy-AND-NOT likelihood functions, under the default assumption that all causes act independently of one another (Cheng, 1997). In both conditions, vegetable C was the generative cause of the illness present on every trial, acting with a causal strength of .80. Vegetables A and B were always preventive causes, but the strength of vegetable B varied depending on the condition participants were assigned to. Participants were randomly assigned to one of two experimental conditions (weak-B or strong-B). In the strong-B condition, vegetable B was assigned a preventive causal strength of $-.75$, whereas in the weak-B condition, vegetable B was assigned a causal strength of $-.25$. In each case, vegetable A prevented the illness with strength $-.50$. Cause A was the focus of the study, as we were interested in whether its judged strength would be influenced by the variation in the strength of cause B. Participants were not told ahead of time which vegetables were generative causes and which were preventive causes.

Participants viewed three blocks of 40 trials each. In each block, the 40 cases were presented sequentially in a different random order for each participant. After completing a learning block, participants were asked questions about each vegetable. For each vegetable, participants were first asked a causal polarity question. Participants were shown a picture of each vegetable along with a question asking whether it caused or prevented the illness. Participants pressed “C” to indicate that they thought the vegetable caused the illness, and “P” to indicate that they thought it prevented the illness.

After they indicated which type of cause they believed a vegetable was, they advanced to the causal strength question. If they believed the vegetable was a generative cause, they were asked, “Suppose 100 people ate this

vegetable, how many will get sick?” If they believed the vegetable was a preventer they were asked, “Suppose 100 people are about to get sick. If they all eat this vegetable, how many of the 100 will not get sick?” Participants then made their rating using a slider, inputting a value between 0 and 100 (see Figure 1, bottom). The polarity question always preceded the strength question, but the order of the three vegetables was randomized for each participant and block. After responding to all questions in a block, participants were shown a summary of their responses and were asked to confirm that they had correctly entered their ratings.

Results

Six participants mistakenly indicated that cause C was a preventer during at least one of the three blocks. As it is unclear how their other responses are to be interpreted in light of this error, data from these participants were excluded from analyses. Participants’ polarity judgments and strength ratings for each vegetable were combined into a single index. When participants indicated that a cause was generative, their strength rating was recorded as their score. When they indicated that a cause was preventive, their score was computed by multiplying this strength rating by -1 .

The data for all cues and conditions were compared with the SS power model (described above). For modeling purposes we simply set $\alpha = 5$ (the value estimated for the data sets reported by Lu et al., 2008), thus avoiding any need to fit a free parameter to the present data. We also compared the human data to an otherwise identical model assuming uniform priors (i.e., $\alpha = 0$). Figure 2 shows the data for the critical A cue, along with the predictions derived from the two models. Participants in the strong-B condition underestimated the strength of cause A relative to

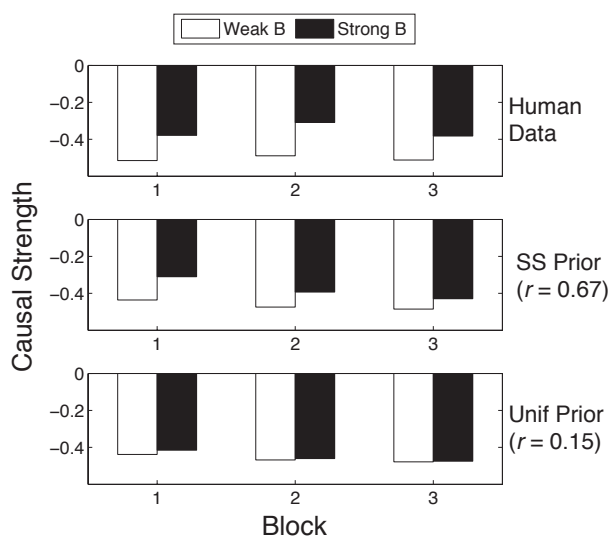


Figure 2: Human causal strength judgments for cause A and predictions from models using sparse-and-strong (SS) and uniform (Unif) priors across blocks.

participants in the weak-B condition, $F(1, 99) = 4.89, p = .029$. There was no significant main effect of block ($F(2, 198) = 0.83, p = .437$) for A ratings, nor was there a significant interaction between block and condition ($F(2, 198) = 0.21, p = .813$).

Qualitatively, the SS power model predicts the observed difference in the judged strength of A in the weak-B versus strong-B conditions more accurately ($r = .67$) than does the model with uniform priors ($r = .15$). However, the SS prior model also predicts that this competition effect will diminish (though not disappear) as participants view a second and third block of trials. The lack of a significant interaction indicates that this second prediction was not confirmed. Rather, the competition effect observed in the human data remained statistically constant across the three blocks.

Additional analyses were performed on responses to all three causal cues. Table 2 presents the mean ratings of causal strength obtained for the three different cues in each condition, along with the predictions derived from the two models based on alternative priors. Analyses of participants’ ratings for cause B revealed a somewhat puzzling pattern of results. As expected, an ANOVA revealed a strong effect of condition, with participants rating cause B as a stronger preventer in the strong-B condition than in the weak-B condition, $F(1, 99) = 31.90, p < .001$. However, this analysis also revealed a significant main effect across blocks ($F(2, 198) = 12.31, p < .01$), as well as a significant interaction between block and condition ($F(2, 198) = 11.68, p < .01$).

As is apparent in Table 2, in the strong-B condition participants’ ratings remained stable across blocks, but in the weak-B condition participants rated cause B as a weaker preventer in later blocks than in earlier blocks. Neither model predicts this trend. The fact that in the weak-B condition the estimated strength of B approached zero across blocks (rather than its objective value of $-.25$) suggests that some participants may have decided at some point that the weaker preventer (i.e., cue B) was simply not a cause at all, and then ceased to attend to the cue. Further, the effect of this error could then be magnified due to a demand pressure to give non-zero responses when presented with the strength question. (Across all vegetables, less than 7% of strength ratings were below an absolute value of 25.)

Discussion

The present experiment provides support for the generality of a parsimony constraint on estimates of causal strength. Previous demonstrations of cue competition between independently-occurring causes only involved generative causes. Our study found similar competition effects among preventive causes, as predicted by a Bayesian model of causal learning that assumes sparse-and-strong priors. A preventive cause of moderate strength was judged to be weaker when a competing (but uncorrelated) preventive cause was strong than when the competing cause was weak.

Table 2. Observed strength ratings for human participants and predicted values based on SS priors and uniform priors. Positive values represent generative strength ratings; negative values represent preventive strength ratings. Individual columns present ratings and estimates from each learning block (1, 2, or 3) for each cause (A, B, or C).

		Cause A			Cause B			Cause C		
		Block			Block			Block		
		1	2	3	1	2	3	1	2	3
Human Data	Weak-B	-.52	-.49	-.51	-.33	-.13	-.03	.72	.75	.71
	Strong-B	-.38	-.31	-.38	-.58	-.59	-.57	.73	.70	.68
SS Prior	Weak-B	-.44	-.47	-.49	-.13	-.16	-.18	.75	.77	.77
	Strong-B	-.31	-.39	-.43	-.73	-.75	-.75	.75	.77	.78
Uniform Prior	Weak-B	-.44	-.47	-.48	-.13	-.19	-.21	.72	.75	.77
	Strong-B	-.42	-.46	-.47	-.67	-.71	-.72	.72	.76	.77

This competition dynamic cannot be explained by naïve Bayesian models that assume uninformative priors (Busemeyer et al., 1993).

The competition effect we observed between two preventive causes was more persistent than that observed in a previous study using a setup involving similar contingencies among generative causes (Powell et al., 2013). For generative causes, competition was observed only in the first block (of 44 trials); for preventive causes, it persisted across three blocks of 40 trials each. The model with SS priors predicted that the competition effect for preventive causes would persist but diminish across three blocks.

The greater persistence of generic priors in the preventive case might reflect the greater cognitive load imposed by a causal model in which polarity of multiple causes is uncertain. Perhaps some participants coped with this complexity by performing model selection (Griffiths & Tenenbaum, 2005) early on, thereby simplifying their causal model. In particular, if the B cause was sometimes dropped from the model in the weak-B condition (i.e., participants decided weak B was not a cause at all), this could account for why competition was persistent (A would always have an advantage if B was dropped, thus favoring A in the weak-B condition relative to the strong-B condition), and also why the judged strength of the weak B cue tended towards 0, rather than its true value of $-.25$. This hypothesis could be tested by adding an explicit non-causal response option for each cause (e.g., “vegetable B is not a cause”). Our findings thus highlight the need for further work on causal learning in complex, multi-causal setups with mixed causal polarities.

References

Busemeyer, J. R., Myung, I. J., & McDaniel, M. A. (1993). Cue competition effects: Theoretical implications for

- adaptive network learning models. *Psychological Science*, 4, 196–202.
- Carroll, C., Cheng, P. W., & Lu, H. (2013). Inferential dependencies in causal inference: A comparison of belief-distribution and associative approaches. *Journal of Experimental Psychology: General*, 142, 845-863.
- Chater, N., & Vitányi, P. (2003). Simplicity: A unifying principle in cognitive science? *Trends in cognitive sciences*, 7, 19-22.
- Cheng, P. W. (1997). From covariation to causation: A causal power theory. *Psychological Review*, 104, 367-405.
- Cheng, P.W., Novick, L.R., Liljeholm, M. & Ford, C. (2007). Explaining four psychological asymmetries in causal reasoning: Implications of causal assumptions for coherence. In M. Rourke (Ed.), *Topics in contemporary philosophy (Vol. 4): Explanation and causation* (pp. 1-32). Cambridge, MA: MIT Press.
- Holyoak, K. J., & Cheng, P. W. (2011). Causal learning and inference as a rational process: The new synthesis. *Annual Review of Psychology*, 62, 135–63.
- Griffiths, T. L., & Tenenbaum, J. B. (2005). Structure and strength in causal induction. *Cognitive Psychology*, 51, 334–384.
- Lombrozo, T. (2007). Simplicity and probability in causal explanation. *Cognitive Psychology*, 55, 232–57.
- Lu, H., Yuille, A. L., Liljeholm, M., Cheng, P. W., & Holyoak, K. J. (2008). Bayesian generic priors for causal learning. *Psychological Review*, 115, 955-984.
- Novick, L. R., & Cheng, P. W. (2004). Assessing interactive causal influence. *Psychological Review*, 111, 455–85.
- Powell, D., Merrick, M. A., Lu, H., & Holyoak, K. J. (2013). Generic priors yield competition between independently-occurring causes. In M. Knauf, M. Pauven, N. Sebanz, & I. Wachsmuth (Eds.), *Proceedings of the 35th Annual Conference of the Cognitive Science Society* (pp. 1157-1162). Austin, TX: Cognitive Science Society.

Yeung, S., & Griffiths, T. L. (2011). Estimating human priors on causal strength. In L. Carlson, C. Hoelscher, & T. F. Shipley (Eds.), *Proceedings of the 33rd Annual Conference of the Cognitive Science Society* (pp. 1709-1714). Austin, TX: Cognitive Science Society.