

# Learning Causal Direction from Transitions with Continuous and Noisy Variables

Kevin W. Soo (kws10@pitt.edu)

Benjamin M. Rottman (rottman@pitt.edu)

Department of Psychology, University of Pittsburgh  
3939 O’Hara Street, Pittsburgh, PA 15260

## Abstract

Previous work has found that one way people infer the direction of causal relationships involves identifying an asymmetry in how causes and effects change over time. In the current research we test the generalizability of this reasoning strategy in more complex environments involving ordinal and continuous variables and with noise. Participants were still able to use the strategy with ordinal and continuous variables. However, when noise made it difficult to identify the asymmetry participants were no longer able to infer the causal direction.

**Keywords:** causal reasoning; causal structure; time

## Introduction

Knowing the direction of causal relations is critical for an agent to be able to act effectively in the world. For example, knowing that James gossips to Julie but not vice versa would allow one to selectively share a secret. Knowing that energy prices influence the price of produce (due to shipping costs) but not vice versa would allow one to predict the prices of energy and produce after an oil embargo or after a poor harvest.

There are four known strategies that lay people use to infer causal direction between two variables, X and Y. The first strategy is *intervention*. Imagine that one intervenes and sets the value of X to either 1 or 0. If the distribution of Y is different when X is set to 1 vs. 0, it implies that X causes Y (e.g. Steyvers et al., 2003)

The remaining three strategies apply to “observational” data, when a learner observes the states of X and Y. Inferring causal direction from observational data is notoriously tricky, and these strategies are usually viewed as heuristics. The next strategy is *causal sufficiency*. In some situations learners have a strong belief that when a cause is present its effect would also be present. In this case, observing  $[x = 0, y = 1]$  is consistent with inferring  $X \rightarrow Y$  but not  $X \leftarrow Y$  (Mayrhofer & Waldmann, 2011; cf. Deverett & Kemp, 2012). Another strategy is to reason using *temporal “delay” or “order”*. If one observes that X usually occurs before Y (e.g., X changes from 0 to 1, and subsequently Y changes from 0 to 1, or Event X occurs and then Event Y occurs) people infer that  $X \rightarrow Y$  instead of  $X \leftarrow Y$  (e.g. Lagnado & Sloman, 2006).

## Learning Causal Direction from Transitions

Rottman & Keil (2012) proposed that people also learn causal direction in a fourth way, by reasoning about changes or “transitions” in variables over time. For example, Rottman & Keil asked participants to observe the moods of

two friends, Friend X and Friend Y over time, and to try to infer whether Friend X’s mood influences Friend Y’s mood, or vice versa. The X and Y columns in Table 1 show the kind of data presented to participants in their experiments. Each person could either be in a positive (1) or negative (0) mood. In order to explain this reasoning process, it is easiest to first think about how a particular causal structure tends to produce certain types of transitions given “shocks” to X or Y. A shock is an economic term that means an exogenous event that produces a change in an observed variable.

Table 1: Transition-based Learning with Binary Variables. “S(X)” means a shock to X.

Time	X	Y	Transition Type	Is transition consistent with causal structure?	
				X→Y	X←Y
1	1	1			
2	0	0	$\Delta X \Delta Y$	Yes: S(X)	Yes: S(Y)
3	0	1	$\Delta Y$	Yes: S(Y)	No: S(Y) would change X
4	0	1	No $\Delta$	Yes: No S	Yes: No S

Assume that  $X \rightarrow Y$  is the true causal structure. When something produces a change in X’s mood (a shock to X), the change in X’s mood will generally produce a change in Y’s mood. Transitions in which both X and Y change are labeled as “ $\Delta X \Delta Y$ ”. Now consider what would happen given a shock that changes Person Y’s mood. If there is reason to believe that the states of X and Y are “temporally dependent” or autocorrelated, that people’s moods are fairly stable over time, then when there is a change in Y’s mood X’s mood would usually remain stable (i.e. in the same state as it was before). Transitions in which there is a change in Y but no change in X are labeled as “ $\Delta Y$ ”. Finally, transitions in which both X and Y stay the same are labeled “No  $\Delta$ ” for no change. No  $\Delta$  transitions are common under both  $X \rightarrow Y$  and  $X \leftarrow Y$  assuming that X and Y are autocorrelated.

This logic implies that if  $X \rightarrow Y$  is the true causal direction, then  $\Delta X \Delta Y$  and  $\Delta Y$  transitions would be common; however,  $\Delta X$  transitions would be rare because if X’s mood changes then there would typically be a change in Y as well ( $\Delta X \Delta Y$  instead of  $\Delta X$ ). Alternatively, if the true causal structure is  $X \leftarrow Y$ , then  $\Delta X \Delta Y$  and  $\Delta X$  would be common, but  $\Delta Y$  would be rare. Given this logic, one can reason backwards to infer the direction of the causal relation. For example, in Table 1, there is a  $\Delta X \Delta Y$  and a  $\Delta Y$  transition but no  $\Delta X$  transitions, implying that  $X \rightarrow Y$  is more likely to be the true causal structure than  $X \leftarrow Y$ . The

reason is that during the  $\Delta Y$  transition there must have been a shock that changed the value of Y, but this shock did not carry over to X, implying that Y does not influence X. Both adults and children learn causal direction using this inference process (Rottman & Keil 2012; Rottman, Kominsky & Keil, 2013).

### Transition-based learning with continuous and noisy variables

In the previous example with binary variables it may seem possible that people would spontaneously categorize the transitions as  $\Delta X$ ,  $\Delta Y$ ,  $\Delta X\Delta Y$ , or No  $\Delta$ . However, real world data (e.g., stock market graphs) contain a number of features that may make the categorization and interpretation process harder, which could make it much more difficult to infer the causal direction. Here we focus on two features, variables with multiple levels opposed to binary, and noise.

Table 2: Examples of Transitions in Experiments 1-3.

Time	Transition Type	Exp 1: No Noise		Exp 2: Noise in $\Delta X\Delta Y$		Exp. 3: Noise in all transitions	
		X	Y	X	Y	X	Y
1		6	4	6	4	6.3	4.5
2	$\Delta X\Delta Y$	4	2	4	3	4.4	3.2
3	$\Delta Y$	4	7	4	7	4.5	7.6
4	No $\Delta$	4	7	4	7	4.3	7.7

Table 2 displays sample data for three experiments that involve progressively more complex environments (noisy and multi-level variables). Next we describe how these environments may obscure the transition types.

**No  $\Delta$**  First, consider the No  $\Delta$  transitions from Times 3 to 4. Experiments 1 and 2 are similar to the binary case in that variables stay in exactly the same state as in the previous trial. In Experiment 3, there is a bit of noise added so that even when there is no shock (No  $\Delta$ ) there will still be slight changes in both X and Y (X changes by -0.2 and Y by +0.1). However, we still consider this a No  $\Delta$  transition because these variations are small and attributed to ‘noise’ in the variables. We will contrast this with ‘shocks’ in the other transition types.

**$\Delta X\Delta Y$**  Next, consider the  $\Delta X\Delta Y$  transition from Times 1 to 2 when there is a shock to X that carries over to Y. In the binary case, transitions where both variables change in the same direction always result in both variables in the same state, implying a positive relationship. In Experiment 1, during  $\Delta X\Delta Y$  transitions X and Y change the same amount (e.g., -2). But because their starting and finishing states are not necessarily the same (e.g., X changes from 6 to 4, Y changes from 4 to 2) it may not be as obvious evidence for a causal relation or as easy to detect as the binary change.

In  $\Delta X\Delta Y$  transitions in Experiment 2, X and Y do not necessarily change the same amount (e.g., from Time 1-2, X decreases by 2 and Y decreases by 1). Presumably, this would make it even less obvious that there is a causal relation between X and Y. Experiment 3 takes this further

using continuous variables (X changes by -1.9, Y by -1.4). If the noise applied to the variables results in  $\Delta X$  and  $\Delta Y$  being sufficiently different in magnitude, that transition may not be easily categorized as a  $\Delta X\Delta Y$  transition.

**$\Delta Y$**  Finally, consider the  $\Delta Y$  transitions (Time 2-3), when there is a shock to Y that does not carry over to X. The crux of transition-based learning is that people notice an asymmetry in the number of  $\Delta X$  and  $\Delta Y$  transitions, and this asymmetry implies the causal direction. In Experiments 1 and 2, X stays exactly the same as in the previous trial, which is similar to the binary case. This should highlight the fact that the change in Y does not have any effect on X.

Experiment 3 adds some noise into X during  $\Delta Y$  transitions (from Time 2-3, X increases by 0.1). We still call this a  $\Delta Y$  transition even though X changes because the change in X due to noise will usually be smaller than the change in Y due to a shock. Being able to attribute small changes in X as due to noise vs. large changes in Y as due to a ‘shock’ is critical for being able to identify  $\Delta Y$  as  $\Delta Y$  as opposed to  $\Delta X\Delta Y$ . However, if the difference between the magnitudes of shocks and noise is sufficiently small on a particular transition, that transition may not be easily categorized as  $\Delta Y$ .

In sum, Experiments 1-3 test how progressively more complex environments make the transitions more difficult to categorize and disrupt the ability to infer causal direction.

### Experiment 1

Experiment 1 tests whether people are able to learn causal direction by observing a pair of ordinal variables over time. This experiment is similar to the binary case because 1) when both variables change,  $\Delta X\Delta Y$ , they change by the same amount ( $\Delta X = \Delta Y$ ), 2) during  $\Delta Y$  transitions X remains the same as the previous trial, and 3) during No  $\Delta$  transitions X and Y stay the same as the previous trial. These aspects of Experiment 1 (relative to the following experiments) should make it easier for reasoners to categorize each transition into the four transition types..

However, ordinal variables add an additional layer of complexity into the categorization and interpretation process compared to binary variables (Table 1). For example, during  $\Delta X\Delta Y$  transitions one could focus on how the variables change with different magnitudes (e.g., sometimes by +3 or by -2). Categorizing all of these transitions as  $\Delta X\Delta Y$  may not be quite so automatic as with binary variables. The same argument applies in  $\Delta Y$  transitions. During one  $\Delta Y$  transition Y may change by +1, and in another Y may change by -6. It is not obvious that people would spontaneously interpret both these transitions in the same way. During all transitions, even during No  $\Delta$  transitions, one could focus on the absolute magnitude at a given time rather than the change in the magnitudes. For example, if X is high and Y is low during a No  $\Delta$  transition, it could imply that they are not causally related. In sum, the ordinal variables introduce a number of layers of complexity into the reasoning process.

Our main question of interest is whether people are able to learn the causal direction in longitudinal sequences of data with ordinal values similar to Table 2 Experiment 1. We call this condition “temporally dependent” because X and Y often remain in the same state from trial to trial. However, in creating the sequences of learning data, since there are  $\Delta Y$  transitions but no  $\Delta X$  transitions, in the long run, Y has higher variance than X. Restated, there is a confound in that there are two differences between X and Y. One difference is that there are  $\Delta Y$  but there are no  $\Delta X$  transitions, and the other is that the variance of Y is often greater than the variance of X.

In order to ensure that the causal direction inference is due to the asymmetry in  $\Delta X$  and  $\Delta Y$ , we also ran a “temporally independent” condition using the exact same trials as the dependent condition, but with a randomized trial order. The randomization destroys the temporal dependence, so participants would not be able to infer causal direction from transitions. Table 3 shows a sample comparing the temporally dependent and independent conditions. Randomizing the trials results in a very high number of  $\Delta X \Delta Y$  transitions, thus there is no longer an asymmetry between  $\Delta X$  and  $\Delta Y$ , so it is impossible to infer causal direction according to the transitions. However, in both conditions Y still has higher variance than X. If the difference in variances explains the effect then it would still occur in the temporally independent condition.

Table 3: Example of Corresponding Temporally Dependent and Independent Conditions in Experiments 1 and 2.

Time	Temporally dependent		Temporally independent			
	Transition	X	Y	Transition	X	Y
1		7	5		5	6
2	No $\Delta$	7	5	$\Delta X \Delta Y$	7	5
3	$\Delta X \Delta Y$	5	3	$\Delta X \Delta Y$	5	3
4	$\Delta Y$	5	6	$\Delta X \Delta Y$	7	5

## Methods

**Participants** 103 participants were recruited using Amazon Mechanical Turk (MTurk) and the experiment was conducted online. The experiment lasted between 10-15 minutes and participants were paid \$1.50 for participation. We intended to recruit 100 participants, but three participants started the study and subsequently returned the HIT before completion; we analyzed all the data.

**Stimuli and Design** Data was presented to participants in the form of graphs depicting the states of two variables on a 1-9 scale for a period of 25 time points (resulting in 24 transitions). In the temporally dependent condition there were 12 No  $\Delta$  transitions, 6  $\Delta X \Delta Y$  transitions, and 6  $\Delta Y$  transitions, randomly ordered. See Figure 1 for a sample graph of the temporally dependent condition.

We generated 1000 temporally dependent graphs in the following way. The initial states for X and Y were each sampled from a normal distribution ( $M = 5.0$ ,  $SD = 1.5$ ), rounded to the closest integer. Every subsequent trial was

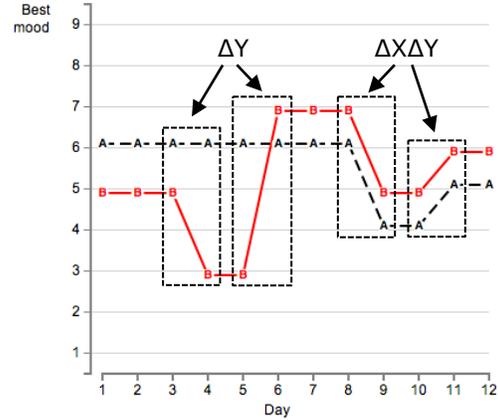


Figure 1: Sample stimuli for Experiment 1. Note: actual stimuli spanned 25 time points. In this graph, A = X (cause) and B = Y (effect). The  $\Delta X \Delta Y$  and  $\Delta Y$  transitions are marked for easy identification above, but were not denoted in any special way in the experiment.

determined in the following way. If the transition was supposed to be a  $\Delta X \Delta Y$  transition (a shock to X), then a new state for X was sampled from the same rounded normal distribution as above. If the sampled value of X happened to be the same as its prior state, a new sample was drawn. The new value for Y was determined by the change in X;  $Y_1 = X_1 - X_0 + Y_0$ . For the  $\Delta Y$  transitions (a shock to Y that does not carry over to X), the new value for Y was chosen from the same rounded normal distribution as above, with the same caveat that  $Y_1$  could not equal  $Y_0$ . X stayed exactly the same as in the previous trial.

We wanted to present data on a 1-9 point scale. However, the generative process explained above sometimes produced data outside the bounds of 1-9. For example, on a  $\Delta X \Delta Y$  transition, if X changed from 4 to 7, and Y started at 7, it would change to 10. Because we did not want to test reasoning about ceiling effects in this experiment, we eliminated data where X or Y exceeded the bounds of 1-9.

For all 1000 temporally dependent sequences of data we created a corresponding temporally independent sequence by randomizing the order of the trials. Each participant in the temporally dependent condition viewed a block of eight graphs chosen randomly from the 1000 sequences. For each participant in the temporally dependent condition, another participant in the temporally independent condition received the 8 corresponding graphs.

**Procedure** Participants read the following cover story:

“Please imagine you are a psychologist studying the moods of married couples. You are trying to figure out how one spouse's mood influences the other. You might find that Spouse A's mood may influence Spouse B's mood, that Spouse B's mood may influence Spouse A's mood, or that neither influences the other. You will observe the moods of married couples, each over a period of 25 consecutive days. On each day, please consider possible events that influenced their moods. For example, Spouse A may have had a good day at work, which put

them into a good mood and spread to Spouse B. Or Spouse B may have had a bad day at work, but their bad mood did not spread to Spouse A. Please remember that moods influence one another on the same day. For example, if Spouse A is in a bad mood on Monday, and if Spouse A's mood affects Spouse B's mood, then Spouse B's mood will also be affected on Monday.”

Participants were randomly assigned (between-subjects) to view graphs of the temporally dependent or temporally independent data. At the start of the scenario, the graph showed the data points for Spouse A and Spouse B visible for Day 1. Participants clicked a button to gradually reveal each subsequent day with a 1 second delay between clicks. The delay was to encourage participants to reason about the data sequentially, rather than clicking through the graph and retrospectively inspecting the data. After the entire graph was revealed, participants were prompted to respond on a scale of 1 (“Confident that Spouse A’s mood influenced Spouse B’s mood”) to 9 (“Confident that Spouse B’s mood influenced Spouse A’s mood”). The midpoint 5 was labeled “Not sure about the direction of the relationship”.

Each participant worked with 8 graphs and made 8 judgments. For each graph, Spouse A and B were randomly assigned to the roles of X (which we predict will be interpreted as the cause) or Y (which we predict will be interpreted as the effect).

## Results

We recoded the inferences so that 9 meant a participant strongly inferred that  $X \rightarrow Y$  and 1 meant that participants strongly inferred  $X \leftarrow Y$ ; 5 meant that a participant did not infer a causal direction. For each participant we took the average of his or her eight judgments. The transition-based learning hypothesis predicts that participants will in general tend to infer  $X \rightarrow Y$ ; so their scores should be greater than 5. Indeed, participants in the temporally dependent condition had scores that were, on average, significantly higher than 5 ( $M = 5.83$ ,  $SE = .20$ ),  $t(49) = 4.10$ ,  $p < .001$ .

An unequal variances t-test also revealed that participants in the temporally dependent condition scored higher than the temporally independent condition ( $M = 4.72$ ,  $SE = .11$ ),  $t(77.5) = 4.75$ ,  $p < .001$ . This ensures that the inferences of causal direction were not due to differences in the variance of X and Y but rather the transitions in the data.

## Experiment 2

Experiment 1 was intended to be as similar as possible to the previous experiments with binary variables. Experiment 2 takes the next step in testing the generalizability of transition-based learning by adding noise into the transitions. In Experiment 2, during  $\Delta X \Delta Y$  transitions, X and Y sometimes changed the same amount and sometimes by different amounts.

## Methods

**Participants** A new group of 129 participants were recruited on MTurk. We intended to recruit 100 participants; however, due to a programming error after the main part of the study was over, 29 participants returned the HIT even though they completed the study. We analyzed all the data. The experiment lasted between 10-15 minutes and participants were paid \$1.50.

**Stimuli and procedure** Data was presented to participants in the same form as Experiment 1. The transition characteristics were the same as in Experiment 1, except in the following ways. First, during the  $\Delta X \Delta Y$  transitions, the change in Y would be within  $\pm 1$  of the change in X. For example, if X changed by +3, Y could change by +2, +3, or +4, with equal probability. Thus, during  $\Delta X \Delta Y$  transitions X and Y always changed in the same direction. Additionally, during  $\Delta X \Delta Y$  transitions X was required to change by at least 2 points in either direction. We eliminated the possibility of X changing by only 1 point, because then it would be possible for X to change by +1 and Y to not change at all (i.e., the noise in Y could cancel out the change in X). As in Experiment 1, 1000 graphs were generated and for each graph a temporally independent version was created by randomizing the trials. All other stimuli and procedural characteristics were similar to Experiment 1.

## Results

As in Experiment 1, participants in the temporally dependent condition continued to infer that  $X \rightarrow Y$  was more likely than  $X \leftarrow Y$  ( $M = 5.74$ ,  $SE = .185$ ),  $t(64) = 3.99$ ,  $p < .001$ . Participants were more likely to infer  $X \rightarrow Y$  in the temporally dependent than independent condition ( $M = 4.79$ ,  $SE = .11$ ),  $t(106.34) = 4.35$ ,  $p < .001$ .

## Experiment 3

The purpose of Experiment 3 was to test whether people continue to infer causal direction with even more noise in the transitions so that it was not always clear what type of transition had occurred. Similar to Experiment 2, during  $\Delta X \Delta Y$  transitions X and Y change to varying degrees. In addition, now during No  $\Delta$  transitions both X and Y change slightly from their previous states. Additionally, during  $\Delta Y$  transitions (a shock to Y), X changes slightly from its previous state. In sum, X and Y are still highly (but not perfectly) dependent on their previous states.

The added noise meant that it was much harder to determine what type of transition had occurred ( $\Delta X \Delta Y$ ,  $\Delta Y$  or No  $\Delta$ ). Or equivalently, it was much harder to identify if any change in a variable was due to a shock that transferred, a shock that did not transfer, or noise (with no shock).

First, consider what would happen with a shock to X,  $S(X)$ . If  $S(X)$  is large, then it would also produce a large change in Y leading to a clear  $\Delta X \Delta Y$  transition. However, suppose there is a small shock increasing X while the noise to Y is negative, cancelling out the increase in Y that would have been produced by  $S(X)$ . In this case it would appear as

if a  $\Delta X$  transition had occurred. Because the asymmetry between  $\Delta X$  vs.  $\Delta Y$  is critical to inferring causal direction, such transitions would make it harder to identify  $X \rightarrow Y$  as the true direction.

Next, consider what happens when there are shocks to  $Y$ . With a large  $S(Y)$ ,  $Y$  would change significantly more than  $X$  (which changes a fairly small amount due to noise) leading to a clear  $\Delta Y$  transition. This can be seen in the first transition (from Time 1-2) in Figure 2. However, suppose that  $S(Y)$  produces a small increase in  $Y$  and the noise to  $X$  is positive. This would make it appear that a  $\Delta X \Delta Y$  transition had occurred because  $X$  and  $Y$  both increase by similar magnitudes (see the final transition from Time 11-12 in Figure 2). If a shock to  $Y$  appears as  $\Delta X \Delta Y$  instead of  $\Delta Y$  it would make inferring causal direction harder because the asymmetry between  $\Delta X$  and  $\Delta Y$  is reduced.

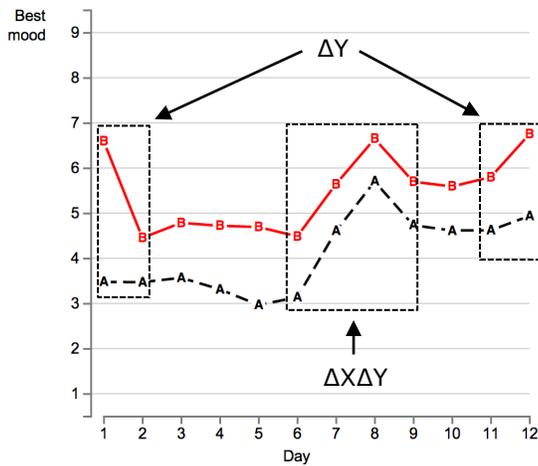


Figure 2: Sample stimuli for Experiment 3, in which shocks do not need to be at least  $\pm 2$ . Note: true stimuli spanned 25 time points. In this graph,  $A = X$  and  $B = Y$ .  $\Delta X \Delta Y$  and  $\Delta Y$  transitions are marked by the boxes, No  $\Delta$  transitions are unmarked (no markings were present in the experiment).

Finally, consider transitions without any shocks. Without any noise these transitions would be interpreted as No  $\Delta$  transitions. However, when there is noise they would not necessarily be interpreted as No  $\Delta$ . If either  $X$  or  $Y$  changes significantly more than the other the transition could be interpreted as  $\Delta X$  or  $\Delta Y$ . If  $X$  and  $Y$  both happen to change a fairly large amount in the same direction, it might appear that a  $\Delta X \Delta Y$  transition has occurred.

In sum, when trying to determine what kind of transition just occurred participants needed to distinguish shocks from noise. We propose that larger differences in the magnitudes of shocks vs. noise will enable participants to better distinguish the transition types that enable inferences of causal direction. To test this, we created two conditions. In the “big shock” condition, the shocks to  $X$  and the shocks to  $Y$  were at least 2 points in either direction, similar to Experiment 2, which ensured that their magnitudes were appreciably greater than the variations expected from noise. In the “any size shock” condition, there was no minimum

magnitude for shocks. Sometimes small shocks occurred which could be confused with noise, making it more difficult to track what kinds of transitions had occurred and infer a causal direction. This manipulation sought to investigate the boundary conditions of transition-based learning by asking three questions: 1) Would people still be able to infer causal direction in the “big shock condition”? This condition is considerably harder than Experiments 1 and 2 which involved considerably less noise. 2) Is it harder to infer causal direction in the “any size shock” condition than the “big shock” condition? 3) Are people still able to infer causal direction in the “any size shock” condition?

## Methods

**Participants** 100 new participants were recruited on MTurk but 3 returned the HIT; we analyzed data for all 103 participants. Similar to Experiments 1 and 2, participants were paid \$1.50 for the 10-15 minute experiment.

**Stimuli.** Data were presented to participants in graphs similar to Experiments 1 and 2, but with several differences in how they were generated. Firstly, the states of variables were continuous (there was no rounding to the closest integer). Secondly, on each transition noise was introduced to variables that did not undergo shocks. Noise, the change from the prior state, was sampled from a normal distribution ( $M = 0, SD = 0.2$ ). This ‘noise distribution’ generally led to smaller changes than those produced by shocks to  $X$  and  $Y$ .

For No  $\Delta$  transitions, the new states  $X_1$  and  $Y_1$  were each equal to their prior states plus noise sampled from the noise distribution (i.e.  $X_1 = X_0 + \text{noise}$ ,  $Y_1 = Y_0 + \text{noise}$ )

For  $\Delta X \Delta Y$  transitions, a new state was sampled for  $X$  (a shock to  $X$ ). In the “big shock” condition, the magnitude of the change was always greater than 2 in either direction. In practice, the state was sampled from the normal distribution ( $M = 5, SD = 1.5$ ), and resampled if the change was less than 2. In the “any size shock” condition, this requirement was dropped, so the shock to  $X$  could result in a change of less than 2. The new state of  $Y$  was equal to the change in  $X$  plus noise sampled from the noise distribution.

For  $\Delta Y$  transitions, a new state for  $Y$  was sampled from the ( $M = 5, SD = 1.5$ ) normal distribution. In the “big shock” condition, similar to the  $\Delta X \Delta Y$  transitions, the new state had to differ from the prior state by at least 2. This requirement was dropped for the “any size shock” condition”. The new state of  $X$  would be equal to its prior state plus noise from the noise distribution (i.e.  $X_0 + \text{noise}$ ).

Figure 2 shows an example from the “any size shock” condition. The first  $\Delta Y$  transition (Time 1-2) shows a shock to  $Y$  that produces a change of roughly -2. The second  $\Delta Y$  transition (Time 11-12) shows a shock to  $Y$  that produces a change of roughly +1. Because  $X$  happens to increase slightly at the same time the transition from Time 11-12 might appear as a  $\Delta X \Delta Y$  transition. In the “big shock” condition, all shocks would appear similar to the former.

Participants were randomly assigned (between-subjects) to either the “big shock” or “any shock” condition and

worked with eight graphs. Experiment 3's procedures were otherwise identical to the previous experiments.

## Results

Participants in the “big shock” condition were still able to infer causal direction above chance; the scores were higher than 5 on average ( $M = 5.34$ ,  $SE = .19$ ),  $t(50) = 1.77$ ,  $p < .05$ . However, participants in the in the “any size shock” condition did not score significantly higher than 5 ( $M = 5.05$ ,  $SE = .15$ ),  $t(51) = .363$ ,  $p > .05$ . An unequal variances  $t$ -test comparing the two conditions was not significant,  $t(93.467) = 1.19$ ,  $p > .05$ . In sum, with big shocks relative to noise participants were still able to infer causal direction; however, compared to previous experiments this ability was reduced.

## General Discussion

Previous work has shown that people infer causal direction by how two variables, X and Y change over time. In series where Y changes but X does not (a ‘shock’ to Y but stable X), people tend to infer that Y does not influence X. If there are other transitions in which both X and Y change simultaneously, people tend to infer that X causes Y. However, the previous research only used simple cases with binary variables (Rottman & Keil, 2012).

Across three experiments we systematically made the inference more complicated by using ordinal and continuous variables and adding noise in how X and Y change. Participants were still able to use this transition-based learning process for inferring causal direction with ordinal variables with no noise (Experiment 1) and with ordinal variables with some noise in the ‘shocks’ (Experiment 2). Experiment 3 involved noise in both ‘shocks’ and in ‘stable’ periods. When the shocks were large and could be distinguished from the stable periods people still inferred causal direction above chance, but when the shocks were small and could be confused with the stable periods people were not able to reliably infer the causal direction.

## Implications and Questions for Causal Inference

This research raises a number of implications and open questions for causal inference. Relatively little causal learning research has investigated ordinal and continuous variables, especially for inferring causal direction. In the current research the explanation of how people inferred causal direction relied upon categorizing transitions as  $\Delta X$ ,  $\Delta Y$ ,  $\Delta X\Delta Y$ , or No  $\Delta$ , reducing continuous stimuli to categorical transition types. It will be important to validate the use of these categories, to understand if and how people spontaneously categorize transitions in these types, and to understand more specifically how noise makes the categorization process harder. In addition, it will be important to understand how people infer the direction of positive vs. negative causal relations; with noise in a positive relation X and Y can sometimes change in opposite directions.

Another question is whether the general type of learning process in this manuscript can be viewed as a normative or rational inference. There is a growing field of machine learning and econometrics devoted to identifying causal direction from time-series data (e.g., Moneta & Spirtes, 2006; also see Rottman & Keil, 2012).

More broadly, this research underscores that people are highly sensitive to transitions in how variables change over time, not just the states of the variables. This was demonstrated by the difference between the temporally dependent vs. independent conditions in Experiments 1 and 2 as well as the explanation for how people learn causal direction. Many theories of causal inference focus exclusively on the states of variables within individual trials, not how the variables change over time, thus ignoring an important part of causal learning.

In conclusion, being able to infer causal direction is a critical ability that allows us to effectively navigate and manipulate the world. There are a variety of ways that people learn causal direction. Manipulation and observing a delay between a cause and effect are very strong cues to causal direction. But presumably there are a host of other more subtle cues that people use in environments in which manipulation is not possible and delays are not present. Uncovering when various cues to causal direction are used and the boundaries of these cues will be vital for a fuller understanding of the breadth of human causal inference.

## References

- Deverett, B., & Kemp, C. (2012). Learning Deterministic Causal Networks from Observational Data. *Proceedings of the 34th Annual Conference of the Cognitive Science Society*.
- Mayrhofer, R., & Waldmann, M. R. (2011). Heuristics in Covariation-based Induction of Causal Models: Sufficiency and Necessity Priors. *Proceedings of the 33rd Annual Conference of the Cognitive Science Society*.
- Moneta, A., & Spirtes, P. (2006). Graphical Models for the Identification of Causal Structures in Multivariate Time Series Models. *Proceedings of the 9th Joint Conference on Information Sciences, 1*, 1–4.
- Lagnado, D. A., & Sloman, S. A. (2006). Time as a guide to cause. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *32*(3), 451–60.
- Rottman, B. M., & Keil, F. C. (2012). Causal structure learning over time: observations and interventions. *Cognitive Psychology*, *64*(1-2), 93–125.
- Steyvers, M., Tenenbaum, J. B., Wagenmakers, E.-J., & Blum, B. (2003). Inferring causal networks from observations and interventions. *Cognitive Science*, *27*(3), 453–489.