

How to Make a Proceedings Paper Submission

Anonymous CogSci submission

Abstract

Keywords:

Introduction

Identifying the causes of an event is a very important ability that people largely rely on in common life. Even if reasoning about causes doesn't generally pose problems and comes up with clear answers about what caused an event, it is sometimes much more difficult to disambiguate an event and answer such a question. When divergent causal intuitions are involved it can be hard to decide for one candidate cause rather than another. This can lead to very different and serious implications, especially legal ones. For example a well-known story in the literature is the "Irish case", a real-world example of a court case: [include here Irish case]. In such legal contexts determining the main causes of an event is an extremely important task that involves finding who or what is at fault in dramatic accidents, property damage, etc. and can pose a veritable challenge. Here massive fines, jail sentences or death penalties can be the outcomes of the inferred causal history of an event. Reasoning about causes, however, is often implied in much more naive and common circumstances and doesn't necessarily pose any difficulty –e.g. identifying what is the cause of the bottle shattering. But even in very simple situations like that, it is sometimes not so easy to theoretically explain and predict intuitive responses that people give. [Include here Susy vs Bob example]. In all these common-life situations we want to know what the *actual cause* of an event is, in a specific context. That is we don't want here to infer some general causal laws responsible for the occurrence of an event but understand what is the main or most intuitive singular causes in particular cases. The current paper focuses on causal judgements of the latter type, that is on the concept of *actual* and not *general* causation¹, following a distinction that has been frequently made in the literature. Causal judgements have been widely studied and there are existing models of actual causation that aim to capture people's inferences about the causes of a single event. Current research on actual causation mainly relies on two different interpretations of causality: the counterfactual account (CF) or the physical process one (PP).

¹also called "token" or "singular" causation

According to the first interpretation, the general idea is that an event A is said to be a cause of a distinct event B if the occurrence of A makes a difference in the occurrence of B. In other words if B would have been different without A, then A is an actual cause. More precisely A is a cause of B if and only if A and B are true, and A hadn't occurred B wouldn't have occurred. [Example of Suzy only] This type of explanation calls upon counterfactual scenario or possible worlds where the presumed cause of the effect is removed from the system while all other relevant factors are kept as unchanged as possible compared to the actual world. However, a first well-known challenge that the CF account encounters is that it is not clear which counterfactual world we have to take into account in order to justify a singular cause, and how to accurately compare those counterfactuals. Again let's say that Suzy throws a rock at a bottle but this time Bob as well throws a rock at the same bottle with the same precision as Suzy. Suzy's rock reaches the bottle just a second before Bob's and breaks it. According to the CF account if Suzy hadn't thrown the rock, the bottle would have still shattered since Bob's rock would have hit it regardless. So following the general idea mentioned above, Suzy's throw cannot be the cause of the bottle shattering – which is counter-intuitive. This is an example of a general type of cases we call *late pre-emption* and refers to the fact that the final event – i.e. the effect – is "overdetermined" [Hall]. A first common attempt to face the problem is to say that the bottle-shattering event that results from Suzy's throw is qualitatively different from the bottle-shattering event that follows from Bob's throw. In other words A is a cause of B if and only if A and B occur, and A hadn't occurred B wouldn't have occurred at all or would have been delayed. Thus Suzy's throw is indeed the cause of the bottle-shattering event that occurs at the time it did.

[Halpern and Pearl solution: graphical models, static, pb of modelling, without time, etc.] [Then, more general problem with counterfactual dependence].

The physical process account of causation gives a completely different interpretation of causality and states that an event A causes B when there is a physical connection between them. Several theories have been proposed to characterise this physical connection. In one of its latest and probably most convincing formulation, we have to distinguish two things: causal process and causation or causal interaction. A causal

process is a physical process involving an object which conserves a certain quantity across space and time. The conserved quantity is typically a measurable physical property of the object like the mass-energy, etc. and is precisely what scientific theories are meant to formalize. In that sense both Suzy's and Bob's throw are causal processes. A causal interaction or causation is an exchange of that conserved quantity, for example the mass-energy of the rock to the glass of the bottle. According to this definition, only Suzy's throw involves a causal interaction with the bottle. However the PP account of causation doesn't consider as genuine causation some very common cases that people intuitively judge as proper example of causation. For example if I hold the head of my enemy under water and make him die, I'm not genuinely the cause of his death; rather I'm actually preventing the possibility of a genuine causation which is the physical process of breathing oxygen in order to live. The main problem is that, as a consequence, the PP account not only goes sometimes against some of our best causal intuitions, but also needs the counterfactual analysis as well to ground cases of *quasi-causation* by omission or prevention like the latter one.

[?? A third account : probability raising ??]

Theoretical proposition

As a result of the preceding analysis it seems that none of the current main theories of actual causation gives a satisfying explanation of the causal judgements that people actually do. This paper aims at solving the problem by interpreting causality in a completely different way. In the CF framework events are represented as mere propositions and causation is thought as a relation between static states. The PP account is meant to give an epistemological definition of the concept of causation and not an explanation of how people actually rely on temporal informations to make causal judgements. By contrast this paper focuses on the role of temporal information for inferring causal relationship between events qua changes of states over time. In other words we hypothesised that people's causal judgements are influenced by not only the values of all relevant variables at the time the effect occurs, but also the temporal order in which these variables took their values. [?? Add a concrete example ??] This approach has some precedent in the literature but has not been developed carefully and has never been tested experimentally as we intend to do. We add further that people mainly identify causation with a *continuous* sequence of changes of states over time. The underlying intuition is that people trace back the history of changes from the immediate one that directly brought about the occurrence of the effect, up to the root change in the system that initiated the series of changes along the path. As we think that in common life it is really rare to see simultaneously two or more events occurring at the very same time and producing some common effect, we also want to postulate a *no coincidence principle*. Relying on this principle we can split the system's time frame in units of time that are

small enough to have no more than one event² per unit. Let's represent by $\mu_{z_i \rightarrow j}^0(t)$ a change of state at time t , from variable $Z = z_i$ to $Z = z_j$, $\forall i, j \in \mathbb{R}_+$. Here Z is the variable we want to reason about and has no children (represented by the empty set $\{\emptyset\}$). Formalizing our above mentioned idea we first have to find $\mathcal{U}_{y_i \rightarrow j}^Z(t-1)$, that is the set of the immediate changes in the parents Y of the variable Z at time $t-1$ that led to the observed change in Z at time t . So we want to find \mathcal{U} such that $\mathcal{U}_{y_i \rightarrow j}^Z(t-1) \rightarrow \mu_{z_i \rightarrow j}^0(t)$. According to our *no coincidence principle* \mathcal{U} is either empty or a singleton – including only one parent whose value changed. Let's say that at $t-1$ we find a change in a parent variable Y , so $\mathcal{U}_{y_i \rightarrow j}^Z(t-1) = \{\mu_{y_i \rightarrow j}^Y(t-1)\}$. Then we want to find the set of changes such that $\mathcal{U}_{x_i \rightarrow j}^Y(t-2) \rightarrow \mu_{y_i \rightarrow j}^Z(t-1)$, that is the set of immediate previous changes in the parents X of the variable Y at time $t-2$ that led to the observed change in Y at time $t-1$. Following the same logic we suggest that if $\mathcal{U}_{w_i \rightarrow j}^X(t-3) \rightarrow \mu_{x_i \rightarrow j}^Y(t-2)$ is such that $\mathcal{U}_{w_i \rightarrow j}^X(t-3) = \emptyset$, then it means that $\mu_{x_i \rightarrow j}^Y(t-2)$ represents the chronologically first change along the path, occurring in X , and we suggest that this change is identified as being the main cause of Z .

Experiments

To test our hypothesis that the main cause of an effect is the root change that initiates a continuous sequences of changes until the occurrence of the effect, participants were presented animations showing activation spreading over networks of nodes up to the final node. We run three different experiments which shared the same plot that we wanted as intuitive as possible: participants were told that they were working in a nuclear control room and that their job was to monitor networks of particle detectors. In the experiments when a detector, depicted by a square or circle (see below), absorbs a radioactive particle it becomes active and turns black, transmitting the activation across the links of the network so that an active detector activates the next one in the chain and so forth. All the networks included a special component, called 'Gauge of Critical Moment', that becomes active only if all of its input from the detectors it is connected to are active. The Gauge of Critical Moment was always represented by a square with "GCM" or "G" above. At the end of an activation sequence, participants had to click on the detector(s) they considered as the main cause(s) of the activation of the Gauge of Critical Moment.

[Add some stuff... For example why choosing this setting, the underlying idea to make something related to physical intuitions, no high order cognition etc. ??]

Experiment 1

Participants. Paragraph...

Method. Paragraph...

Results. Paragraph...

Discussion. Paragraph...

² Again we insist on the definition of *event* as *change of state*.

Experiment 2

Participants. Paragraph...

Method. Paragraph...

Results. Paragraph...

Discussion. Paragraph...

Experiment 3

Participants. Paragraph...

Method. Paragraph...

Results. Paragraph...

Discussion. Paragraph...

General discussion