# Supplementary Material to "Causal Explanations and XAI"

**Sander Beckers**                                                                SREKCEBREDNAS@GMAIL.COM
*University of Tübingen*

**Editors:** Bernhard Schölkopf, Caroline Uhler and Kun Zhang

## Appendix A.

**Proposition 8** *If $X = x$ is directly sufficient for $Y = y$ then $X = x$ is strongly sufficient for $Y = y$ along some $N$, and if $X = x$ is strongly sufficient for $Y = y$ along some $N$ then $X = x$ is weakly sufficient for $Y = y$.*

**Proof:** Follows directly from the definitions. ▌

**Theorem 12** *If a causal model $M$ that agrees with $h$ satisfies* **Independence** *then the following statements are all equivalent:*

- $X = x$ *is weakly sufficient for $Y = y$ in $M$.*

- $X = x$ *is strongly sufficient for $Y = y$ in $M$.*

- $X = x$ *is directly sufficient for $Y = y$ in $M$.*

**Proof:** The implications from bottom to top are a direct consequence of Proposition 8.

Assume $X = x$ is weakly sufficient for $Y = y$ in $M$. This means that for all $\boldsymbol{u} \in \mathcal{R}(\mathcal{U})$ we have that $(M, \boldsymbol{u}) \models [\boldsymbol{X} \leftarrow \boldsymbol{x}]Y = y$. Let $\boldsymbol{C} = \mathcal{R}(\mathcal{V} - (\boldsymbol{X} \cup \{Y\}))$.

Given that $M$ agrees with $h$, either the equation for $Y$ is of the form $Y = U$ for some $U \in \mathcal{U}$, or $Y$ only has parents in $\mathcal{V} \setminus \{Y\}$. Since the former contradicts our assumption that in all contexts $(M, \boldsymbol{u}) \models [\boldsymbol{X} \leftarrow \boldsymbol{x}]Y = y$, it has to be the latter.

As a consequence, interventions on all endogenous variables make the particular context irrelevant, i.e., for all $\boldsymbol{c} \in \mathcal{R}(\boldsymbol{C})$ and all $\boldsymbol{u_1}, \boldsymbol{u_2} \in \mathcal{R}(\mathcal{U})$, we have that $(M, \boldsymbol{u_1}) \models [\boldsymbol{X} \leftarrow \boldsymbol{x}, \boldsymbol{C} \leftarrow \boldsymbol{c}]Y = y$ iff $(M, \boldsymbol{u_2}) \models [\boldsymbol{X} \leftarrow \boldsymbol{x}, \boldsymbol{C} \leftarrow \boldsymbol{c}]Y = y$.

Further, for each $V_i \in \mathcal{V} \setminus \{Y\}$ the equation is of the form $V_i = U_i$. Although technically one could choose to define $\mathcal{R}(V_i)$ such that $\mathcal{R}(V_i) \not\subseteq \mathcal{R}(U_i)$, this comes down to defining a variable with values that it cannot obtain, which serves no purpose. Therefore we can assume that for each $\boldsymbol{c} \in \boldsymbol{C}$ there exists a context $\boldsymbol{u'}$ such that $(M, \boldsymbol{u'}) \models \boldsymbol{C} = \boldsymbol{c}$. Given that no members of $\boldsymbol{X}$ are parents of members of $\boldsymbol{C}$, it also holds that $(M, \boldsymbol{u'}) \models [\boldsymbol{X} \leftarrow \boldsymbol{x}]\boldsymbol{C} = \boldsymbol{c}$. Since $X = x$ is weakly sufficient for $Y = y$, we also have that $(M, \boldsymbol{u'}) \models [\boldsymbol{X} \leftarrow \boldsymbol{x}]\boldsymbol{C} = \boldsymbol{c} \wedge Y = y$, from which it follows that $(M, \boldsymbol{u'}) \models [\boldsymbol{X} \leftarrow \boldsymbol{x}, \boldsymbol{C} \leftarrow \boldsymbol{c}]Y = y$. Taken together with the previous observation that the particular context is irrelevant, we get that for all $\boldsymbol{c} \in \mathcal{R}(\mathcal{V} - (\boldsymbol{X} \cup \boldsymbol{Y}))$ and all $\boldsymbol{u} \in \mathcal{R}(\mathcal{U})$ we have that $(M, \boldsymbol{u}) \models [\boldsymbol{X} \leftarrow \boldsymbol{x}, \boldsymbol{C} \leftarrow \boldsymbol{c}]Y = y$, which is what we had to prove. ▌

**Theorem 16** *Given a causal setting $(M, \boldsymbol{u})$, the following two statements are equivalent:*

- *there exist $W_1$ and $w_1 \in \mathcal{R}(W_1)$ so that $Y = y$ counterfactually depends on $X = x$ rather than $X = x'$ relative to $W_1 = w_1$ in $(M, u)$.*

- *there exist $W_2$, $w_2 \in \mathcal{R}(W_2)$, and $N$ so that $(X = (x, x'), W_2 = w_2, N)$ is a good counterfactual explanation of $Y = y$.*

**Proof:**

**Observation 1** *Recall from Definition 3 that exogenous variables only appear in equations of the form $V = U$. Say $R \subseteq \mathcal{V}$ are all variables which have such an equation, and call these the* root *variables. It is clear that if we intervene on all of the root variables, they take over the role of the exogenous variables. Concretely, given strong recursivity, for any setting $r \in \mathcal{R}(R)$ there exists a unique setting $v \in \mathcal{R}(\mathcal{V})$ so that for all contexts $u \in \mathcal{R}(\mathcal{U})$ we have that $(M, u) \models [R \leftarrow r]\mathcal{V} = v$.*

Assume that $Y = y$ counterfactually depends on $X = x$ rather than $X = x'$ relative to $W_1 = w_1$ in $(M, u)$. This means that $(M, u) \models X = x \wedge W_1 = w_1 \wedge Y = y$, and $(M, u) \models [X \leftarrow x', W_1 \leftarrow w_1]Y \neq y$.

Let $S = R \setminus (W_1 \cup X)$, and let $s \in \mathcal{R}(S)$ be the unique values so that $(M, u) \models S = s$.

As $R \subseteq (S \cup W_1 \cup X)$, we have that $(X = x, S = s, W_1 = w_1)$ is strongly sufficient for $Y = y$ along $N = \mathcal{V} \setminus (X \cup W_1 \cup S)$, and thus $((X = x, S = s, W_1 = w_1), N)$ is an actual sufficient explanation of $Y = y$.

Furthermore, changing $X$ from $x$ to $x'$ obviously has no effect on any of the other values in $R$. Therefore $(M, u) \models [X \leftarrow x', W_1 \leftarrow w_1]S = s$, and thus we get that $(M, u) \models [X \leftarrow x', W_1 \leftarrow w_1, S \leftarrow s]Y = y'$ for some $y' \neq y$. As before, we can conclude that $(X = x', S = s, W_1 = w_1)$ is strongly sufficient for $Y = y'$ along $N$, and thus $((X = x', S = s, W_1 = w_1), N)$ is a sufficient explanation of $Y = y'$.

Combining the two previous paragraphs, we get that $(X = (x, x'), (S = s, W_1 = w_1)), N)$ is a counterfactual explanation of $Y = y$. Let $(X = (x, x'), W_2 = w_2), N_2)$ be a dominating counterfactual explanation that is not strictly dominated by any other explanation that contains $X = (x, x')$. (It is easy to see that such an explanation must exist: one can simply keep removing elements from $W_1$, $S$, and $N$ until no further element can be removed while still remaining a counterfactual explanation of $Y = y$.)

Now assume that there exist $W_2$, $w_2 \in \mathcal{R}(W_2)$, and $N$ so that $((X = x, W_2 = w_2), N)$ is an actual sufficient explanation of $Y = y$ and $((X = x', W_2 = w_2), N)$ is a sufficient explanation of some $Y = y'$ with $y' \neq y$. Since the first explanation is actual, it follows immediately that $(M, u) \models X = x \wedge W_2 = w_2 \wedge Y = y$. Combining the second explanation with Proposition 8 we get that $(M, u) \models [X \leftarrow x', W_2 = w_2]Y = y'$.

Note that we did not require $X$ to be minimal in either direction, and thus the conditions as stated without minimality of $X$ are equivalent. Therefore the conditions that include the minimality of $X$ are also equivalent, which is what we had to prove. ■

**Proposition 18** *If $((X = x, W = w), N)$ is a sufficient explanation of $Y = y$ and there exists a dominating explanation $((X = x', A = a), B)$ for some values $x'$ and $a \subseteq w$, then $X = x'$ can replace $X = x$.*

**Proof:** Assume $((X = x, W = w), N)$ is a sufficient explanation of $Y = y$ and $((X = x', A = a), B)$ is a dominating explanation of $Y = y$, with $a \subseteq w$. We show that $((X = x', W = w), B)$ is a sufficient explanation of $Y = y$, from which the result follows.

Let $C = \mathcal{V} \setminus (X \cup A \cup B)$. From the definition of sufficient explanations, we know that for all $u \in \mathcal{R}(\mathcal{U})$ and all $c \in \mathcal{R}(C)$, we have that $(M, u) \models [X \leftarrow x', A \leftarrow a, C \leftarrow c]B = b$ for some $b \in \mathcal{R}(B)$ that includes $y$.

Let $D = \mathcal{V} \setminus (X \cup W \cup B)$, $F = W \setminus A$, and let $f$ be the restriction of $w$ to $F$. Note that $C = F \cup D$. From the previous paragraph it follows that for all $u \in \mathcal{R}(\mathcal{U})$ and all $d \in \mathcal{R}(D)$, we have that $(M, u) \models [X \leftarrow x', A \leftarrow a, F \leftarrow f, D \leftarrow d]B = b$, and thus $(M, u) \models [X \leftarrow x', W \leftarrow w, D \leftarrow d]B = b$, which is what had to be shown. ∎

**Theorem 20** *If $X_1 = x_1$ rather than $X_1 = x_1'$ is a counterfactual explanation of $Y = y$ in $(M, u)$ (relative to some $(W = w, N)$) then for some $X_2 \subseteq X_1$, $X_2 = x_2$ rather than $X_2 = x_2'$ is an actual cause of $Y = y$ in $(M, u)$ (where $x_2$ and $x_2'$ are the relevant restrictions to $X_2$).*

**Proof:** Assume $X_1 = x_1$ rather than $X_1 = x_1'$ is a counterfactual explanation of $Y = y$ in $(M, u)$ relative to $(W = w, N)$. This means that $((X_1 = x_1, W = w), N)$ is an actual sufficient explanation of $Y = y$ and $((X_1 = x_1', W = w), N)$ is a sufficient explanation of $Y = y'$ with $y' \neq y$.

Let $(T = t, S)$ be a good sufficient explanation of $Y = y$, i.e., an actual sufficient explanation of $Y = y$ that dominates $((X_1 = x_1, W = w), N)$ and cannot itself be dominated by another actual sufficient explanation of $Y = y$. (It is easy to see that such an explanation must exist: one can simply keep removing elements from $((X_1 = x_1, W = w), N)$ until no further element can be removed while still remaining a sufficient explanation of $Y = y$.)

Let $X_2 = X_1 \cap T$. We now show that $X_2 \neq \emptyset$ by a reductio.

Assume $X_2 = \emptyset$. This means that $T \subseteq W$. Also, $S \subseteq N$. Let $s \in \mathcal{R}(S)$ and $n \in \mathcal{R}(N)$ be the actual values of $S$ and $N$ in $(M, u)$.

Let $C = \mathcal{V} \setminus (X_1 \cup W \cup N)$. Given that $((X_1 = x_1', W = w), N)$ is a sufficient explanation of $Y = y'$, we have that for all $c \in \mathcal{R}(C)$, $(M, u) \models [X_1 \leftarrow x_1', W \leftarrow w, C \leftarrow c]Y = y'$.

Let $D = \mathcal{V} \setminus (T \cup S)$. Given that $(T = t, S)$ is a sufficient explanation of $Y = y$, we have that for all $d \in \mathcal{R}(D)$, $(M, u) \models [T \leftarrow t, D \leftarrow d]Y = y$.

By our assumption, $X_1 \subseteq (W \setminus T)$. Thus $D = X_1 \cup C \cup (W \setminus (T \cup X_1)) \cup (N \setminus S)$. Therefore from the previous paragraph we get that for all $c \in \mathcal{R}(C)$, $(M, u) \models [X_1 \leftarrow x_1', W \leftarrow w, C \leftarrow c]Y = y$. This contradicts the paragraph before the previous, and therefore $X_2 \neq \emptyset$.

Let $W_2 = T \setminus X_2$. Then we can conclude that $(X_2 = x_2, W_2 = w_2, S)$ is a good sufficient explanation of $Y = y$, where $x_2$ is the restriction of $x_1$ to $X_2$, and $w_2$ is the restriction of $W$ to $W_2$.

Remains to be shown that $x_2'$ cannot replace $x_2$ in this explanation, where $x_2'$ is the restriction of $x_1'$ to $X_2$. The former just means that $((X_2 = x_2', W_2 = w_2), S_2)$ is not a sufficient explanation of $Y = y$ for any $S_2 \subseteq S$.

Again we proceed by a reductio: assume that $((X_2 = x_2', W_2 = w_2), S_2)$ is a sufficient explanation of $Y = y$. Let $F = \mathcal{V} \setminus (X_2 \cup W_2 \cup S_2)$. We have that for all $f \in \mathcal{R}(F)$, $(M, u) \models [X_2 \leftarrow x_2', W_2 \leftarrow w_2, F \leftarrow f]Y = y$. In particular, we have that $(M, u) \models [X_1 \leftarrow x_1', W \leftarrow w]Y = y$.

Recall that $((X_1 = x_1', W = w), N)$ is a sufficient explanation of $Y = y'$ with $y' \neq y$. Using Proposition 8, we get that $(M, u) \models [X_1 \leftarrow x_1', W \leftarrow w]Y = y'$. This contradicts the result in the previous paragraph, which concludes the proof. ∎

**Proposition 22** *If $X = x$ is a direct cause of $Y = y$ in $(M, u)$ then there exist values $x'$ such that $X = x$ rather than $X = x'$ is an actual cause of $Y = y$.*

**Proof:** Assume $X = x$ is a direct cause of $Y = y$ in $(M, u)$, i.e., it is part of a direct good sufficient explanation $(X = x, W = w)$ of $Y = y$ in $(M, u)$. This means that all that remains to be shown, is that there exist values $x' \in \mathcal{R}(X)$ such that $(X = x', W = w)$ is not a direct sufficient explanation of $Y = y$.

We know that $(X = x, W = w)$ is directly sufficient for $Y = y$, and this does not hold if we remove any subset from either $X$ or $W$. Let $C = \mathcal{V} \setminus (X \cup W \cup \{Y\})$. Then we have that for all $c \in \mathcal{R}(C)$, $(M, u) \models [X \leftarrow x, W \leftarrow w, C \leftarrow c]Y = y$. From the minimality of $X$, it follows that there exists a $c \in \mathcal{R}(C)$ and $x' \in \mathcal{R}(X)$ so that $(M, u) \models [X \leftarrow x', W \leftarrow w, C \leftarrow c]Y \neq y$. Therefore $(X = x', W = w)$ is not a direct sufficient explanation of $Y = y$. ∎

**Theorem 23** *If a causal model $M$ satisfies **Independence** then for any $W_1, W_2, N_1, N_2$, it holds that there exist values $w_1 \in \mathcal{R}(W_1)$ so that $(X = (x, x'), W_1 = w_1, N_1)$ is a good counterfactual explanation of $Y = y$ iff there exist values $w_2 \in \mathcal{R}(W_2)$ so that $(X = (x, x'), W_2 = w_2, N_2)$ is a good counterfactual explanation of $Y = y$.*

**Proof:** Given the symmetry, it suffices to prove the implication in one direction. We invoke Theorem 16 so that we can use counterfactual dependence instead of good counterfactual explanations.

Assume $Y = y$ counterfactually depends on $X = x$ rather than $X = x'$ relative to $W_1 = w_1$. This means that $X$ is a minimal set such that $(M, u) \models X = x \land W_1 = w_1 \land Y = y$ and $(M, u) \models [X \leftarrow x', W_1 \leftarrow w_1]Y \neq y$. Take any set $W_2 \subseteq (\mathcal{V} \setminus (X \cup \{Y\}))$, and let $w_2$ be the values of $W_2$ in $(M, u)$. First, note that we have $(M, u) \models X = x \land W_2 = w_2 \land Y = y$. Second, since none of the members of $X$ are parents of any of the members of $W_2$, we also have that $(M, u) \models [X \leftarrow x']W_2 = w_2 \land Y = y'$. Thus we also have that $(M, u) \models [X \leftarrow x', W_2 \leftarrow w_2]Y = y'$. As we did not use the fact that $X$ is a minimal set, and thus both statements are equivalent when ignoring minimality on both sides, they are also equivalent when including minimality on both sides, which is what we have to prove. ∎

**Theorem 24** *If a causal model $M$ satisfies **Independence** then the following statements are all equivalent:*

- *$X = x$ is a direct cause of $Y = y$ in $(M, u)$.*

- *there exist values $x'$ so that $X = x$ rather than $X = x'$ is an actual cause of $Y = y$ in $(M, u)$.*

- *$X = x$ is part of a good sufficient explanation of $Y = y$ in $(M, u)$.*

**Proof:** The implication from the first statement to the second is a direct consequence of Proposition 22.

The implication from the second statement to the third follows from the definition of actual cause.

The implication from the third statement to the first follows from Theorem 12, which shows that under **Independence** we may replace sufficiency with direct sufficiency, and thus the result follows from the definition of direct cause. ∎