

Causal Inference with Graphical Models

Jose M. Peña
STIMA, IDA, LiU

Lecture 5: Counterfactuals

Contents

- ▶ Counterfactuals
- ▶ The Twin Network Method
- ▶ Counterfactual Back-Door Criterion
- ▶ Counterfactuals in Linear-Gaussian Causal Models
- ▶ Axiomatic Characterization of Counterfactuals
- ▶ Necessary and Sufficient Causes
- ▶ The Ladder of Causation

Literature

- ▶ Main sources

- ▶ Pearl, J. *Causality: Models, Reasoning, and Inference* (2nd ed.). Cambridge University Press, 2009. Sections 7.1, 7.3, and 9.1-9.3.
- ▶ Pearl, J., Glymour, M. and Jewell, N. P. *Causal Inference in Statistics: A Primer*. Wiley, 2016. Chapter 4.

- ▶ Additional sources

- ▶ Balke, A. and Pearl, J. [Probabilistic Evaluation of Counterfactual Queries](#). *AAAI 94*, 230-237, 1994.
- ▶ Tian, J. and Pearl, J. [Probabilities of Causation: Bounds and Identification](#). *UAI 2000*, 589-598, 2000.

Counterfactuals

- ▶ A causal model consists of a DAG over V (a.k.a. causal structure), a set of functions $x_i = f_i(pa_i, u_i)$ for each $X_i \in V$ (a.k.a. structural equations), and a joint distribution $p(u)$ over the error variables U_i (a.k.a. unobserved causes).
- ▶ A causal model defines a joint distribution over V :

$$p(v) = \prod_i p(x_i | pa_i) = \sum_u [\prod_i \mathbf{1}_{x_i=f(pa_i, u_i)}] p(u) = \sum_{\{u \mid \wedge_i x_i=f(pa_i, u_i)\}} p(u).$$

- ▶ **An assignment $U = u$ represents an individual or unit in the population, or a situation in nature.**
- ▶ Given some variables $X \subseteq V$, the submodel M_x of a causal model M is the model resulting from deleting the edges into the variables in X , and replacing the functions f_i corresponding to the variables in X with the constant functions $X_i = x_i$.
- ▶ The effect of an action $do(x)$ on M is given by M_x .
- ▶ Given two sets of variables $X, Y \subseteq V$, the value of $Y(u)$ in M_x is the potential response of Y to an action $do(x)$ in situation $U = u$, which we denote as $Y_x(u)$.
- ▶ The counterfactual sentence “The value that Y would have obtained in situation $U = u$, had X been x ” is interpreted as $Y_x(u)$.
- ▶ **Population intervention** = $do(x) = Y_x \neq Y_x(u) =$ **unit intervention**.

Counterfactuals

- ▶ Recall that we used counterfactuals to define direct and indirect effects:

$$TE(x, x^*, Y) = E[Y_x - Y_{x^*}]$$

$$CDE(x, x^*, Y) = E[Y_{xz} - Y_{x^*z}]$$

$$NDE(x, x^*, Y) = E[Y_{xz_{x^*}} - Y_{x^*}]$$

$$NIE(x, x^*, Y) = E[Y_{x^*z_x} - Y_{x^*}]$$

- ▶ The expressions above are population-level effects. The unit-level effects are given by the evaluating the expressions under the expectations as functions of U , because **the expectations are over U** .

Counterfactuals

- ▶ The following quantities are well-defined:

$$p(Y = y) = \sum_{\{u|Y(u)=y\}} p(u)$$

$$p(Y_x = y) = \sum_{\{u|Y_x(u)=y\}} p(u)$$

$$p(Y_x = y, X = x') = \sum_{\{u|Y_x(u)=y \wedge X(u)=x'\}} p(u)$$

$$p(Y_x = y, Y_{x'} = y') = \sum_{\{u|Y_x(u)=y \wedge Y_{x'}(u)=y'\}} p(u)$$

despite Y_x and $Y_{x'}$ cannot be observed simultaneously (multiple worlds).

- ▶ The following quantity is particularly interesting:

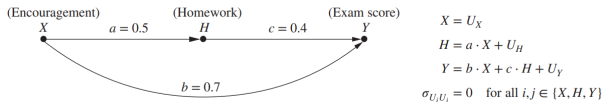
$$p(Y_{x'} = y' | x, y) = \sum_{\{u|Y_{x'}(u)=y'\}} p(u|x, y)$$

because **the probability that $X = x$ caused $Y = y$** may be interpreted as the probability that Y would not be y had X not been x , given that $X = x$ and $Y = y$ have in fact occurred. This is important for diagnosis, planning and determination of liability.

- ▶ Evaluating the expression above can be divided into three steps:
 1. Update $p(u)$ to obtain $p(u|x, y)$ which is assumed to be **invariant** to the hypothetical action in the next step,
 2. modify M to obtain $M_{x'}$, and
 3. compute Y in $M_{x'}$ with $p(u|x, y)$ instead of $p(u)$.

Counterfactuals

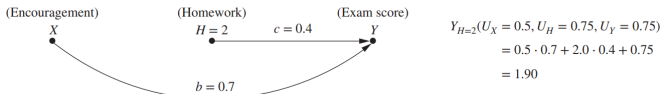
- Consider the following causal model:



- For a student named Joe, we have that $X = 0.5$, $H = 1$ and $Y = 1.5$. What would Joe's score have been, had he doubled his study time?
- Note that the question is about (a subpopulation of) **one individual**, not about the whole population (cf. *do*-calculus). However, the values of a , b and c apply to the population, and U_X , U_H and U_Y account for all variation among the individuals in the population.
- The values of the unobserved variables for Joe (i.e., **his specific characteristics**) are:

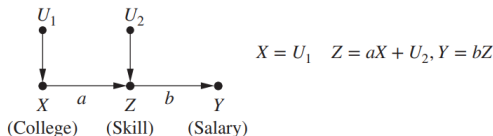
$$\begin{aligned} U_X &= 0.5, \\ U_H &= 1 - 0.5 \cdot 0.5 = 0.75, \text{ and} \\ U_Y &= 1.5 - 0.7 \cdot 0.5 - 0.4 \cdot 1 = 0.75. \end{aligned}$$

- The answer is:



Counterfactuals

- ▶ Consider the following causal model:



where X , U_1 and U_2 are binary, whereas Z and Y are continuous.

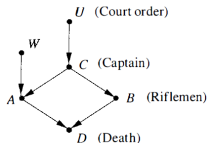
- ▶ What is the expected salary of individuals with skill level $Z = 1$, had they received a college education ?
- ▶ Note that the question is about a **subpopulation**, not about the whole population (cf. *do*-calculus). The values of a and b apply to the population, and U_1 and U_2 account for all variation among the individuals.
- ▶ Assume that $a = 1$. Note that $Z = 1$ occurs only for $(U_1 = 0, U_2 = 1)$ and $(U_1 = 1, U_2 = 0)$ with $p(U_1 = 0)p(U_2 = 1)$ and $p(U_1 = 1)p(U_2 = 0)$. Then:

$$E[Y_{X=1}|Z = 1] = b \left(1 + \frac{P(u_1 = 0)P(u_2 = 1)}{P(u_1 = 0)P(u_2 = 1) + P(u_1 = 1)P(u_2 = 0)} \right)$$

$$E[Y_{X=0}|Z = 1] = b \left(\frac{P(u_1 = 0)P(u_2 = 1)}{P(u_1 = 0)P(u_2 = 1) + P(u_1 = 1)P(u_2 = 0)} \right)$$

which implies that the expected salary of individuals with skill level $Z = 1$ would have been higher had they gone to college than if not. This is to be expected since some of these individuals did not attend college and, had they done it, their skill would have been $Z = 2$ and their salary $Y = 2b$.

Counterfactuals



- ▶ If the prisoner is dead, what is the probability that he/she would be dead even if rifleman A had not shot? That is, we want $p(D_{\bar{a}} = d | D = d)$.
- ▶ Assumptions: The rifleman B is accurate and shoots only if commanded. The rifleman A is accurate and shoots when commanded or due to other reasons, e.g. nervousness (W). Finally, U is unknown.

Model $(M, P(u, w))$

$$\begin{array}{ll}
 C = U & (C) \\
 A = C \vee W & (A) \\
 B = C & (B) \\
 D = A \vee B & (D)
 \end{array}
 \quad (U, W) \sim P(u, w)
 \quad P(u, w) = \begin{cases} pq & \text{if } u = 1, w = 1, \\ p(1 - q) & \text{if } u = 1, w = 0, \\ (1 - p)q & \text{if } u = 0, w = 1, \\ (1 - p)(1 - q) & \text{if } u = 0, w = 0. \end{cases}$$

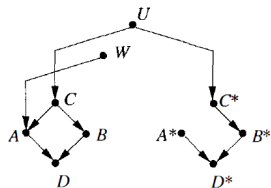
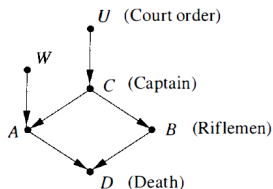
Model $(M_{\neg A}, P(u, w | D))$

$$\begin{array}{ll}
 C = U & (C) \\
 \neg A & (A) \\
 B = C & (B) \\
 D = A \vee B & (D)
 \end{array}
 \quad (U, W) \sim P(u, w | D)
 \quad P(u, w | D) = \begin{cases} \frac{P(u, w)}{1 - (1 - p)(1 - q)} & \text{if } u = 1 \text{ or } w = 1, \\ 0 & \text{if } u = 0 \text{ and } w = 0. \end{cases}$$

$$P(\neg D_{\neg A} | D) = P(\neg U | D) = \frac{q(1 - p)}{1 - (1 - q)(1 - p)}.$$

The Twin Network Method

- ▶ As seen, $p(D_{\bar{a}} = d|D = d)$ can be computed from $p(u, w|d)$ in $M_{\bar{a}}$.
- ▶ Alternatively, it can be computed as $p(D^* = d|D = d, A^* = \bar{a})$ in the twin network. For this, we interpret the twin network as a causal model, which defines a joint distribution over $V \cup V^*$. Marginalization and conditioning do the rest.



- ▶ More efficient methods exist (based on message passing, a.k.a. sum-product algorithm).
- ▶ The advantage of the twin network method is that it avoids computing $p(u|e)$, which may be time and space consuming, because the variables in U may become dependent conditioned on $E = e$.

Counterfactual Back-Door Criterion

- ▶ A set of variables Z satisfies the back-door criterion wrt an ordered pair of sets of variables (X, Y) in a causal structure G which may include unobserved variables if
 - ▶ Z contains no descendants of X , and
 - ▶ Z blocks every path between X and Y that contains an arrow into X .
- ▶ If Z satisfies the back-door criterion wrt (X, Y) , then

$$p(Y_x = y | Z = z, X = x') = p(Y_x = y | Z = z)$$

for all x, x' and z .

- ▶ Proof: Assume to the contrary that there is a path **in the twin network** between X and $Y^* = Y_x$ that is not blocked by Z . The path must be of the form $X \cdots \leftarrow U_S \rightarrow \cdots Y^*$ or $X \cdots \leftarrow U_S \leftrightarrow U_T \rightarrow \cdots Y^*$. So, the original network has a path $X \cdots \leftarrow U_S \rightarrow \cdots Y$ or $X \cdots \leftarrow U_S \leftrightarrow U_T \rightarrow \cdots Y$ that is not blocked by Z . So, the original network has a path between X and Y that contradicts the back-door criterion.
- ▶ The original front-door criterion can also be extended to counterfactuals.

Counterfactual Back-Door Criterion

Example 4.4.1 *A government is funding a job training program aimed at getting jobless people back into the workforce. A pilot randomized experiment shows that the program is effective; a higher percentage of people were hired among those who finished the program than among those who did not go through the program. As a result, the program is approved, and a recruitment effort is launched to encourage enrollment among the unemployed, by offering the job training program to any unemployed person who elects to enroll.*

Lo and behold, enrollment is successful, and the hiring rate among the program's graduates turns out even higher than in the randomized pilot study. The program developers are happy with the results and decide to request additional funding.

- ▶ Critics may say “The program works for randomly chosen people, but those that enroll of their own initiative are more resourceful and, thus, more likely to find a job regardless of the program. One needs to estimate the hiring rate of the enrolled and compare it with what it would have been had they not enrolled”.
- ▶ That is, they want to know the so-called effect of treatment on the treated:

$$ETT = E[Y_{X=1} - Y_{X=0} | X = 1]$$

where X represents enrollment and Y is the hiring rate. Note that

$$ETT = E[Y | X = 1] - E[Y_{X=0} | X = 1]$$

by linearity of the expectation and $E[Y_{X=1} | X = 1] = E[Y | X = 1]$ by consistency of counterfactuals, i.e. **for the individuals that are observed to have $X = 1$, setting $X = 1$ should produce no change.**

Counterfactual Back-Door Criterion

- ▶ If we can find a set of nodes Z that satisfies the back-door criterion wrt (X, Y) , then

$$\begin{aligned} p(Y_x = y | X = x') &= \sum_z p(Y_x = y | X = x', Z = z) p(Z = z | X = x') \\ &= \sum_z p(Y_x = y | Z = z) p(Z = z | X = x') \\ &= \sum_z p(Y_x = y | X = x, Z = z) p(Z = z | X = x') \\ &= \sum_z p(Y = y | X = x, Z = z) p(Z = z | X = x') \end{aligned}$$

by the counterfactual back-door criterion and the consistency of counterfactuals, which implies that

$$E[Y_{X=0} | X = 1] = \sum_z E[Y | X = 0, Z = z] p(Z = z | X = 1).$$

- ▶ Z is obtained by inspecting the causal model. Typically, it is age, education, disposition, etc.

Counterfactual Back-Door Criterion

- ▶ Consider an intervention that adds q mg/l of insulin to a group of patients with varying levels of insulin already in their systems. Can we estimate the effect of this intervention ?
- ▶ Let X denote the current level of insulin, and Y the outcome variable. Note that

$$E[Y|do(x+q)] - E[Y|do(x)]$$

does not answer the question, because it corresponds to giving the same amount of insulin to everyone. Precisely for this reason, *do*-calculus cannot answer queries about **personalized actions**, i.e. queries where the actions depend on the individual.

- ▶ The answer to the question above is

$$\sum_{x'} E[Y_{X=x'+q}|X=x']p(X=x') - E[Y].$$

where $E[Y_{X=x'+q}|X=x']$ can be estimated as seen before letting $x = x' + q$, if we find a set of nodes that satisfy the the back-door criterion wrt (X, Y) . Again, Z is obtained by inspecting the causal model.

Typically, it is age, weight, genetic disposition, etc.

- ▶ A similar problem occurs when we want to estimate the effect of actions on a subpopulation characterized by features that are affected by the actions. In the homework encouragement example, the effect on test score of sending lazy students to the encouragement program is $E[Y_{X=1}|H \leq H_0]$, not $E[Y|do(X=1), H \leq H_0]$.

Counterfactuals in Linear-Gaussian Causal Models

- ▶ In non-parametric causal models, $E[Y_x|z]$ may not be identifiable even if we run experiments: We need the model's functional equations to compute $p(u|z)$ and the desired quantity in M_x .
- ▶ For that reason, in linear-Gaussian causal models, any counterfactual is identifiable whenever the model parameters are identifiable, since these fully define the model. The parameters are identifiable in experimental studies by applying the definition of direct effect. The question is then whether counterfactuals are identifiable in observational studies.
- ▶ If $TE(X, Y)$ is identifiable, then $E[Y_x|z]$ is identifiable as

$$E[Y_x|z] = E[Y|z] + TE(X, Y)(x - E[X|z])$$

for any arbitrary set of nodes Z . In other words, it equals the best prediction of Y given z plus the change expected in Y when X is shifted from the predicted value to the hypothetical one.

- ▶ The result above allows us to compute the ETT in the homework encouragement example:

$$\begin{aligned} ETT &= E[Y_{X=1} - Y_{X=0}|X = 1] = E[Y_{X=1}|X = 1] - E[Y_{X=0}|X = 1] \\ &= E[Y|X = 1] + TE(X, Y)(1 - E[X|X = 1]) \\ &\quad - E[Y|X = 1] - TE(X, Y)(0 - E[X|X = 1]) = TE(X, Y). \end{aligned}$$

In words, the ETT is equal to the effect of treatment on the entire population. This is a general result for linear-Gaussian causal models, i.e. it holds if we replace the evidence $X = 1$ with any $Z = z$.

Axiomatic Characterization of Counterfactuals

- ▶ Composition: $W_x(u) = w \Rightarrow Y_{xw}(u) = Y_x(u)$, i.e. if we force W to take a value that it would have had anyway, then this has no effect on the other variables. Note that **consistency follows from composition**.
- ▶ Effectiveness: $X_{xw}(u) = x$, i.e. if we force X to take value x , then it will take that value.
- ▶ Reversibility: $Y_{xw}(u) = y \wedge W_{xy}(u) \Rightarrow Y_x(u) = y$, i.e. multiple solutions due to feed-back loops are precluded. It follows from composition in causal models without feed-back loops. It may not hold in a model with feed-back loops if the model is too coarse. Adding missing factors to the model restores reversibility.
- ▶ The three **properties** above are **sound** and **complete** for the identification of counterfactuals.
- ▶ To apply the properties above for counterfactual identification, we have first to translate the causal model into the language of counterfactuals using the following two rules (which follow from $y = f_Y(pa_Y, u_Y)$):
 - ▶ Exclusion restrictions: $Y_{pa_Y}(u) = Y_{pa_Y, z}(u)$ for all $Y \in V$ and $Z \subseteq V$.
 - ▶ Independence restrictions: $Y_{pa_Y} \perp \{Z_{pa_{Z^1}}^1, \dots, Z_{pa_{Z^k}}^k\}$ for all $Y \in V$ and $\{Z^1, \dots, Z^k\} \subseteq V$ that is not connected to Y via paths containing only U variables.
- ▶ Unfortunately, unlike for *do*-calculus, there is no algorithm to apply the properties above.

Axiomatic Characterization of Counterfactuals

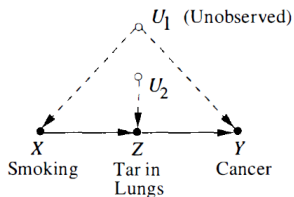


Figure 7.5 Causal diagram illustrating the effect of smoking on lung cancer.

Applying these two rules to our example, we see that the causal diagram in Figure 7.5 encodes the following assumptions:

$$Z_x(u) = Z_{y,x}(u), \quad (7.27)$$

$$X_y(u) = X_{z,y}(u) = X_z(u) = X(u), \quad (7.28)$$

$$Y_z(u) = Y_{z,x}(u), \quad (7.29)$$

$$Z_x \perp\!\!\!\perp \{Y_z, X\}. \quad (7.30)$$

Compute $P(Z_x = z)$ (i.e., the causal effect of smoking on tar).

$$\begin{aligned} P(Z_x = z) &= P(Z_x = z \mid x) \quad \text{from (7.30)} \\ &= P(Z = z \mid x) \quad \text{by composition} \\ &= P(z \mid x). \end{aligned} \quad (7.31)$$

Necessary and Sufficient Causes

- ▶ Let X and Y denote two binary variables, and $x = y = \text{true}$ and $x' = y' = \text{false}$. The probability of **necessity** is

$$PN = p(y'_{x'} | x, y)$$

i.e., the probability that the effect would be absent in the absence of the cause, given that the cause and the effect were present, i.e. how necessary the cause is for the production of the effect. E.g., how necessary the treatment was for no recurrence of the disease.

- ▶ The probability of **sufficiency** is

$$PS = p(y_x | x', y')$$

i.e., the probability that the effect would be present in the presence of the cause, given that the cause and the effect were absent, i.e. how sufficient the cause is for the production of the effect. E.g., how sufficient the treatment would have been for no recurrence of the disease.

- ▶ The probability of **necessity and sufficiency** is

$$PNS = p(y_x, y'_{x'}).$$

- ▶ Moreover, $PNS = p(x, y)PN + p(x', y')PS$. Proof:

$$y_x \wedge y'_{x'} = (y_x \wedge y'_{x'}) \wedge (x \vee x') = (y \wedge x \wedge y'_{x'}) \vee (y_x \wedge x' \wedge y')$$

by consistency of counterfactuals. Finally, take probabilities on both sides.

Necessary and Sufficient Causes

- ▶ In general, we have that

$$\max[0, p(y_x) - p(y_{x'})] \leq PNS \leq \min[p(y_x), p(y_{x'})].$$

- ▶ Under no-confounding (a.k.a. exogeneity), we have that

$$\max[0, p(y|x) - p(y|x')] \leq PNS \leq \min[p(y|x), p(y'|x')]$$

because $p(y_x) = p(y|x)$, and

$$PN = \frac{PNS}{p(y|x)} \text{ and } PS = \frac{PNS}{p(y'|x')}$$

which provide corresponding bounds for PN and PS.

- ▶ **Monotonicity:** $Y_x(u) \geq Y_{x'}(u)$ for all u , i.e. the effect would be present in the presence of the cause, given that the effect is present in the absence of the cause. E.g., the treatment does not produce the disease.
- ▶ Under no-confounding and monotonicity, we have that

$$PNS = p(y|x) - p(y|x')$$

and

$$PN = \frac{p(y|x) - p(y|x')}{p(y|x)} \text{ and } PS = \frac{p(y|x) - p(y|x')}{1 - p(y|x')}.$$

Necessary and Sufficient Causes

- Under just monotonicity, we have that PNS, PN and PS are identifiable whenever $p(y_x)$ and $p(y_{x'})$ are identifiable:

$$PNS = p(y_x) - p(y_{x'})$$

$$PN = \frac{p(y) - p(y_{x'})}{p(x, y)}$$

$$PS = \frac{p(y_x) - p(y)}{p(x', y')}$$

- Moreover,

$$\begin{aligned} PN &= \frac{p(y) - p(y_{x'})}{p(x, y)} = \frac{p(y|x)p(x) + p(y|x')p(x') - p(y_{x'})}{p(y|x)p(x)} \\ &= \frac{p(y|x) - p(y|x')}{p(y|x)} - \frac{p(y|x') - p(y_{x'})}{p(y|x)p(x)} \end{aligned}$$

where the first term is PN under no-confounding, and the second term is a correction for confounding, i.e. $p(y|x') \neq p(y_{x'})$. E.g., the first term measures how much likely it is to recover under treatment, and the second corrects for any confounding between treatment and recovery, e.g. socio-economic factors.

- The expressions above become lower bounds in the non-monotonic case.

Necessary and Sufficient Causes

Consider the following data (Table 9.1, adapted¹⁰ from Finkelstein and Levin 1990) comparing leukemia deaths in children in southern Utah with high and low exposure to radiation from the fallout of nuclear tests in Nevada. Given these data, we wish to estimate the probabilities that high exposure to radiation was a necessary (or sufficient, or both) cause of death due to leukemia.

| | Exposure | |
|--------------------|--------------|--------------|
| | High (x) | Low (x') |
| Deaths (y) | 30 | 16 |
| Survivals (y') | 69,130 | 59,010 |

Assuming monotonicity – that exposure to nuclear radiation had no remedial effect on any individual in the study – the process can be modeled by a simple disjunctive mechanism represented by the equation

$$y = f(x, u, q) = (x \wedge q) \vee u, \quad (9.44)$$

where u represents “all other causes” of y and where q represents all “enabling” mechanisms that must be present for x to trigger y . Assuming that q and u are both unobserved, the question we ask is under what conditions we can identify the probabilities of causation (PNS, PN, and PS) from the joint distribution of X and Y .

Necessary and Sufficient Causes

Assuming no-confounding, i.e. $X \perp \{Q, U\}$:

$$\text{PNS} = P(y | x) - P(y | x') = \frac{30}{30 + 69,130} - \frac{16}{16 + 59,010} = 0.0001625, \quad (9.45)$$

$$\text{PN} = \frac{\text{PNS}}{P(y | x)} = \frac{\text{PNS}}{30/(30 + 69,130)} = 0.37535, \quad (9.46)$$

$$\text{PS} = \frac{\text{PNS}}{1 - P(y | x')} = \frac{\text{PNS}}{1 - 16/(16 + 59,010)} = 0.0001625. \quad (9.47)$$

Statistically, these figures mean that:

1. There is a 1.625 in ten thousand chance that a randomly chosen child would both die of leukemia if exposed and survive if not exposed;
2. There is a 37.544% chance that an exposed child who died from leukemia would have survived had he or she not been exposed;
3. There is a 1.625 in ten thousand chance that any unexposed surviving child would have died of leukemia had he or she been exposed.

Necessary and Sufficient Causes

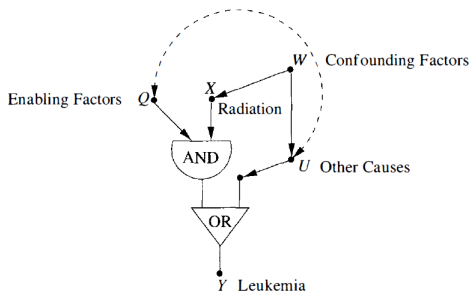
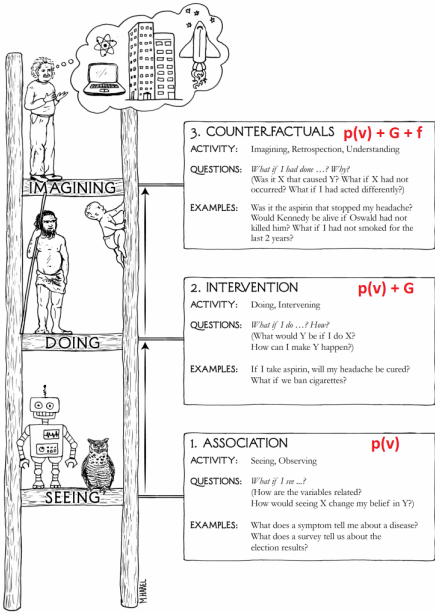


Figure 9.2 Causal relationships in the radiation–leukemia example, where W represents confounding factors.

Finally, Theorem 9.2.15 assures us that PN and PS are identifiable even when x is not independent of $\{u, q\}$, provided only that the mechanism of (9.44) is embedded in a larger causal structure that permits the identification of $P(y_x)$ and $P(y_{x'})$. For example, assume that exposure to nuclear radiation (x) is suspect of being associated with terrain and altitude, which are also factors in determining exposure to cosmic radiation. A model reflecting such consideration is depicted in Figure 9.2, where W represents factors affecting both X and U . A natural way to correct for possible confounding bias in the causal effect of X on Y would be to adjust for W , that is, to calculate $P(y_x)$ and $P(y_{x'})$ using the standard adjustment formula (equation (3.19))

$$P(y_x) = \sum_w P(y | x, w) P(w), \quad P(y_{x'}) = \sum_w P(y | x', w) P(w) \quad (9.48)$$

The Ladder of Causation



Summary

- ▶ Counterfactuals
- ▶ The Twin Network Method
- ▶ Counterfactual Back-Door Criterion
- ▶ Counterfactuals in Linear-Gaussian Causal Models
- ▶ Axiomatic Characterization of Counterfactuals
- ▶ Necessary and Sufficient Causes
- ▶ The Ladder of Causation

Thank you