

The diagram shows a directed graph with the following nodes and edges:

- Nodes:**
 - N_A (grey circle)
 - A (black circle with black border)
 - N_C (grey circle)
 - \dot{A} (grey circle)
 - B (white square with black border)
 - C (grey circle with grey border)
 - \dot{B} (white diamond with black border)
 - N_F (grey circle)
 - F (grey circle)
 - \dot{C} (green circle with green border)
 - \dot{F} (grey circle)
- Edges:**
 - Grey arrows: $N_A \rightarrow A$, $N_A \rightarrow \dot{A}$, $A \rightarrow B$, $B \rightarrow F$, $C \rightarrow F$, $N_C \rightarrow C$, $N_C \rightarrow \dot{C}$, $\dot{A} \rightarrow \dot{C}$, $\dot{B} \rightarrow \dot{F}$, $N_F \rightarrow F$, $N_F \rightarrow \dot{F}$, $C \rightarrow \dot{C}$ (green arrow).
 - Black arrow: $A \rightarrow C$.
 - Blue arrow: $N_C \rightarrow C$.
 - Green arrow: $C \rightarrow \dot{C}$.

Lars Kai Hansen

Preface

This document was made as part of the course 02463 Active Machine Learning and Agency at the Technical University of Denmark, but works independently the rest of the course. It can be used by anyone who would like an introduction to causality and it has very few prerequisites.

While most of the document only requires a basic understanding of mathematics and probability (percentages etc.), a few parts use more advanced probability theory. For ease of use we have marked all sections and exercises if they require an advanced understanding of probability theory. Furthermore we mark exercises with difficulty, as well as whether they require knowledge of the programming language Python.



Section/subsection/exercise uses probability theory



Exercise uses Python



Exercise difficulties from easiest to hardest

The exercises provided are split into two parts. The first set of exercises concerns the topics covered in sections 1-6. The second set of exercises concerns the topics covered in 7-10.

Version: 1.0, April 22, 2020

Acknowledgement

We have been heavily inspired by, borrowed ideas and examples from, and learned extensively from the following resources.

Books

- Judea Pearl and Dana Mackenzie. *The Book of Why*. Basic Books, 2018
- Jonas Martin Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of Causal Inference: Foundations and Learning Algorithms*. MIT Press, 2017
- Peter Spirtes, Richard Scheines, and Clark Glymour. *Causation, Prediction, and Search*. Springer New York, 1993

Introductory **presentations**

- Ferenc Huszar. Causal Inference in Everyday Machine Learning, January 2019. Machine Learning Summer School
- Bernhard Schölkopf. Causality, January 2019. Machine Learning Summer School
- Ricardo Silva. Causality, August 2019. Cambridge Advanced Tutorial Lecture Series on Machine Learning


A great **blog** on machine learning

- Ferenc Huszar. inFERENCe, 2012. URL <https://www.inference.vc/>

And a few **articles** with great examples

- Pedro A. Ortega. Bayesian Causal Induction. *preprint*, 2013
- Pedro A. Ortega. Subjectivity, Bayesianism, and causality. *Pattern Recognition Letters*, 64:63–70, 2015
- A. Balke and J. Pearl. Probabilistic evaluation of counterfactual queries. *UCLA Department of Statistics Papers*, 2011
- R.C. Robinson. Enumeration of acyclic digraph. *Combinatorial mathematics and its applications (R.C. Bose et al.)*, 1970

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 1 |
| 1.1 | What is Causality? | 1 |
| 1.2 | Causality and Statistics | 3 |
| 1.3 | Confounders | 4 |
| 1.4 | Serious Problems in the Real World | 5 |
| 2 | Interventions | 9 |
| 2.1 | The Barometer | 11 |
| 2.2 | The Barometer - a Challenging Hypothesis | 13 |
| 3 | The Do-operator | 17 |
| 3.1 | Interventional Distribution for Barometer | 19 |
| 3.2 | Interventional Distributions in General | 21 |
| 3.3 | Causal Patterns | 26 |
| 3.4 | Randomized Trials | 27 |
| 4 | Causal Inference - The Switch  | 31 |
| 4.1 | Agent 1: An Observational Sample | 34 |
| 4.2 | Agent 2: An Interventional Sample | 35 |
| 4.3 | Prediction On Observational Sample | 37 |
| 4.4 | Prediction On Interventional Sample | 38 |
| 4.5 | More Samples for Agent 2 | 39 |
| 5 | Conditioning | 41 |
| 5.1 | Conditioning on Confounders | 41 |
| 5.2 | Large Conditioning Space | 44 |
| 5.3 | Causes, Mediators and Colliders | 45 |

Contents

5.4 Conditioning on Children 48

6 Causal Inference - In General 51

6.1 Inference 51

6.2 Information and Causal Graphs 53

6.3 Causal Patterns - Revisited 54

6.4 Solving a Graph 56

7 Knowledge and Causation 59

7.1 Who Intervened? 59


7.2 No Free Lunch Theorem 61


7.3 Ladder of Causation 62

8 Causality and AI 65

8.1 Mathematics, Science and Machine Learning 65


8.2 Models and Physics - an example 67

8.3 Machine Learning Recap  70

9 Full Causal Models  75

9.1 Generative Process Models 75

9.2 Structural Equation Models 77

10 Counterfactuals  79

10.1 A Counterfactual Question 79

10.2 Party Time! 81

10.3 Counterfactuals in General 87

10.4 Counterfactuals and Interventions 91

11 Perspectives on Causality 93

11.1 Why? 93

11.2 Explainability 97



































11.3 Bias and Fairness 98

11.4 Causality in Physics 99































11.5 Do-notation 105

Exercises 107

Python Code 109

| | |
|--|----------------|
| 1st Exercise | 111 |
| 1.1   Problems and Graphs | 112 |
| 1.2   A Paradox | 115 |
| 1.3   Intervention on Programs  | 120 |
| 1.4   A New Switch   | 123 |
| 1.5   Information Flow | 126 |
| 2nd Exercise | 129 |
| 2.1   Signal Through Variables | 131 |
| 2.2   Programming Models   | 133 |
| 2.3   Modelling Programs   | 138 |
| 2.4   Counterfactuals 1  | 145 |
| 2.5   Counterfactuals 2  | 147 |
| 2.6   Counterfactuals 3  | 151 |
| Project | 153 |
| Project in Causal Inference | 155 |
| 3.1 Project Description | 155 |
| 3.2 Experiments | 156 |
| 3.3 Inferring the Graph | 157 |
| 3.4 Making a Guess | 158 |
| 3.5 Competition | 159 |
| Project in Causal Inference - Teacher | 161 |
| 4.1 Project Scope | 161 |
| 4.2 Server | 162 |
| 4.3 Causal Model | 164 |
| 4.4 Running Server | 166 |
| 4.5 Results | 168 |
| Exercises - With Solutions | 173 |
| 1st Exercise | 175 |
| 1.1   Problems and Graphs | 176 |

Contents

| | | |
|--------------------------------|--|------------|
| 1.2 |   A Paradox | 179 |
| 1.3 |   Intervention on Programs  | 188 |
| 1.4 |   A New Switch   | 191 |
| 1.5 |   Information Flow | 194 |
| 2nd Exercise | | 197 |
| 2.1 |   Signal Through Variables | 199 |
| 2.2 |   Programming Models   | 201 |
| 2.3 |   Modelling Programs   | 207 |
| 2.4 |   Counterfactuals 1  | 214 |
| 2.5 |   Counterfactuals 2  | 216 |
| 2.6 |   Counterfactuals 3  | 220 |

Introduction

1.1 What is Causality?

This note introduces the subject of *causality* and *causal inference*. Most people have an intuitive idea of what a *cause* is, and what an *effect* is, and what it means when something *causes* something else. However, we will need some more rigid definitions. The definition of causality that we will use in this note is

Causality: Causality is a relationship by which one process or state, a *cause*, contributes to the production of another process or state, an *effect*, where the cause is partly responsible for the effect, and the effect is partly dependent on the cause.

The first key thing to note from the definition is that we have two processes or states; the cause and the effect. Examples of these two could be: the neighbours cat outside and my dog barking, a traffic jam and a road accident, a rooster crowing and the sun rising etc. The definition also sets a directional relationship between the cause and the effect. For example, it may be that the sun rising caused the rooster to crow, but it seems very unlikely that the rooster crowing caused the sun to rise. Finally, we note that the cause is *partly* responsible for the cause and the effect is *partly* dependent on the cause. This is important as there may be many causes to an effect and a cause may not guarantee an effect. For example the dog may be barking at something else than the neighbours cat, and sometimes when the neighbours cat is close the dog may be asleep (and probably not barking).

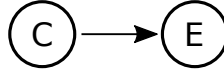
Furthermore we define

Causal Inference: Causal inference is the process of drawing a conclusion about the causal relationship between processes or states.

Causal inference is simply when we attempt to figure out what causes what. This will often be the goal of, for example, looking out the window to see what the dog is barking at.

One great tool that's readily available for explaining assumptions and hypotheses is to illustrate them with *causal graphs*

Causal Graph: A causal graph is a graph describing causal relationships. An example of a causal graph is



where the *cause* (C) has an edge directed towards the *effect* (E).

One example is shown below where the neighbours cat walking by (C) causes the dog to start barking (B).

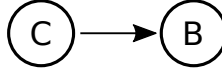
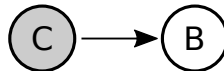


Figure 1: Causal graph for the dog-barking-at-cat situation.

When analysing systems we will often be measuring some variables, while others may not be available even though we know that they are important. We thus make the following definitions

Observable Variable: A variable whose state we can measure/observe. We visualize these by white nodes in graphs.

Hidden Variable: A variables whose state we can not measure/observe. We visualize these by grey nodes in graphs. One example of the graphical representation is



where the can hear the dog bark (B), but we can not see the cat (C) from where we are comfortably sitting on the living room couch.

1.2 Causality and Statistics

Statistics is an important tool that has been vital for the success of science and engineering. It is the foundation we use to determine whether experiments show the results we expect, or whether we need to change our minds.

Still, the statistics community has had quite a bit of trouble with causality. It all boils down to the very famous phrase

Correlation does not imply causation.

One example is the following fallacious conclusion.

As ice cream sales increase, the rate of drowning deaths increases sharply.

Therefore, ice cream consumption causes drowning.

This hypotheses seems implausible. The root of the problem is that more ice cream is sold during the summer and more people go swimming during the summer. If more people go swimming, then more people are in risk of drowning. Thus ice cream and downing are both caused by a third thing: the season. Figure 2 graphs the two different explanations (hypotheses) of the phenomena.

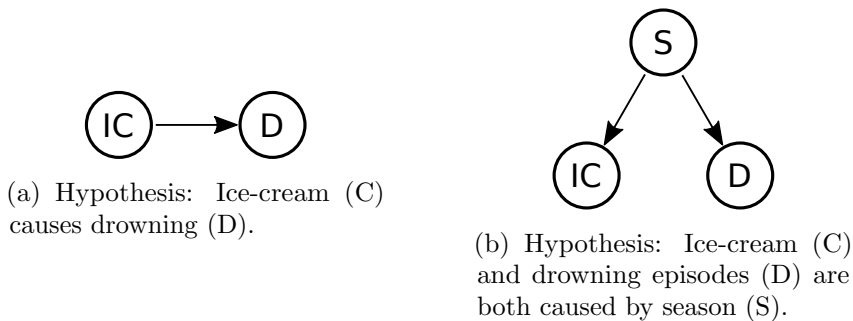


Figure 2: Two hypotheses for the correlation between ice-cream and drowning episodes.

The phrase "*correlation does not imply causation*" is important because people have believed very arbitrary, and sometimes problematic, things due

to correlation. For example in Europe during the the Middle Ages people believed that lice where beneficial to your health, because sick people rarely had lice. Turns out lice leaves people when their body temperature increases and lice did not at all have health benefits.

For a long time, many statisticians were opposed to the idea that we can even discuss causality. But causality is still quite important for science and engineering. If I put a burning Bunsen burner under a pot of water, the water temperature will increase over time. I can do this experiment millions of times with the same result: the temperature of the water is very correlated with the flame. So did the flame cause the increase in temperature or not? In science we accept experimental results as a way of showing causal relationships, which motivates the revised version of the phrase

Correlation does not imply causation. But sometimes it does.

We are going to learn how to identify those sometimes.

1.3 Confounders

The ice-cream-and-drowning problem provides an example of a *confounder*.

Confounder: A confounder is a variable that influences both dependent and independent variables in a statistical/causal analysis.

In probability/statistical theory, a dependent variable is a variable that we are attempting to determine based on independent variables. In causality we take that relationship even more serious, where the dependent variable is *caused* by the independent variables.

In the aforementioned problem the dependent variable is drowning episodes while the only independent variable is ice-cream sales, because we are testing a hypothesis in which ice-cream sales cause drownings. The confounding

1.4. *Serious Problems in the Real World*

variable is season, which influences both drowning episodes and ice-cream sales, and creates a spurious correlation which in turn motivates the incorrect hypothesis.

Confounders turn out to be absolutely essential for the correct analysis of causal relations and is the most common problem for such an analysis.

1.4 Serious Problems in the Real World

To motivate the study of causality we give two examples with serious consequences to the world.

A Dark Cloud

The first example is well known: *Does smoking cause lung cancer?*

Despite the agreement on the topic today, this was a very controversial problem in the late 1950s and early 1960s. Scientists had previously agreed on causal relationships that were simpler. For example it was known that the lack of vitamin C caused scurvy. There was a big problem with smoking though; some people didn't smoke and still got lung cancer, and some people smoked all their lives and did not get lung cancer. As we will see later, handling these problems can be done using *randomized controlled trials* (experiments). In this case though, creating an experiment where we ask non-smokers to start smoking and take the cigarettes away from smokers, was not an option. While many scientists argued that smoking did cause lung cancer, others argued that there was a gene that caused lung cancer, while also increasing the craving for smoking. Mathematically it was very difficult to explain which model was correct.

The graph for the current theory is shown in Figure 3a. Here smoking has a causal link to lung cancer. Both smoking and cancer are affected by environment variables, such as genetic heritage influencing risk of cancer and culture influencing probability of being a smoker. The environmental factors will usually be very difficult to determine and is thus considered a hidden variable. The alternative hypothesis is illustrated in Figure 3b.

1. INTRODUCTION

Here, we also assume a causal relationship from environment to both smoking and cancer, while we assume no causal link between smoking and cancer.

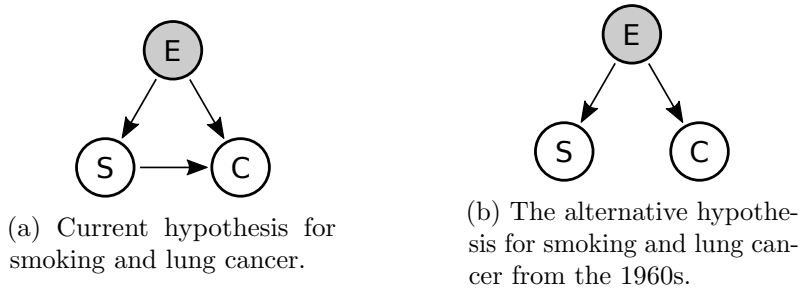


Figure 3: Two hypotheses for the probability of lung cancer (C) depending on whether the persons smokes (S). We contract other variables, such as genetic heritage, outdoors environment, culture etc., into one variable called environment (E).

The proponents of the smoking-industry abused the problem of confounders to their own advantage! They attempting to discredit the important work done by scientists to determine the health risks of smoking.

Using causal graphs can help us understand the underlying assumptions of hypotheses. It would be very good practice for more researchers to graph all assumptions for their work, as it becomes easier to specify where people disagree. As we will see later, the graphs can also be used to decide how to design experiments and analyse data, in order to test various hypotheses.

A Much Darker Cloud

Today most people seem to accept that smoking does cause lung cancer, but here is a problem with less agreement: *Is climate change caused by humans?*¹

Here we have the same problem as with smoking, but worse; we cannot do experiments on the outcome of the planets health, because we only have one! So how can we tell if CO₂ is causing climate change? Back in 2006 Al Gore attempted to warn the world about the problems of climate change in his documentary *An Inconvenient Truth*. One argument he used (with a bit a

humour) was to compare the graphs of the estimated CO₂ and temperature over the past several hundred thousand years. The graph can be seen in Figure 4, where he humorously uses a lift to be able to point to the top of the chart.

Unfortunately this approach is not bullet-proof. In Figure 5 we show that the age of Miss America varies quite like the number of murders by steam, hot vapours and hot objects in the US. This brings us back to the problem that simple correlation does not prove causal relationships, which will be the first argument used against someone like Al Gore. We therefore need to carefully consider our causal arguments and experiments.

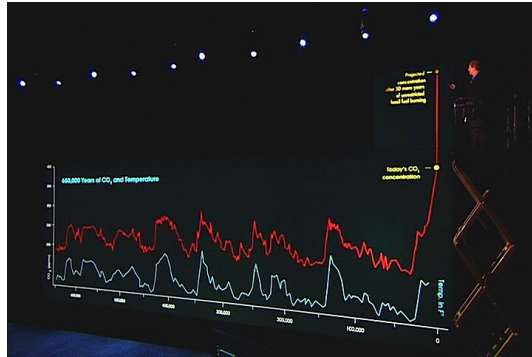


Figure 4: Al Gore comparing estimates of CO₂ and temperature.

¹it is at least to some extent

² There is a collection of these laughable, spurious correlations at <http://www.tylervigen.com/spurious-correlations>

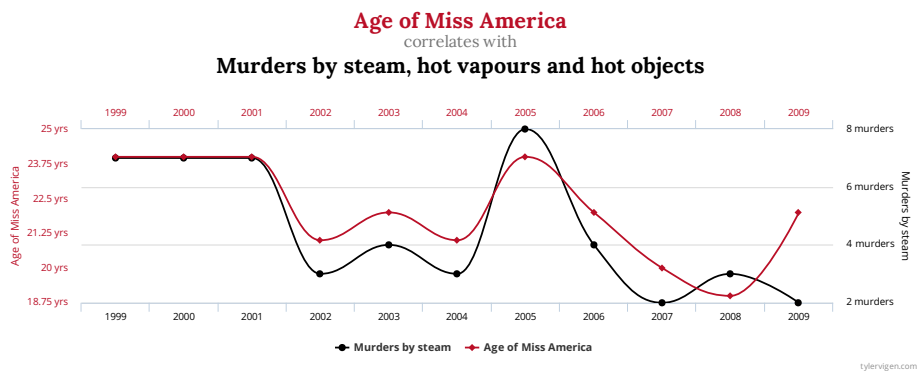


Figure 5: A spurious correlation².

Interventions

In science we make plenty of causal conclusions. I don't think there exists a universally accepted definition of the *scientific method*, but here is what Wikipedia has to say

Scientific Method: The scientific method is an empirical method of acquiring knowledge. It involves careful observation, applying rigorous scepticism about what is observed, given that cognitive assumptions can distort how one interprets the observation. It involves formulating hypotheses, via induction, based on such observations; experimental and measurement-based testing of deductions drawn from the hypotheses; and refinement (or elimination) of the hypotheses based on the experimental findings.

Notice first that the method is *empirical*; it makes observations and bases its inference on those observations. Secondly it is concerned with *knowledge*. While knowledge is difficult to define precisely one thing is for sure; knowing the causal structures of systems provides a LOT more knowledge than simply knowing the observational distributions. One can only state that more people in the 1960s got lung cancer, while the other can tell us WHY. Thirdly it is concerned with formulating hypotheses and performing *experimental tests* to eliminate and confirm such hypotheses. Making causal hypotheses and testing them is precisely what this section will introduce you to.

When we perform scientific experiments we construct a controlled environment, in which specific variables are set to specific values/conditions, and then we observe what happens to other variables. For discussing this we will use the term *intervention*.

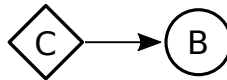
2. INTERVENTIONS

Intervention: An intervention is when we forcefully set a state or process to be specifically what we want, without affecting any other part of the causal system. We show an intervention graphically by



where a variable (V) has been changed by intervention.

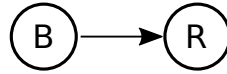
For a physical experiment we have therefore intervened on all factors which we control and set when creating the experiment. All states and processes that we are not in control of, including those we measure, are not intervened on. One example could be to get a cat and place it near our house to see if the dog starts barking, which is illustrated below.



2.1. The Barometer



(a) Hypothesis 1: Rain changes the state of the barometer.



(b) Hypothesis 2: Barometer changes the state of the rain.

Figure 6: The barometer problem.

2.1 The Barometer

Let us consider a simple situation. Over a year (Year 1) we note down rainy days and barometer readings. Figure 7 shows the recorded values and from the scatter plot in the top it seems like the two values are very correlated. Note that we don't inherently know which one should be plotted above which (we don't even know which causes which) and so this was randomly chosen. We also plotted a histogram for each quantity to see how they behave independently of the other.

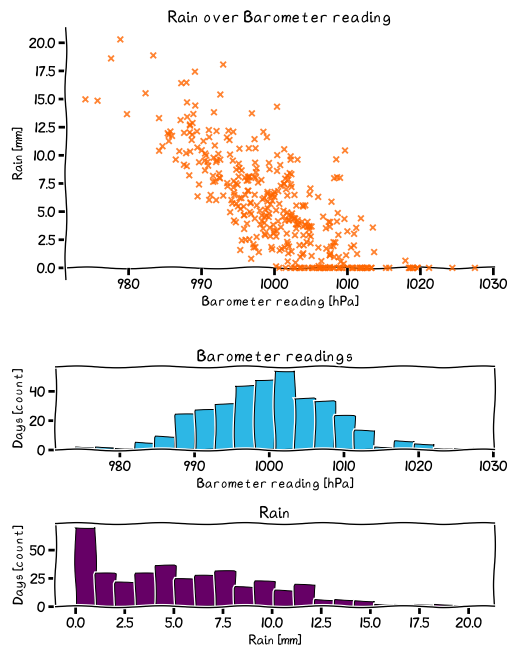


Figure 7: Data gathered about rain and barometer readings during Year 1.

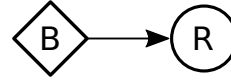
We now create two competing hypotheses; the rain causes changes in the barometer (Hypothesis 1) vs. the barometer causes changes in the weather (Hypothesis 2). The hypotheses are graphed in Figure 6. This example is of cause quite trivial, but we can use it to ensure that our methodology in causality will correctly produce the

expected result.

In order to determine which hypothesis is correct we create the following experiment; we again note down rainy days and barometer readings over a year (Year 2), but this time we physically hold the barometer reading to 1000hPa (we do intervention on the barometer). Now there is no longer anything that can affect the barometer (because we are holding it tight). Therefore we know that all causes of the barometer readings must be cancelled (all ingoing arrows). The graphs of the two hypotheses after intervention is seen in Figure 8.



(a) Hypothesis 1: Rain changes the state of the barometer.



(b) Hypothesis 2: Barometer changes the state of the rain.

Figure 8: The barometer problem after intervention on the barometer.

Under hypothesis 1 we expect the rain to do what it always does (nothing has changed for that node), but we expect the barometer to be different from normal behaviour, because we choose its reading. Under hypothesis 2 we again expect the barometer to behave differently, but this time we also expect different behaviour of the rain, because the rain depends on the barometer.

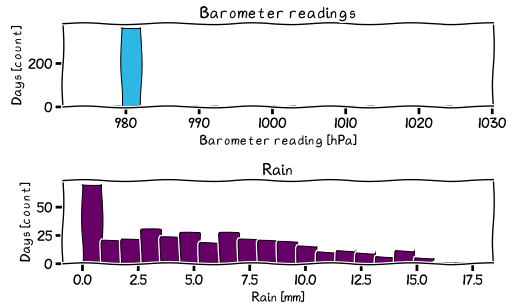


Figure 9: Histogram of barometer readings and rain amount after intervention on barometer during Year 2.

2.2. The Barometer - a Challenging Hypothesis

Now we look at our new results from the experiment (Figure 9). The histogram of rainy days seem similar as the one from Figure 7 and it seems plausible that the rain has followed the same distribution this year. The barometer though showed exactly the reading we set it to, so our intervention was certainly successful. If the rain did not change from the intervention while the barometer did, then we can accept hypothesis 1 and reject hypothesis 2 - we inferred a causal relationship!

2.2 The Barometer - a Challenging Hypothesis

We could also have attempted to make an experiment by intervening on the rain. This is of course challenging to do, but let's say someone has invented a machine that can do just that: cause rain no matter what the conditions were before and without changing anything else³. How did we avoid having to use this machine for a year? We avoided that because we made one important assumption: that only one of the two hypotheses were possible. Thus falsifying one was equivalent of verifying the other.

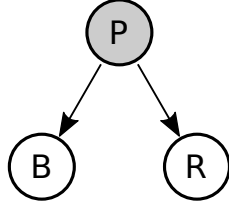
But what if someone introduced a third hypothesis (Hypothesis 3): rain and barometer readings are both caused by atmospheric pressure. Our previous experiment still falsifies hypothesis 2, but it does not verify hypothesis 1. The new hypothesis is graphed in Figure 10a.

Under hypothesis 3 we expect the same results of our experiment as with hypothesis 1. One further complication may be that we can not observe the atmospheric pressure (we measure the atmospheric pressure is measured by some kind of barometer, which is what we are currently investigating).

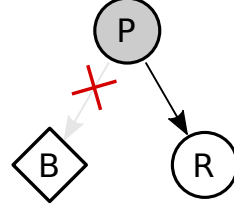
Let's therefore try out the machine that controls rain! If we set the machine to make constant rain for a year (Year 3 - depressing I know) we expect the following causal graphs for our two competing hypotheses (hypothesis 1 and 3).

³I'm guessing the machine pumps a ton of water into the air like a garden hose.

2. INTERVENTIONS

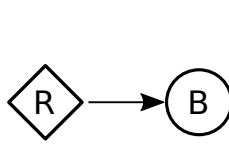


(a) Hypothesis 3: Pressure changes the state of rain and barometer.

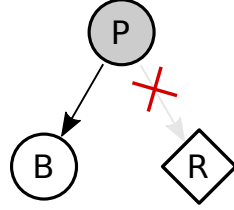


(b) Hypothesis 3: Pressure changes the state of rain and barometer after intervention.

Figure 10: The alternative hypothesis.



(a) Hypothesis 1: Rain changes the state of the barometer.



(b) Hypothesis 3: Pressure changes the state of rain and barometer.

Figure 11: Hypothesis 1 and 3 after intervening on rain.

This time we expect rain to have constant value for both of the hypothesis. For hypothesis 1 we expect the barometer to be constant as caused by the rain, but for hypothesis 3 we expect the barometer to behave normally throughout the year (have varying values depending on pressure).

After running this experiment

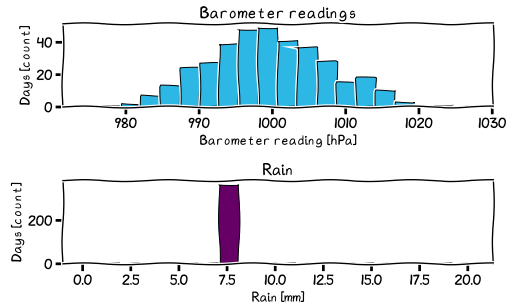


Figure 12: Histogram of barometer readings and rain amount after making it rain during Year 3.

2.2. The Barometer - a Challenging Hypothesis

we see the histograms in Figure 12. We note that the barometer is in fact *not* constant throughout the year and that rain thus does *not* cause barometer change. We have falsified hypothesis 1 while verified hypothesis 3: there is a third cause (pressure) that causes both rain and barometer reading.

The important take-home message from the barometer problem is this;

*When we consider a set of hypotheses or set of causal graphs, we are automatically assuming that **anything** not mentioned is irrelevant.*

This is quite critical, because experiments can be correctly made, while drawing wrong conclusions because of left-out variables. It also means that even if we make a causal graph with potential links between all nodes (everything influences everything else to some degree), then we are still making the assumption that anything not in the graph is irrelevant - we can never consider all factors in the universe.

2. INTERVENTIONS

SECTION 3

The Do-operator

In the previous section we have shown how to do interventions in order to test various causal hypotheses. In order to analyse this more thoroughly we need to introduce some mathematical notation and pair it with our knowledge of probability theory.

The table on the next page is a quick recap of probability theory.

Table 1: Probability theory recap.

| | | |
|--|---|---|
| $P(x)$ | Probability (sometimes called <i>marginal</i> probability) of observing event x . | |
| $P(\neg x)$ | Probability of <i>not</i> observing event x . | |
| $P(x, y) = P(x \cap y)$ | <i>Joint</i> probability of observing both events x and y . | |
| $P(x y)$ | <i>Conditional</i> probability of observing event x given that we have already observed y . | |
| $P(x \cup y)$ | Probability of observing x and/or y . | |
| Normality | For any x | $0 \leq P(x) \leq 1$ |
| Compliment | For any x | $P(\neg x) = 1 - P(x)$ |
| Additivity | x and y are mutually exclusive iff. | $P(x \cup y) = P(x) + P(y)$ |
| Totality | x and y are mutually exclusive and collectively exhaustive iff. | $P(x \cup y) = P(x) + P(y) = 1$ |
| Independence | x and y are independent iff. | $P(x, y) = P(x) P(y)$ |
| Overlap | For any x and y | $P(x \cup y) = P(x) + P(y) - P(x, y)$ |
| Multiplication | For any x and y | $P(x, y) = P(x y) P(y)$ |
| Conditional | For any x and y | $P(x y) = \frac{P(x, y)}{P(y)}$ |
| Monotonicity | For any x and y | $P(x, y) \leq P(x)$ |
| Bayes rule | The fundamental rule of learning! | $P(y x) = \frac{P(x y) P(y)}{P(x)}$ |
| <p>Uppercase letters like A, X and Y denote stochastic variables (parts of systems which we do not know the value of). These symbols are also used for nodes in causal graphs, because nodes in causal graphs represent variables which could take on different values at each experiment. Small case letter like a, x and y denote specific values for the related stochastic variables (for example x is shorthand for $X = x$).</p> | | |

3.1. Interventional Distribution for Barometer

In order to handle interventions mathematically we introduce the *do*-operator

Do-operator⁴: \dot{x} denotes that we have intervened to force event x to happen (we *do* the event). $P(y|\dot{x})$ is the probability of y given that we have intervened and forced event x . We can also set a variable to a specific value by $P(y | x \leftarrow 2)$.

Here's a couple of things to note about the do-operator. First of all it's *very* different from conditioning - the next section makes that clear. Secondly it is a form of *assignment*, similar to that of computer science. We *causally assign* a value to a variable. We can also do other do-operations like assigning a new probability distribution to a variable. For example here we assign a normal distribution to a variable

$$p(x) \leftarrow \mathcal{N}(0, 1). \quad (1)$$

This is what we do in randomized trials (more on that later), where we for example randomly select people for a test. We can also make interventions to create restrictions on the system. For example

$$P(z | x \leftarrow y), \quad (2)$$

is the probability of z , given that we force variable x to be equal to variable y (for some system) - y becomes in charge of x .

3.1 Interventional Distribution for Barometer

The first thing we need to realize is that generally

Conditional distribution does not equal interventional distribution:

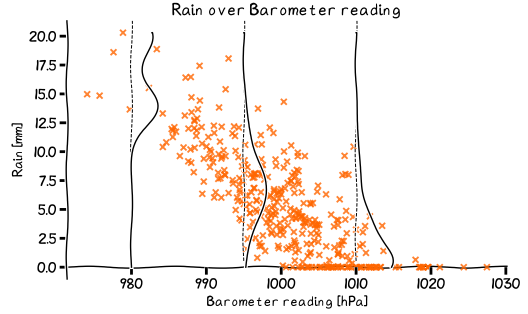
$$P(y | x) \neq P(y | \dot{x})$$

⁴ Note that in some literature the do-operator has different notation - see section 11.5.

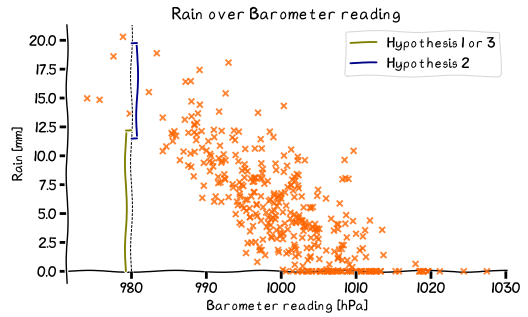
3. THE DO-OPERATOR

The conditional distribution is the distribution of a variable given that we have observed another variable, while remaining a passive observer. The interventional distribution on the other hand is the distribution of a variable given that we actively change another variable in the system.

In order to realize this consider the barometer problem. In Figure 13a we see three conditional distributions plotted on top of our data. They show the probability density of rain conditioned on our barometer readings at 980hPa, 995hPa and 1010hPa. We can see that the expected rainfall is different depending on our reading of the barometer because the two variables correlate.



(a) $P(\text{rain} \mid \text{barometer})$ at Year 1.



(b) The ranges of values expected for the interventional distribution under Hypothesis 1 and 2.

Figure 13: Analysis of Year 1.

Now we assume we are going to intervene on the value of the barometer (like Year 2) and set it to 980. If Hypothesis 2 is true (barometer reading causes rainfall), then we expect that our intervention will change the behaviour of rainfall for that year. In this case we expect the rainfall to be mostly in the range shown in blue in Figure 13b, because that is where most rainfall is conditioned on the barometer reading being 980 - the conditional from Figure 13a.

3.2. Interventional Distributions in General

If on the other hand Hypothesis 1 or 3 is true (rain or pressure causes the barometer's reading), then we expect our intervention to have no effect on the rainfall. We thus expect the rainfall to be mostly in the green range of 13b, because that is generally where most rainfall is.

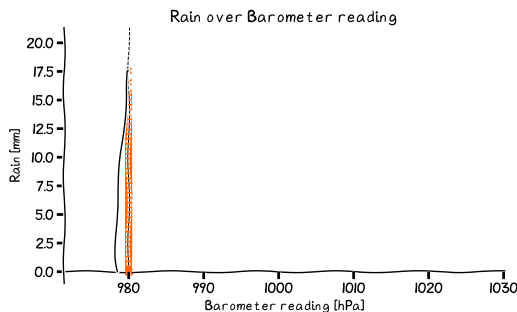


Figure 14: Interventional distribution from Year 2. $P(\text{rain} \mid \text{barometer} = 980)$

Since we did collect interventional data during Year 2, we can plot the interventional distribution as in Figure 14. Here all the datapoints are on-top of each other, because they all have the same barometer reading. The result corresponds well with the expectations of hypothesis 1 or 3 (the barometer reading does not cause anything).

3.2 Interventional Distributions in General

Let's consider the interventional distribution in a more generic setting. We have two variables X and Y and we believe X causes Y . We have contracted all shared causes of X and Y into one node called C - confounders. We have contracted all causes of X which does not affect Y directly into one node called A - ancestors. We have contracted all variables affected by y into one node called D - descendants. Finally we have contracted all nodes between x and y into one node called M - mediators. The situation is shown the Figure 15.

Note that there could be additional links between C and D , A and D etc., but it would not change the outcome of the following analysis. One assumption that is important though, is that the graph does not have cycles (at least not any involving x and y)⁵.

⁵Causal systems with cycles is in general very difficult to analyse mathematically and is part of contemporary research.

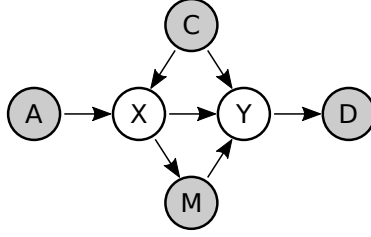


Figure 15: A quite generic causal graph investigating causal relationship between x and y .

We notice the following conditional independencies

$$P(y \mid x, a) = P(y \mid x) \quad (3)$$

$$P(x \mid y, d) = P(x \mid y). \quad (4)$$

The conditional probability of y becomes

$$P(y \mid x) = \sum_{acm} P(y \mid x, a, c, m) P(a, c, m \mid x) \quad (5)$$

$$= \sum_{cm} P(y \mid x, c, m) P(m \mid c, x) P(c \mid x) \quad (6)$$

$$= \sum_{cm} P(y \mid x, c, m) P(m \mid x) P(c \mid x), \quad (7)$$

where we can remove a because it is independent of y conditioned on x .

We now intervene on x resulting in the following interventional graph

The interventional probability of y is

$$P(y \mid \dot{x}) = \sum_{acm} P(y \mid x, a, c, m) P(m \mid x) P(a, c) \quad (8)$$

$$= \sum_{cm} P(y \mid x, c, m) P(m \mid x) P(c). \quad (9)$$

We now see the difference between the interventional distribution and the

3.2. Interventional Distributions in General

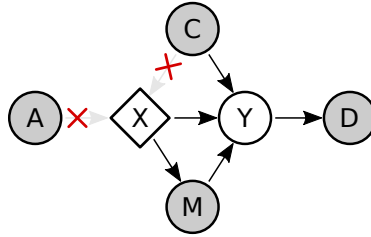


Figure 16: The interventional (on x) version of Figure 15.

conditional distribution by

$$P(y | x) = \sum_{cm} P(y | x, c, m) P(m | x) P(c | \mathbf{x}), \quad (10)$$

$$P(y | \dot{x}) = \sum_{cm} P(y | x, c, m) P(m | x) P(c). \quad (11)$$

The conditional distribution conditions all variables in the system on X , while the interventional distribution only conditions the *downstream* variables on X . This difference in conditioning is exactly what makes the interventional distribution special.

Perhaps more importantly,

*If there are no confounders for the causal problem at hand, then the conditional distribution **equals** the interventional distribution, when intervening on the **cause**. (assuming no cycles)*

If we know what is the cause, we can therefore reduce our problem of inferring causal relationships to that of correctly handling the confounders of the problem.

We also conclude is that any ancestor (A) of X , which is only an ancestor of Y *through* X , is irrelevant for both distributions, because all information about A passed through X which is known in both cases. Also the

3. THE DO-OPERATOR

descendants (D) of Y , that are only descendants of X *through* Y , are irrelevant because all information they have are already captured in Y .

3.2. Interventional Distributions in General

In this example we knew that the causal direction was from X towards Y , and definitely not the other way. This is not always the case though. The *disambiguation problem* concerns situation in which we don't know the direction. If two events are always observed together, then which of them caused the other? An even more generic causal graph is thus

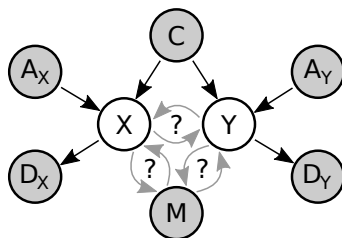


Figure 17: An even more generic causal graph investigating causal relationship between x and y .

Here we have ancestors and descendants for both X and Y , while still only having a single confounders C . The directions of the links between X and Y and the mediators are not know. For this problem we now need to both handle confounders *and* disambiguation of direction (while we are still assuming *no* cycles).

When inferring a causal relationship between two variables we need to handle all confounders as well as the disambiguation problem. (assuming no cycles)

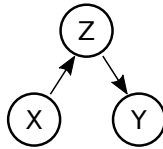
Another complication can be if we want to determine the strength of the *direct* causal relationship between X and Y , in which case the mediators M can also be trouble.

3.3 Causal Patterns

When analysing causal graphs it's very useful to consider three types of arrangements of three causal variables.

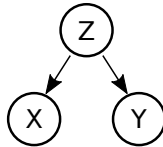
The first arrangement we are considering is the

Chain Junction: An arrangement where one variable, X , is *directly* causing another, Z , which is *directly* causing a third, Y . We also say that X causes Y (just not directly) and we call Z a *mediator*.



The second arrangement is one we considered in the barometer problem, where pressure the confounding variable.

Confounding Junction / Fork: Two variables X and Y share a common ancestor Z . Both X and Y holds information about Z , and information on Z provides information about both X and Y .

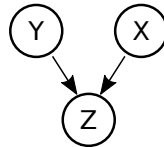


Since confounding is a very central concept in causal analysis, this pattern is quite crucial.

3.4. Randomized Trials

The third arrangement can be a bit more tricky in causal analysis - for reason that will later become apparent. We call the arrangement a

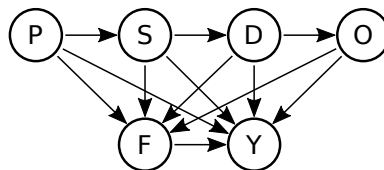
Collider: is composed of three causal variables, where two of them, X and Y , directly cause the third one, Z . Z thus holds information about both X and Y , while X and Y each holds some information about Z .



3.4 Randomized Trials

One essential part of scientific studies has been the concept of randomized trials. Let us consider a situation in which a scientist wants to help a set of farmers. The farmers have different fertilizers F and wants to know which one produces highest yield Y based on the subjected plant P , soil fertility S , water drainage D of the area and other factors O . This is actually a situation which the famous statistician R. A. Fisher was concerned with⁶. The farmers have worked in their field for many years and have some guesses for what to use where, but they don't agree and have made no shared documentation of strategies.

We segment their fields into various segments, all of which we can monitor. If we start monitoring the segments we get data on all parameters based on the following causal model.



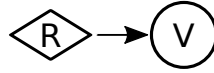
⁶Judea Pearl and Dana Mackenzie. *The Book of Why*. Basic Books, 2018.

3. THE DO-OPERATOR

From this model we can see a problem. We only want a single connection between F and Y , but there are 4 confounders!. Previously we saw that we can solve this with intervention. For example we could simply choose a single fertilizer to use everywhere and then we could determine the effect of that one fertilizer to the yield. This does not solve our task though, as we want to compare the fertilizers and not just evaluate a single one. We therefore need to test all the fertilizer, but also make sure that the confounders are not a problem. This is where the randomized trial becomes relevant.

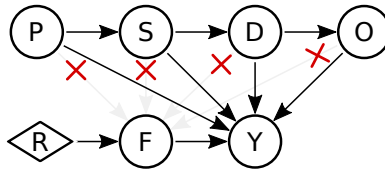
Randomized Trial: An experiment with randomization on causal variables in order to determine their effect.

We depict a randomized variable in a graph by



to illustrate that the variable V is still a stochastic variable, but that it now follows the distribution from the randomization R . All incoming links to V can now be ignored.

If we randomly choose a fertilizer for each segment, then we eliminate all relationships that affect the farmers choice of fertilizer. The randomization is a special case of intervention, where we do not intervene with a single value, but rather we replace the original distribution of a variable with a custom made one. By randomization of the fertilizer we get data from a different causal model



In this causal model there is only one causal link between F and Y , which is the one we want to measure.

Randomized trials/experiments are extremely important for scientific research and is known by many names. A/B testing for example is a special

3.4. Randomized Trials

case with two possible values for a causal variable. In medicine A/B testing will typically be done by randomly assigning a new type of medicine (A) or a placebo (B) to patients and measure their improvement.

Let's Do It!

Let us try and use our do-notation for the above mentioned randomized trial. We wish to determine the quantity

$$P(Y | F) = \frac{P(Y, F)}{P(F)}. \quad (12)$$

Originally the joint distribution of the yield and fertilizer is

$$P(Y, F) = \sum_{P, S, D, O} P(Y | F, P, S, D, O) \textcolor{red}{P}(F | \textcolor{red}{P}, \textcolor{red}{S}, \textcolor{red}{D}, \textcolor{red}{O}) P(P, S, D, O),$$

where P, S, D and O are confounding, as marked by red.

Let's instead employ a randomized trial. We want to test two different fertilizers, f_1 and f_2 - as well as the control case where we do not use any fertilizer f_0 . Say we use a die to determine the fertilizer used. The value for F now only depends on our die:

$$F \leftarrow \begin{cases} f_0 & \text{1 or 2} \\ f_1 & \text{3 or 4} \\ f_2 & \text{5 or 6} \end{cases} \quad (13)$$

Using a die we intervene on the *distribution* of F

$$P_F \xleftarrow{\circ} \mathcal{C}(1/3, 1/3, 1/3), \quad \mathring{P}_F(f_0) = \mathring{P}_F(f_1) = \mathring{P}_F(f_2) = \frac{1}{3}, \quad (14)$$

where $\mathcal{C}(1/3, 1/3, 1/3)$ is a categorical distribution with $1/3$ chance of being either f_0 , f_1 or f_2 . The important property of the new distribution of F is that it is *independent* of all other variable.

3. THE DO-OPERATOR

The interventional distribution now becomes

$$\begin{aligned}\dot{P}(Y, F) &= \sum_{P, S, D, O} P(Y | F, P, S, D, O) \dot{P}(F | P, S, D, O) P(P, S, D, O) \\ &= \sum_{P, S, D, O} P(Y | F, P, S, D, O) \dot{P}(F) P(P, S, D, O)\end{aligned}\tag{15}$$

$$= \dot{P}(F) \sum_{P, S, D, O} P(Y | F, P, S, D, O) P(P, S, D, O)\tag{16}$$

$$= \dot{P}(F) P(Y | F).\tag{17}$$

Now the marginal conditional probability $P(Y | F)$ stands alone, which is exactly the quantity we wish to measure! Now we can divide by our known probability distribution of F and find

$$P(Y | F) = \frac{\dot{P}(Y, F)}{\dot{P}(F)}.\tag{18}$$

Therefore if we get samples and measure the interventional distribution $\dot{P}(Y, F)$ - by doing the randomized trial with a die - we can determine the effects of $P(Y | F)$.

Causal Inference - The Switch



In this section we will get to the goal of all this causality-talk; the ability to infer causal relationships. We will use our knowledge about intervention and the do-operator to solve a problem similar to the barometer problem, but where the answer is not intuitively known: the switch problem.

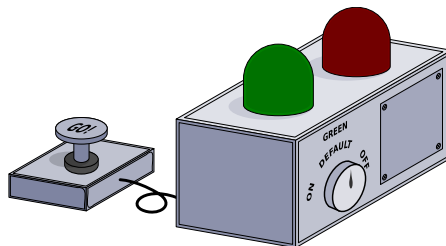


Figure 18: The switch problem.

We have a box with a green and a red light on it. If we press GO on the related contact, one of four things happen. Either both light turn on, both stay off, the green light turns on or the red light turns on. The process of the outcome is stochastic, but we can get a lot of sample by keep pressing GO and write down what happens.




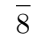
| | $\neg g$ | g |
|----------|---|---|
| $\neg r$ |  |  |
| r |  |  |
| | $\frac{3}{8}$ | $\frac{1}{8}$ |
| | $\frac{1}{8}$ | $\frac{3}{8}$ |

Table 1: Recorded outcomes of switch experiment in ratios.

We use the box a ton of times and write down all outcomes. We compute how often each outcome occurs and write the results in a confusion matrix seen in Table 1. It seems that the light are

usually on together and off together, while sometimes they differ.

We use this to hypothesise that one light is the cause of the other light - with some noise. At this point though, we do not know which light causes which.

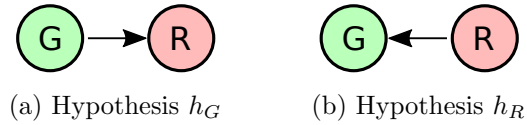


Figure 19

In Figure 20 we have illustrated the hypotheses and possible outcomes.

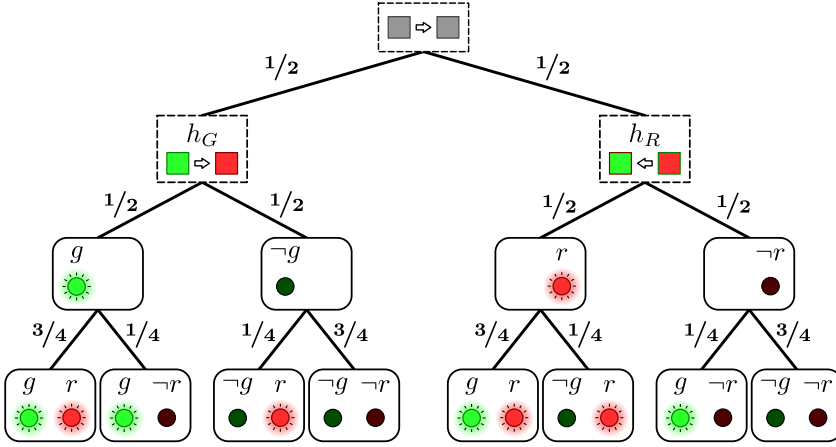


Figure 20: Hypotheses and outcomes of the switch problem.

The dashed-lined boxes are hypotheses. The top one hypothesises that one light causes the other. On the left we hypothesise that green causes red (h_G), and on the right we hypothesise that red causes green (h_R). Below the hypotheses we have illustrated all outcomes in a hierarchical manner where cause precedes effect (if green causes red, then the state of green is above the state of red).

At each edge we have noted a probability. At the top we indicate that we believe the two hypotheses are equally probable. For both hypotheses there is a 50/50 chance of the cause-light turning on, and then some probabilities for the behaviour of the effect-light. The probabilities are made so that they corresponds to the ratios in Table 1.

4.1 Agent 1: An Observational Sample

We now press GO, see that *both lights turn on* and give this datapoint to *Agent 1*. Agent 1 will attempt to determine which hypothesis is most probable, by applying the probabilities from Figure 20 with the sample and Bayes theorem.

$$P(h_G | g, r) = \frac{P(g, r | h_G) P(h_G)}{P(g, r | h_G) P(h_G) + P(g, r | h_R) P(h_R)} \quad (1)$$

$$= \frac{P(r | g, h_G) P(g | h_G) P(h_G)}{P(r | g, h_G) P(g | h_G) P(h_G) + P(g | r, h_R) P(r | h_R) P(h_R)} \quad (2)$$

$$= \frac{\frac{3}{4} \cdot \frac{1}{2} \cdot \frac{1}{2}}{\frac{3}{4} \cdot \frac{1}{2} \cdot \frac{1}{2} + \frac{3}{4} \cdot \frac{1}{2} \cdot \frac{1}{2}} = \frac{1}{2} = P_{A1}(h_G), \quad (3)$$

$$P(h_R | g, r) = \frac{P(g, r | h_R) P(h_R)}{P(g, r | h_G) P(h_G) + P(g, r | h_R) P(h_R)} \quad (4)$$

$$= \frac{P(g | r, h_R) P(r | h_R) P(h_R)}{P(r | g, h_G) P(g | h_G) P(h_G) + P(g | r, h_R) P(r | h_R) P(h_R)} \quad (5)$$

$$= \frac{\frac{3}{4} \cdot \frac{1}{2} \cdot \frac{1}{2}}{\frac{3}{4} \cdot \frac{1}{2} \cdot \frac{1}{2} + \frac{3}{4} \cdot \frac{1}{2} \cdot \frac{1}{2}} = \frac{1}{2} = P_{A1}(h_R). \quad (6)$$

Based on this datapoint, Agent 1 concludes that the two hypotheses are equally probable - just like what we started out with. Actually no matter what the outcome is from the lights, and no matter how many datapoints we get, we will not come any closer to determining the causal structure. Since you have made it this far you should of cause have expected this; we can not make causal inference simply from observing the data.

4.2 Agent 2: An Interventional Sample

There is a switch on the side of the box, which allows us to control the state of the green light. It can be set to either DEFAULT, for the currently observed behaviour, or to ON or OFF, which forces the green light to be on or off respectively. If we use the switch to force the green light on, we intervene on the system. The space of hypotheses is the same, but the space of states is now reduced, which is illustrated in Figure 21. Some probabilities are now zero and one, because of the forced green light.

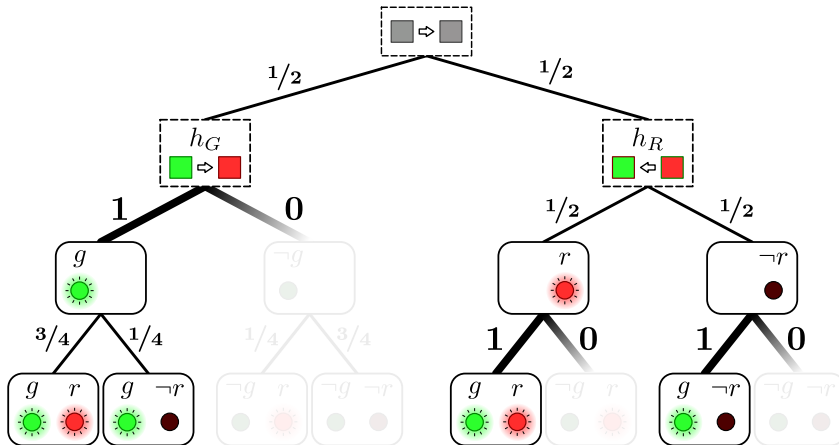


Figure 21: Hypotheses and outcomes of the switch problem during intervention on the green light.

We force green to be on, press GO and see both lights on. The outcome is identical to the observational sample, but we know it is interventional and hand it to *Agent 2*. *Agent 2* also computes the probability of each

hypothesis (differences are highlighted in blue)

$$P(h_G | \overset{\circ}{g}, r) = \frac{P(g, r | \overset{\circ}{g}, h_G) P(h_G)}{P(g, r | \overset{\circ}{g}, h_G) P(h_G) + P(g, r | \overset{\circ}{g}, h_R) P(h_R)} \quad (7)$$

$$= \frac{P(r | \overset{\circ}{g}, h_G) P(g | \overset{\circ}{g}, h_G) P(h_G)}{P(r | \overset{\circ}{g}, h_G) P(g | \overset{\circ}{g}, h_G) P(h_G) + P(g | \overset{\circ}{g}, r, h_R) P(r | h_R) P(h_R)} \quad (8)$$

$$= \frac{\frac{3}{4} \cdot \textcolor{blue}{1} \cdot \frac{1}{2}}{\frac{3}{4} \cdot \textcolor{blue}{1} \cdot \frac{1}{2} + \textcolor{blue}{1} \cdot \frac{1}{2} \cdot \frac{1}{2}} = \frac{\textcolor{blue}{3}}{\textcolor{blue}{5}} = \textcolor{blue}{0.6} = P_{A2}(h_G), \quad (9)$$

$$P(h_R | \overset{\circ}{g}, r) = \frac{P(g, r | \overset{\circ}{g}, h_R) P(h_R)}{P(g, r | \overset{\circ}{g}, h_G) P(h_G) + P(g, r | \overset{\circ}{g}, h_R) P(h_R)} \quad (10)$$

$$= \frac{P(g | \overset{\circ}{g}, r, h_R) P(r | h_R) P(h_R)}{P(r | g, h_G) P(g | \overset{\circ}{g}, h_G) P(h_G) + P(g | \overset{\circ}{g}, r, h_R) P(r | h_R) P(h_R)} \quad (11)$$

$$= \frac{\textcolor{blue}{1} \cdot \frac{1}{2} \cdot \frac{1}{2}}{\frac{3}{4} \cdot \textcolor{blue}{1} \cdot \frac{1}{2} + \textcolor{blue}{1} \cdot \frac{1}{2} \cdot \frac{1}{2}} = \frac{\textcolor{blue}{2}}{\textcolor{blue}{5}} = \textcolor{blue}{0.4} = P_{A2}(h_R). \quad (12)$$

Now we see a difference in the belief in the two hypotheses! When we intervene to force the green light on and red also turns on, Agent 2 concludes that h_G more probable than h_R .

4.3 Prediction On Observational Sample

Agent 1 believes that the two hypotheses concerning the light-box are equiprobable. Agent 2 on the other hand, believes that the probability of h_G being correct is 60%, while the probability of h_R is only 40%.

We now get an observational sample (no intervention) and observe the green light only - it turns on. We then ask the two agents to compute the probability of the red light being on/off.

The probability of the red light being on for h_G can be read off of Figure 20.

$$P(r | g, h_G) = \frac{3}{4}. \quad (13)$$

The probability of the red light being on for h_R can be computed using Bayes theorem

$$P(r | g, h_R) = \frac{P(g | r, h_R) P(r | h_R)}{P(g | r, h_R) P(r | h_R) + P(g | \neg r, h_R) P(\neg r | h_R)} \quad (14)$$

$$= \frac{\frac{3}{4} \cdot \frac{1}{2}}{\frac{3}{4} \cdot \frac{1}{2} + \frac{1}{4} \cdot \frac{1}{2}} = \frac{\frac{3}{8}}{\frac{4}{8}} = \frac{3}{4}. \quad (15)$$

The agents' predictions are simply a weighted mean of these two values, but since they are identical we trivially have

$$P_{A1}(r | g) = P(r | g, h_G) P_{A1}(h_G) + P(r | g, h_R) P_{A1}(h_R) \quad (16)$$

$$= \frac{3}{4} \cdot \frac{1}{2} + \frac{3}{4} \cdot \frac{1}{2} = \frac{3}{4} \quad (17)$$

$$P_{A2}(r | g) = P(r | g, h_G) P_{A2}(h_G) + P(r | g, h_R) P_{A2}(h_R) \quad (18)$$

$$= \frac{3}{4} \cdot \frac{3}{5} + \frac{3}{4} \cdot \frac{2}{5} = \frac{3}{4}. \quad (19)$$

The two agents thus makes identical *observational predictions*, because both hypotheses are capable of perfectly explaining the observational distribution.

4.4 Prediction On Interventional Sample

Now we get an interventional sample, by forcing the green light on and again ask the two agents to predict the probability of the red light.

The probability of the red light being on, under hypothesis h_G , during intervention is the same as if we were to condition on the green light

$$P(r \mid \overset{\circ}{g}, h_G) = \frac{3}{4}. \quad (20)$$

The probability under hypothesis h_R changes to

$$\begin{aligned} P(r \mid \overset{\circ}{g}, h_R) &= \frac{P(g \mid \overset{\circ}{g}, r, h_R) P(r \mid h_R)}{P(g \mid \overset{\circ}{g}, r, h_R) P(r \mid h_R) + P(g \mid \overset{\circ}{g}, \neg r, h_R) P(\neg r \mid h_R)} \\ &= \frac{1 \cdot \frac{1}{2}}{1 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2}} = \frac{1}{2}. \end{aligned} \quad (21)$$

That is, the interventional distribution is equal to the original marginal distribution because under H_R there is no influence on the red light no matter what we do to the green light.

The interventional predictions of the two agents are

$$P_{A1}(r \mid \overset{\circ}{g}) = P(r \mid \overset{\circ}{g}, h_G) P_{A1}(h_G) + P(r \mid \overset{\circ}{g}, h_R) P_{A1}(h_R) \quad (22)$$

$$= \frac{3}{4} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = \frac{3}{8} + \frac{2}{8} = \frac{5}{8} = 0.625 \quad (23)$$

$$P_{A2}(r \mid \overset{\circ}{g}) = P(r \mid \overset{\circ}{g}, h_G) P_{A2}(h_G) + P(r \mid \overset{\circ}{g}, h_R) P_{A2}(h_R) \quad (24)$$

$$= \frac{3}{4} \cdot \frac{3}{5} + \frac{1}{2} \cdot \frac{2}{5} = \frac{9}{20} + \frac{4}{20} = \frac{13}{20} = 0.65. \quad (25)$$

The two agents has *different interventional predictions*, because they have learned different causal structures! The difference between these two numbers is quite small, but remember this comes from a *single* training point :)

4.5 More Samples for Agent 2

To finish off the switch problem we here show how the causal inference would go about with more data. If we pressed GO a total of 10 times, while forcing green light on, and observed that the red light was on 7 times and off 3 times. We can further improve Agent 1's knowledge about the hypothesis space. This time we remove some trivial terms⁷ to simplify the expressions

$$P(h_G \mid \mathring{g}, \{(\neg r)^3, r^7\}) \quad (26)$$

$$= \frac{P(\{(\neg r)^3, r^7\} \mid \mathring{g}, h_G) P(h_G)}{P(\{(\neg r)^3, r^7\} \mid \mathring{g}, h_G) P(h_G) + P(\{(\neg r)^3, r^7\} \mid \mathring{g}, h_R) P(h_R)} \quad (27)$$

$$= \frac{\left(\frac{3}{4}\right)^7 \cdot \left(\frac{1}{4}\right)^3 \cdot \frac{1}{2}}{\left(\frac{3}{4}\right)^7 \cdot \left(\frac{1}{4}\right)^3 \cdot \frac{1}{2} + \left(\frac{1}{2}\right)^7 \cdot \left(\frac{1}{2}\right)^3 \cdot \frac{1}{2}} = \frac{\frac{2187}{2097152}}{\frac{2187}{2097152} + \frac{1}{2048}} \approx 0.68. \quad (28)$$

Using this data, Agent 1 have further increased its belief that h_G is the correct hypothesis and not h_R . Of cause it has still not completely ruled out hypothesis h_R , because that system can also create the same data - it's just less likely.

⁷and abuse mathematical notation a bit

Conditioning

5.1 Conditioning on Confounders

We have looked at how to manipulate a causal graph using intervention and randomization (which is also an intervention). We will now look at a statistical tool which can also help, *conditioning*. Conditioning can help us with confounding variables, which as mentioned is critical for causal analysis (the barometers problem, ice-cream-drowning problem and fertilizer problem are all concerned with confounder issues).

Consider again the problem of summertime confounding the relationship between ice-cream sales and drowning incidents. Let's say we picked some data for days in the winter and days in the summer. It may seem a bit harsh that there are drownings every day, but let's assume this is a very big country with a population that really likes extreme water-sports.

Figure 22a shows the data for summer and winter. We have also plotted a regression line, which on the face of it explains the relationship between ice-cream sales and drowning episodes. We see a clear relationship between the two variables which creates the mentioned spurious correlation.

In Figure 22b we condition on summer data-points and winter data-points respectively. When we analyse the two groups independently we do not see the spurious correlation any more. Actually, for the two groups apart it looks like there is no correlation at all! We have thus solved the confounding problem by conditioning on the confounding variable.

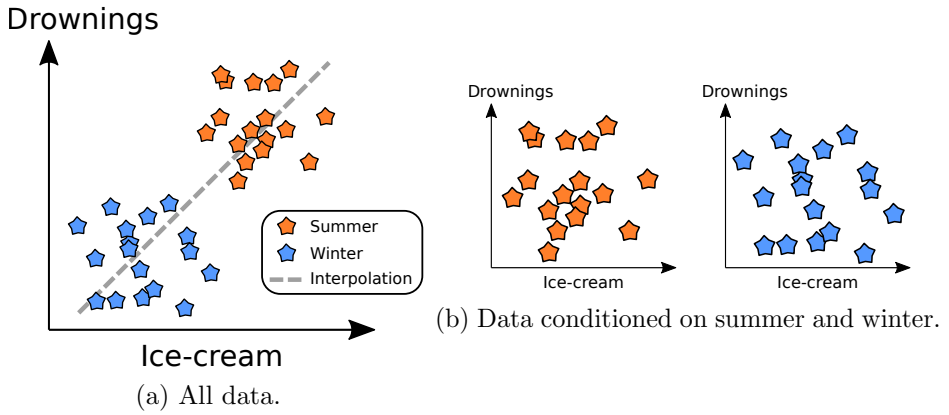


Figure 22: Data for drowning episodes and ice-cream sales during summer and winter.

Conditioning can solve the problem of confounders!

As we have seen in section 3.2, solving confounders is both essential and sufficient for inferring causal relationships, which makes conditioning an extremely important tool.

5.1. Conditioning on Confounders

Consider what conditioning actually is; we set a variable to a constant value, based on observational data. It is like we make one causal model for each value of the conditioned variable. We have illustrated this interpretation in Figure 23, where we condition on variable Z . The square node of Z indicates that this node is a *constant*, while X and Y remain stochastic variables (round). For each possible values of Z we have a square box and a related set of stochastic nodes for X and Y , because these may behave differently for each value of Z . The ancestor A of Z is also duplicated, because the distribution of A is very different depending on what it caused Z to be. When we intervened on nodes we removed their causal links from their ancestors, but when we condition on Z we not **not** remove them, because A is still the cause of Z . We are merely selecting the A for which Z is the conditioned value.

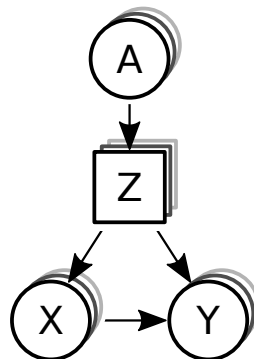


Figure 23: Graphical interpretation of conditioning.

The causal relationships potentially change dramatically depending on the value of Z . It is possible that X causes Y for one value of Z , while there is no causal relationship between X and Y for another value of Z . It could even be that the causal relationship reverses so that Y causes X for a third value of Z . We have further illustrated this situation in Figure 24. For each value of the conditioned variable Z the causal graph has different distributions and potentially different causal structures.

While conditioning is an extremely important tool there is a couple of catches, which are investigated in the next sections.

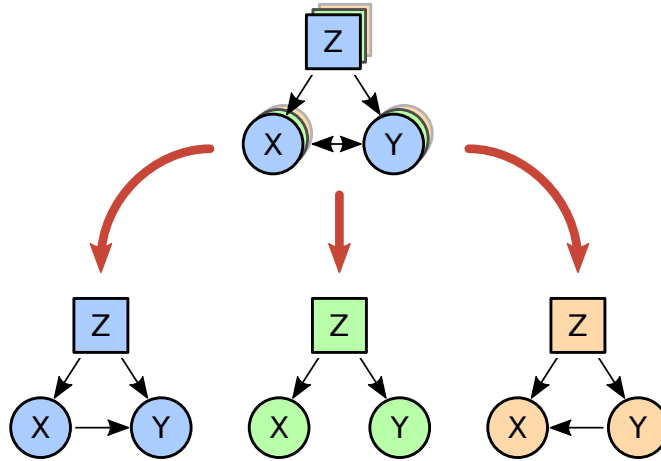


Figure 24: Expanded example of graphical interpretation of conditioning.

5.2 Large Conditioning Space

Let us say you have two confounding variables, which can each take on 10 values. If we condition on both of these we will have a total of $10 \times 10 = 100$ different combinations. In order to make a proper analysis of each conditioned group, we would therefore roughly need 100 times as much data! This problem becomes even more complicated if we are considering continuous variables, because how do you correctly condition on a variable which (in theory) has infinitely many possible values? How do you discretize? Let us say we also collected some data on the ice-cream-drowning problem from spring. Some of these datapoints will be similar to winter, while some of them will be similar to summer. Figure 25 illustrates this data.

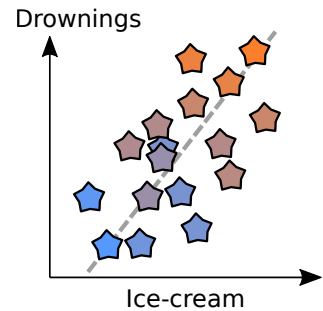


Figure 25: Data conditioned on spring.

We have colourized this data, so that the early spring-datapoints have a

5.3. Causes, Mediators and Colliders

colour similar to winter-datapoints and the late spring-datapoints have a colour similar to summer-datapoints. The spurious correlation between ice-cream sales and drowning episodes have now reappeared! The datapoints close to summer are in the top-right corner, while the winter-like datapoints are in the lower-left corner. This is because we haven't correctly conditioned on our season. In reality spring appears to be a transition period between winter and summer, and can therefore not really be conditioned on by simple measures.

In conclusion

Conditioning can become a problem if the space of conditioned variable-states becomes large or continuous.

5.3 Causes, Mediators and Colliders

Direct Cause and Mediators

Consider the simple causal graph in Figure 26a. If we condition on $X = x$ in this graph we conceptually make X a constant as in Figure 26. If X is constant then Y will always take on the same value. If we sample from this graph we will get various values for Y and always the same value for X . By conditioning on the cause we are trying to analyse we therefore create a problem that wasn't there.

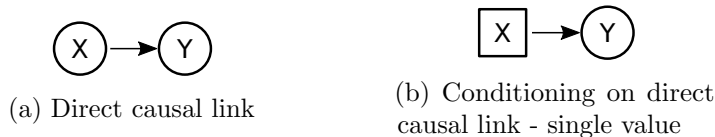


Figure 26: Unfortunate conditioning on direct causal link.

While the above example seems a bit silly, consider instead the *chain junction* from section 3.3. If we condition on the mediator Z , we conceptually make it a constant. No information is ever passed from X to Y , so it

may look like X does not cause Y - even though it does! If we have a theory about Z being a *confounder*, while in reality Z is a *mediator*, then conditioning on it ruins our analysis. We must therefore be sure about Z being a confounder before we condition on it.

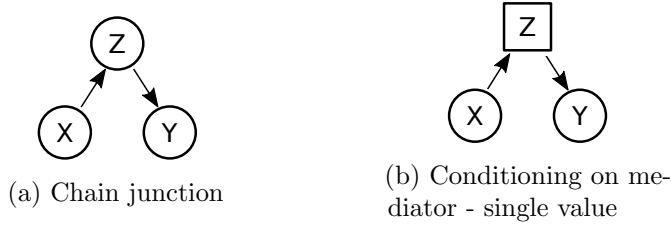


Figure 27: Unfortunate conditioning on mediator.

Conditioning on a mediator disables our ability to see causal relationships through that mediator.

Collider

Now we turn our attention to the *collider* from section 3.3. If we condition on the collider Z we conceptually make Z into a constant z , but we do so by carefully selecting the correct samples where $Z = z$. These samples are the ones where X and Y created the correct circumstances for which Z became z .

In order to understand colliders better we make a simple example with two coins; X and Y . We flip both coins and count the number of heads, Z . Clearly $z \in \{0, 1, 2\}$. If we know X and Y then we know Z exactly. If we know that coin X is heads ($x = 1$), then $z \in \{1, 2\}$, while if X is tails ($x = 0$), then $z \in \{0, 1\}$. Knowing x gives us no way of knowing y unless we also know z .

When the X -coin is known, then the remaining variability in Z comes from Y .

5.3. Causes, Mediators and Colliders

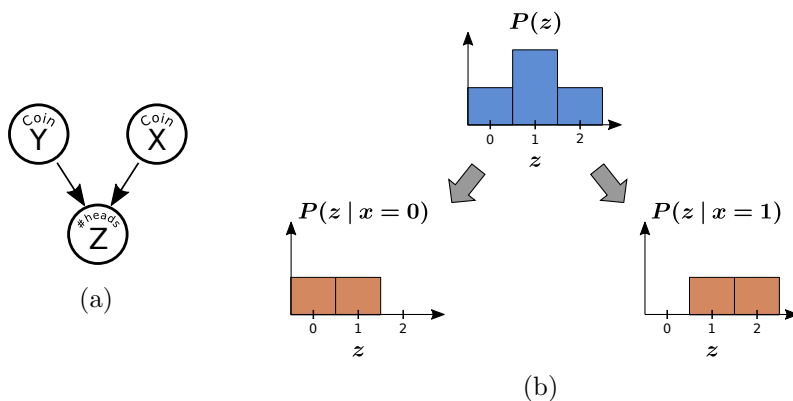


Figure 28: Double coin flip example.

If we intervene on Z then all causal relationship in the collider are eliminated. No information flows anywhere. X and Y can not predict Z and Z can not help us determine X and Y .

But what happens when we condition on Z ? If we condition the number of heads to $Z = 1$, then we can not determine X and Y , but we know that $X \neq Y$, so knowing one gives us the other. If we condition on $Z = 1$ it will look like X and Y are *inversely proportional*, because every time one of them is 1 the other will be 0. It thus looks like there is a causal relationship between X and Y , which is clearly wrong!

Our conclusion is therefore

*Conditioning on a collider can **create** spurious correlations.*

We will always want to avoid conditioning on a collider.

5.4 Conditioning on Children

There is another problem we have to consider when conditioning on variables. We return to the problem of flipping two coins, but we add a variable which is an exact copy of Y .

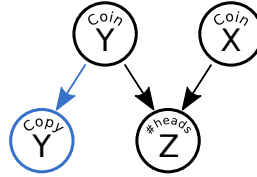
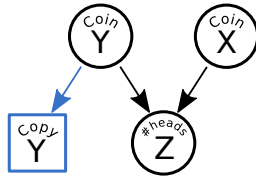
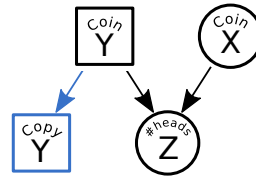


Figure 29: Another case of coin-flipping.
New variable is highlighted in blue.

We still want to determine the causal relationship between X , Y and Z . Let's see what happens if we condition on the copy of Y . We expect to conceptually make Y -copy a constant, but since this variable is identical to Y , then in practise we end up also conditioning on Y !



(a) What we think we get from conditioning on the copy of Y .



(b) What we actually get from conditioning on the copy of Y .

Thus we have ruined our experiment, because we cannot determine Z 's dependence of Y if we keep Y constant.

Let's make another experiment where we introduce a die D ⁸. We also introduce a "topscore"-variable T for X and D . This variable is 1 if we get heads with X and a 6 with D . Otherwise it is 0.

What happens if we condition on $T = 0$? Everytime $X = 0$ then nothing really happens, the relationship between Y , X and Z remains the same.

⁸If you didn't know: "dice" is plural :)

5.4. Conditioning on Children

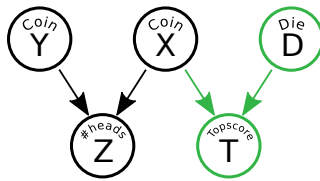


Figure 31: A die has been introduced.

But if $X = 1$ then there is a chance that we will remove that sample. It will therefore look like $X = 0$ more often than $X = 1$, which is not the case. So by conditioning on T we have *partly* conditioned on X , which is not what we wanted.

Conditioning on a child of a node X can partly or entirely condition on X as well. We generally want to avoid conditioning on children of the variables we are investigating

For the coin-and-dice examples these concepts may seem trivial. But consider the situation where we are testing a new drug. Perhaps we want to disregard various types of blood-pressure, sugar-levels etc. and consider how the medicine works conditioned on these. But if any of these are actually *causal children* of the drug (for example if the drug affect blood-pressure), then we end up partly conditioning on the drug itself, which can make us draw incorrect conclusions.

Conditioning can be used to solve confounder problems, but should be handled with care. Avoid conditioning on a node you are investigating, or a descendent of a node you are investigating, or on a collider.

Causal Inference - In General

6.1 Inference

Induction inference (here simply called inference) is the process of learning from evidence. It is the opposite process of deduction, in which a learned model is applied to information to produce a conclusion.

If we know the atmospheric pressure and we know how our barometer works, then we can *deduce* the distribution of the barometer state. This seems obvious, because the only thing that influences the barometer is the pressure - knowing the pressure allows us to know the barometer state.

Usually we use barometers the other way around though. We *infer* the distribution of the pressure from looking at the barometer. We know the barometer may not be exact, but we believe the actual atmospheric pressure will be pretty close to what is shown on the barometer.

Inference is governed by Bayes Theorem, which underpins all machine learning⁹

Bayes Theorem:
$$P(x | y) = \frac{p(y | x) p(x)}{p(y)}$$

When we apply Bayes theorem we use a known probability distribution $P(y | x)$, the *likelihood*, together with some additional information $P(x)$ (the *prior*) and $P(y)$, to compute $P(x | y)$, the *posterior* probability. When $P(x | y)$ is computed from $P(y | x)$, we sometimes refer to $P(y | x)$ as *forward*

⁹In fact it underpins ALL learning!

6. CAUSAL INFERENCE - IN GENERAL

probabilities and to $P(x | y)$ as *inverse probabilities*.

When we do causal inference we make use of all our tools from probability theory, but we need more than that. The causal tools we need are interventions, randomized interventions and conditioning. We will use these to determine *interventional distributions* - if we know the interventional distributions, then we know the causal structure. From section 3.2 we know that we need to handle confounders as well as disambiguation in order to infer causal relationships. Furthermore, we need to make sure we do not deteriorate our causal structure with the tools.

If we are for example interested in the causal relationship between X and Y in the graph on the right, how do we determine what tools to use?

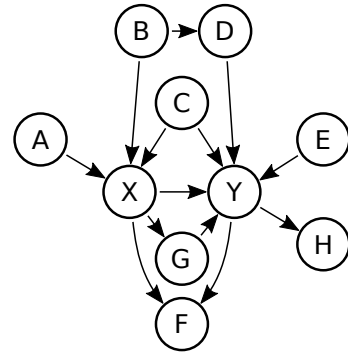


Figure 32: A big causal graph.

6.2 Information and Causal Graphs

One way of analysing causal graphs is from the perspective of information flow. If we know the value of a cause, then that tells us something about the effect from forward probabilities. Similarly if we know the effect then we know something about the cause from inverse probabilities.

If Figure 33 we have revisited the causal graph of the barometer problem.

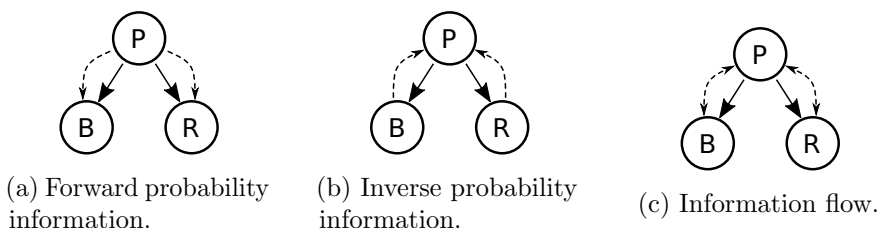


Figure 33: Information flow in the causal model governing atmospheric pressure (P), barometer reading (B) and rainfall (R).

The nodes are still causal variables and the solid links remain causal relationships. Figure 33a shows how we can use the forward probability distributions $P(\text{barometer} \mid \text{pressure})$ and $P(\text{rain fall} \mid \text{pressure})$ to deliver information from pressure to the other variables. Figure 33b shows how we can use the inverse probability distributions $P(\text{pressure} \mid \text{barometer})$ and $P(\text{pressure} \mid \text{rain fall})$ to deliver information from rainfall and the barometer to atmospheric pressure. The total flow of information is therefore as shown in 33c.

We can even distribute information through the pressure node as well - which confirms our first experiments where we noticed that our barometer and the rainfall shared some information.

Thinking about information "flowing" through a causal graph can be a useful way of understanding a problem.

6.3 Causal Patterns - Revisited

We will now revisit the causal patterns defined in section 3.3 and discuss how information flows through them.

Chain Junction

For the *chain junction* seen in Figure 34a we can use forward probabilities to propagate information from X to Z and onwards to Y . We can use inverse probabilities to infer information in the other direction, and so information runs freely through this graph (as shown in Figure 34b).

If we intervene or condition on Z , then conceptually we make it "constant". If it is constant then it carries no information. This will therefore block all information through Z .

Chain junctions lets information through, but can be blocked by intervention or conditioning.

Confounding Junction

For the *confounding junction / fork* seen in Figure 34c we can use inverse probabilities to infer information about Z from X and/or Y . If we can infer information about Z from X , then we can use forward probabilities to get information on Y - and vice versa. Therefore information also runs freely through this type of arrangement (Figure 34d).

Similarly to the chain junction, if we intervene or condition on Z we block of all information through Z . This is what we used when we conditioned on a confounder to help us determine a causal relationship.

Confounding junctions / forks lets information through, but can be blocked by intervention or conditioning.

6.3. Causal Patterns - Revisited

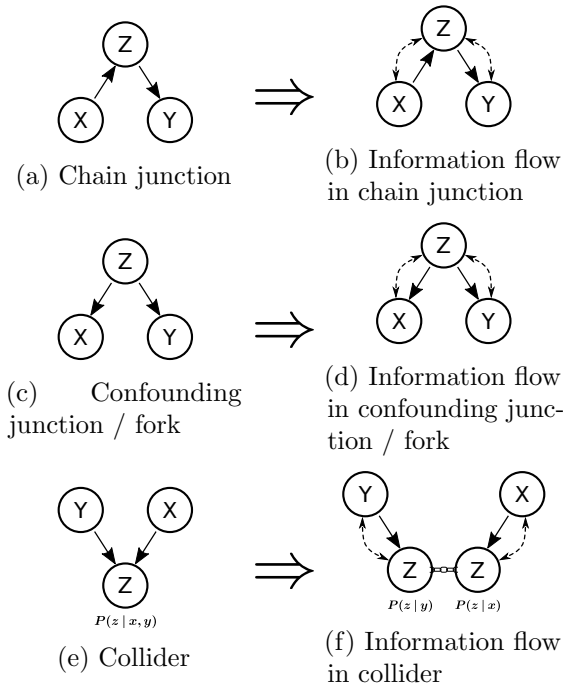


Figure 34: Three important arrangements of causal variables.

Collider

We now get to the third arrangement in Figure 34e, the *collider*, which is a bit more tricky. If we know X and/or Y we can say something about Z . If we know Z we can say something about both X and Y . The catch is that we can no longer use X to say anything about Y ! The colliding pattern *blocks* the information flow across the node.

One way to think of this is that the collider node holds information from both sides of the graph. If we know that $X = x$, then Z will follow the distribution $P(z | x)$, where Y has been marginalized because we do not know it. $P(z | x)$ therefore holds all the variance that Z gets from Y .

It is like the Z -node is split for each side of the graph, where one side holds the variability of being caused by Y and the other the variability of being

caused by X . In Figure 34f we have illustrated this by splitting the Z -node into two parts chained together. The left node follows the distribution $P(z | y)$, because it is on the side where Y is known but X is unknown and therefore marginalized. The right side similarly follows $P(z | x)$. If we know Z then we know both of these chained nodes and can therefore propagate information to both X and Y .

One peculiar thing about colliders is what we saw in section 5.3, conditioning on a collider then *opens up information*!

Colliders block information, but are opened up if you condition on them. They still block if intervened upon.

6.4 Solving a Graph

Using these ideas of information flow, we can solve most causal graphs! Where solve means determining how to find causal relationships. For example consider again the graph on the right where we wish to determine the causal relationship between X and Y .

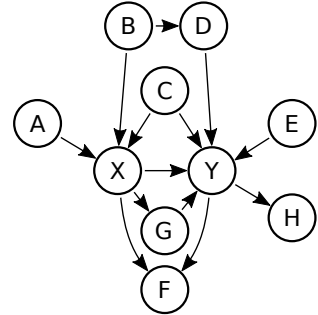


Figure 35: A difficult problem.

There are a lot of nodes, but let us draw all paths between X and Y where information is carried. The nodes A , E and H definitely do not carry any information between X and Y . Furthermore we know that colliders like F blocks information, so F is also not a problem. Below we have sloppily drawn the information paths between X and Y in blue

There are four paths. The BD -path and the C -path are confounders, which we need to handle. The remaining paths are the direct path and the path through G . We want to know how X causes Y and so both of these paths

6.4. Solving a Graph

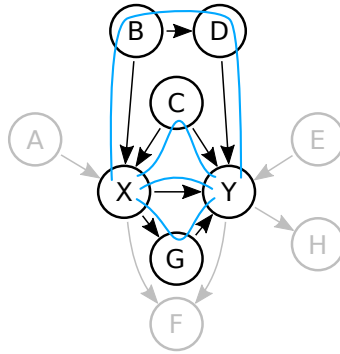


Figure 36: Information flow in the difficult problem.

are needed. If we want to get rid of the confounders we know that we can either condition on all of them or use intervention.

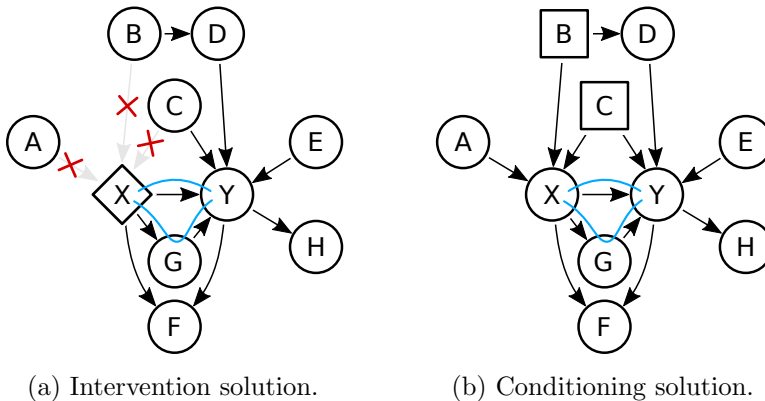


Figure 37: Solutions to the difficult problem.

If we made an experiment in which we intervened on X we could therefore test for the causal relationship between X and Y . We could also test for the relationships by gathering an observational dataset and conditioning on B and C . We can also conclude that data on A , E and H is not important for the problem. We now know what kind of experiment we should make, what data we should gather and how we should analyse it afterwards.

6. CAUSAL INFERENCE - IN GENERAL

Using intervention, randomized intervention and conditioning we can block off information in causal graphs, so that the only information flow remaining is the one we wish to measure.

Knowledge and Causation

In this section we discuss some aspects of causation which may not be intuitive at first.

7.1 Who Intervened?

If we want to infer causal relationships we will generally be making interventional experiments. But what if someone else makes an intervention? If another person makes an interventional experiment and hands you all data and information, then we might as well use those instead of doing our own experiment - no harm done. But what if you get your hands on a dataset without knowing that some of the data is interventional? First of all we saw in section 4 that not having interventional samples put you at a disadvantage. Also you may believe that all the data is observational (because you have no reason to believe otherwise) and your conclusions will probably be wrong. The data does not represent the system you are analysing, because you are analysing a system without interventions.

Let us say you try to increase employees rating of the work environment, by buying a new coffee machine. At first you ask the employees to rate each day for one month. Then you install a new coffee machine to see if the next month's ratings are higher. Hopefully you will see an improvement and your experiment is a success.

What if someone else decided to celebrate the company's anniversary, by buying Starbucks for the whole crew for the first month of your experiment? In this case you have a system where two groups are both intervening

7. KNOWLEDGE AND CAUSATION

on the same variable; coffee quality. If your new coffee machine does not make better coffee than Starbucks, then you might get worse ratings in the second month, which completely opposes your hypothesis.

*You need to **know** that data is interventional or observational.*

If you knew about the anniversary and the Starbucks coffee, then you could have planned your experiment differently. You could have collected ratings during the Starbucks-month and used that as the coffee-intervention - even though "you didn't do it". Afterwards you can decide whether a new coffee machine is worth it.

If you know about other's interventions you can utilize that information.

This illustrates that for interventions we also need to *know* that an intervention is going on. While this example is contrived there are many interesting multiagent systems (such as for example a country's economy) for which experiments are always messy and difficult to perform in isolation. Systems with many actors who all intervene differently at the same time, and each tries to learn from their own experiments.

7.2. No Free Lunch Theorem

A famous saying in science is by Isaac Newton *"If I have seen further it is by standing on the shoulders of Giants."*¹⁰ He contributes some of his scientific progress to the people before him, who taught him, wrote books for him to read, and in general gave him the knowledge needed for him to do his work. This has an interesting perspective to causality and intervention. Scientists like Isaac Newton make interventional experiments in order to infer causal relationships about our world. They did that by settings up experiments, intervening on systems, collecting data and analysing this data. When we today base our understanding, science and work on their findings, we use their conclusions as our assumptions. Therefore

Other's interventional data becomes our causal knowledge.

7.2 No Free Lunch Theorem

The *No Free Lunch Theorem*¹¹ can be found in multiple versions, but is here defined as

No Free Lunch Theorem: Any two optimization algorithms are equivalent when their performance is averaged across all possible problems.

This theorem makes a very important foundation for machine learning and AI. It states that for an algorithm that works well on a set of problems, we can always find a similar set of problems with opposite characteristics, for which the algorithm works bad.

The most important consequence of the No Free Lunch Theorem is that we have to make some assumptions in order to learn something. We cannot learn everything only from data - we can not get the perfect learner for free. Luckily for most problems we are interested in there are some shared characteristics which we can use as assumptions. Thus we may (perhaps)

¹⁰Although "on the shoulders of giants" can be traced back to before Newton as well.

¹¹David H. Wolpert and William G. Macready. No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation*, 1997

be able to come up with algorithms that works well for all problems of our interest.

In causal inference we also have to make some assumptions. As noted in section 2.2, whenever we specify a causal graph we are automatically assuming that everything else in the world does not matter. We restrict our focus to only be about the variables of our model. Similarly in a causal model we make a set of links, which implements our assumptions about the relationships between the variables. Again all links that we do not include are automatically assumed non-existing, which could potentially be wrong.

Therefore in order to do any sort of causal inference we also have to make some assumptions, similarly to the No Free Lunch Theorem of learning. No matter how much and how good our data is, at some point we need to make some assumptions, which currently means we need a human to think about the problem.

Causal inference requires making assumptions. Causal graphs can visualize those assumptions.

7.3 Ladder of Causation

The Ladder of Causation is a conceptual model by Judea Pearl¹², which distinguishes three types of reasoning based on the questions they are able to answer.

We have illustrated the Ladder of Causation in figure 38.

The lowest level of causation is called the association level (seeing) - the level of modern AI and of most statistical models. It is able to observe data and make conclusions like, *if this happened, will this other thing also happen?* It is all about learning probability distributions and conditioning

¹²Judea Pearl and Dana Mackenzie. *The Book of Why*. Basic Books, 2018.

7.3. Ladder of Causation

on what you see. Animals (and humans) do this all the time; if there is usually food here, there will likely be food again.

The second level of causation is the level of intervention (doing), which is what we have been dealing with in this note. Intervention allows the learner to learn causal models, which in turn can answer questions like *what happens if I do this?* The association level only conditions on observational data, while the association level conditions on interventional data.

The next level of causation concerns imagining and is what section 10 will be dealing with. It considers questions like *what if X had been different?*, in which we *imagine* a case where X is a different value than what it is. It also concerns questions like *why did this happen?* What *caused* something to happen. It allows models to explain themselves and explain decisions.

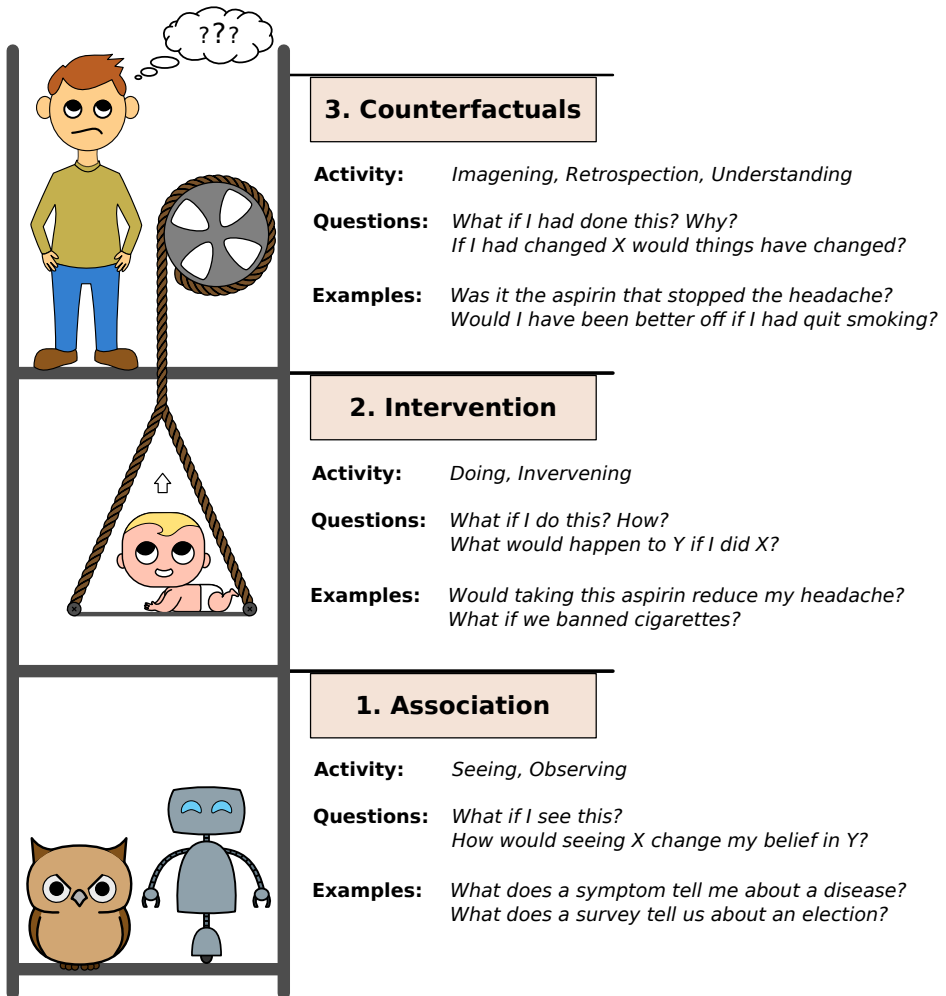


Figure 38: The Ladder of Causation.

Causality and AI

This section is to ensure an overall understanding of mathematical models (specifically machine learning) and how they relate to causal models. It is useful to solidify these concepts before advancing our causal models.

8.1 Mathematics, Science and Machine Learning

The exact definitions of mathematics and science is a controversial topic, but here we discuss some general properties of the fields and how machine learning relates.

In science we attempt to learn about the physical universe. We do so by making models (theories/hypotheses) of how we believe the universe works. We then *test* these models in the universe through experiments, which is the cornerstone of science. With experiments we control systems we want to analyse, we intervene on variables we want to examine, we collect data on outcomes, and we use the data to verify/falsify the models. Models that are verified we keep and improve, while models that are falsified we stop using. This makes science *empirical*, as we use experiments and data to decide what to believe. Furthermore scientific models are *causal*. They help us predict future events, while also explain their causes.

Mathematics is fundamentally quite different. In mathematics we *prove* our ideas. We do not rely on empirical experiments and data in order to verify/falsify ideas. Newly proposed ideas remain *conjectures* until they

are proven or disproven. If disproven we disregard them, while if we prove them we consider them theorems or lemmas. Mathematics is therefore not empirical and also not inherently causal.

While mathematics is not build on experiments, it is common to make experiments to get a sense of what might be provable. If you have a conjecture that you believe to be true, then you could experiment by testing it a ton of times and see if it always works. If it does always work, then it seem plausible that you may be able to prove it. Still only when proven do we consider it a mathematical theorem/lemma.

Furthermore, mathematics is perhaps the most useful "tool" know to humanity. As we prove more and more theorems in mathematics they become tools that we can use in science and all kinds of other fields. Therefore basically all experiments will be *using* mathematics, but not *extending* mathematics.

One field of mathematics that is crucial to science is statistics. In statistics we prove properties of statistical distributions/models. We use these models as hypotheses for how data behaves. If we assume a distribution/model for a data process, we can make provable conclusions about that process. Also we can utilize proven tests for determining whether a distribution/model is plausibly *correct* for the given problem.

In machine learning we have given up on finding the *true* model of the world. We have come to realize that it may be too time-consuming to correctly model every detail of the world. Instead we make very flexible models, which can hypothetically learn anything, and then *train* them for their purpose using data. We therefore make the models do the inference for us, which saves us a lot of time. This is why some defined machine learning as "improving from experience without being explicitly programmed". We do not prove the correctness of the model as we do in mathematics. Machine learning could thus be considered *empirical mathematics*, where we make mathematical models, but verify them empirically.

8.2. Models and Physics - an example

*Science is empirical and mathematics is provable.
Machine learning could be considered a mix - empirical
mathematics.*

In science *we* (people) make *explicit causal models* about the universe. In machine learning *computers* makes *automatic models* for *practical purposes*. The formalisation of causality in this document introduces ideas which could enable machine learning to also make causal models. In this case machine learning will be able to make automated science - it can help us discover causal models in huge datasets and problems too complicated for us to grasp.

8.2 Models and Physics - an example

Consider the following scenario. We have a set of rubber balls that we roll from a starting point out onto the floor of a gym. The balls have the same weight, but vary in size, and we roll them with varying angle and force.

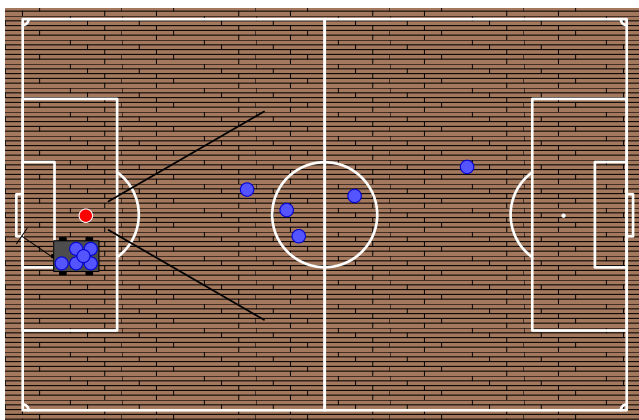


Figure 39: We are at the red point rolling out blue rubber balls across the floor of the gym.

We now collect a dataset, where we provide the initial speed, the size, the angle of the roll, as well as the final position of each ball we roll. We train

a fancy machine learning model, say a neural network, to predict the final position in (x, y) -coordinates, based on the other information. Sometimes the angle of the ball will deviate a bit and the floor friction will not be completely constant, so the balls are not completely deterministic. Still, after 200 throws the neural network predicts the final position with very high precision.

We also write out a formula for the final position, based on our understanding of physics. Using the teachings of physics for rigid bodies, we can derive such a formula. We only need to determine a couple of constants for the friction of the floor, air resistance etc. which we can fit using our dataset. This formula is also capable of predicting the final position of the balls quite precisely and so we now have two different models to predict positions: a machine learning model and a physical model.

Let's get another type of ball which is of the same varying sizes, but weighs more. Weighing more will give the ball more initial, kinetic energy for a given initial speed.

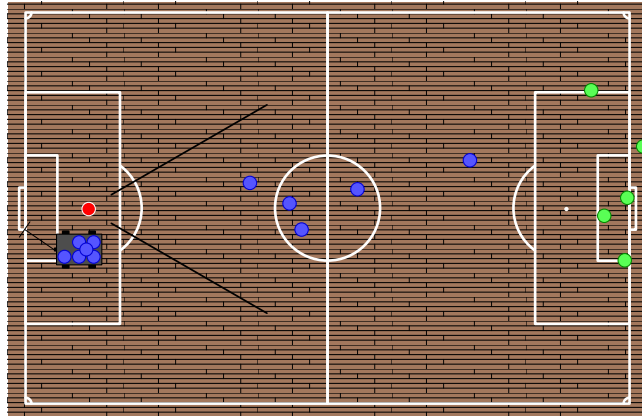


Figure 40: We now also have green balls.

If we directly predict with both models, they both predict incorrectly - we need to fix these models.

With the physical model we can easily incorporate the new weight. With

8.2. Models and Physics - an example

almost no effort we can make the physical model handle the new situation and predict correctly again! What do we do with the neural network? We could try to explain to it that the balls have changed weight, but how exactly do we do that? We would probably have to roll quite a few of the new balls, in order to get data to retrain the neural network.

How about if we go outside and roll balls across a parking with asphalt.

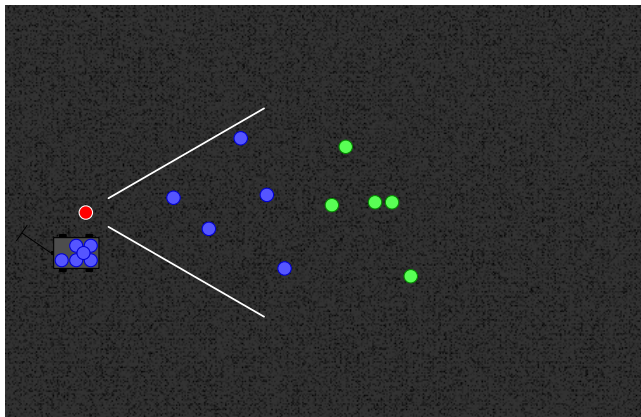


Figure 41: Experiment on asphalt.

Again both models are off, because the friction of the asphalt is completely different from that in the gym. We could get more data for refitting the two models. Alternatively we could look up information about friction of asphalt and apply to the physical model. Using this approach we are again able to manipulate the physical model to handle new circumstances, without getting new data. And perhaps even more interestingly, we have used *others experience* (the friction info on asphalt) for using in our own problem.

So what is the big difference between two models? Well, the first thing to note is that we broke the most important rule of machine learning: the test data *must* come from the same distribution as the training data! This assumption underpins everything machine learning builds on, and we broke it when we tested the models on different systems than what we started with. But physical models *can* help us in new situations. With physical models we can even compute how systems behave on Mars, even though

the "training data" we used to learn about physics is not from Mars. The second thing to note is that technically the neural network is wrong. No matter how great neural networks are, we have to face the fact that a ball rolling on a surface *are not neural networks*. They are mechanical systems following the laws of physics. In fact, the best we can hope for when using the neural network, is that it *learns the physical model*. The situation in which the neural network will be most precise, is the situation where it contains the physical model of the balls. The balls *must follow the laws of physics*.

The physical model is a causal model. It does not simply try to model the distribution of final positions, but actually models the physics of the system, and if our understanding of physics is correct, then a ball rolling across a surface *is that modelled system*. The causal model is capable adapting to interventions in the system, as well as answer other questions like, why did the ball not roll very far? - because of a higher friction of asphalt.

8.3 Machine Learning Recap

Here we make a quick recap of models used in machine learning (and other branches too). It's useful for expanding to more detailed causal models. While not all previous sections are needed to understand this section, it is a good idea to read sections 8.1 and 8.2. For the rest of this section we use the following notation: x is inputs, y is a real-valued output and c is a discrete-valued output. $P(c)$ denotes a probability of discrete event c , while $p(y)$ denotes the probability-density of real-value y .

In AI and ML we commonly distinguish between two types of problems:

Classification: Models the relationship from inputs to a **discrete**-valued output.

Regression: Models the relationship from inputs to a **real**-valued output.

Basic Models

The most basic version of these two methods chooses the most probable/-expected value of the problem. The *pure discriminative model* determine the most probable class for an input

$$\arg \max_c P(c | x). \quad (1)$$

The model makes discrete decisions, but does not tell us about any possible alternative. The basic regression model computes the expectation of the output variable based on the inputs

$$\mathbb{E}[y | x] = \int y \cdot p(y | x) dx. \quad (2)$$

This model provides one good estimate of the output, but does not model the confidence of the prediction.

Figure 42 illustrates the basic models. The classifier divides the space into parts according to class and the regression model fits a line through the data.

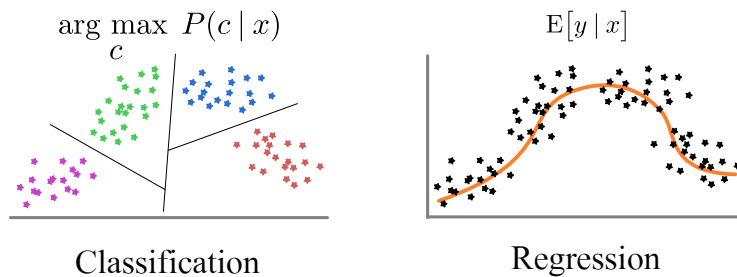


Figure 42: Basic models

Probabilistic Models

If we take more probabilistic approach to classification get a *conditional discriminative model*, which models the probability mass of all classes

$$P(c | x). \quad (3)$$

Using this type of model we can easily provide the same predictions as with the pure discriminative model, but for each sample we also know if other classes are also possible, as well as how likely each class is.

Similarly by upgrading the regression model we get the conditional density of the output

$$p(y | x). \quad (4)$$

Again we can use this model to select the most probable value, but we can also say something about the much the output may vary from this estimate.

The probabilistic models are illustrated in Figure 43. The classifier makes boundaries between classes with soft transitions to illustrates areas of uncertainty. Notice that it is quite arbitrary what the classifier decides to do in regions with no data (the centre part). The regression model again fits a line, but this time with some notion of uncertainty.

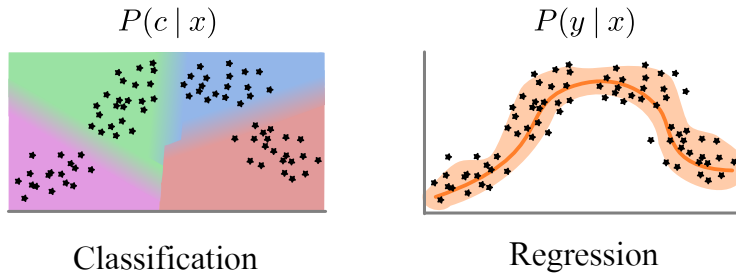


Figure 43: Probabilistic models

Knowing the confidence and variance of outputs is crucial for many tasks, and honestly there are very few reasons to settle for a basic model today, as we can usually make probabilistic versions.

Generative Models

We now get to the final level of statistical models

Generative Model: A generative model models both inputs and outputs - it models the whole system. Crucially a generative model is capable of generating new samples.

8.3. Machine Learning Recap

The generative model provides the joint probability of inputs and output

$$P(c, x) \text{ or } p(y, x). \quad (5)$$

We can use this model to provide conditional probabilities (for both classes and regression targets) and it becomes a very natural starting point for designing conditional discriminative models. Since we model the joint distribution, the difference between inputs and outputs become less relevant. The difference between regression and classification also becomes blurry as a regression model with a discrete input could be identical to a classification model.

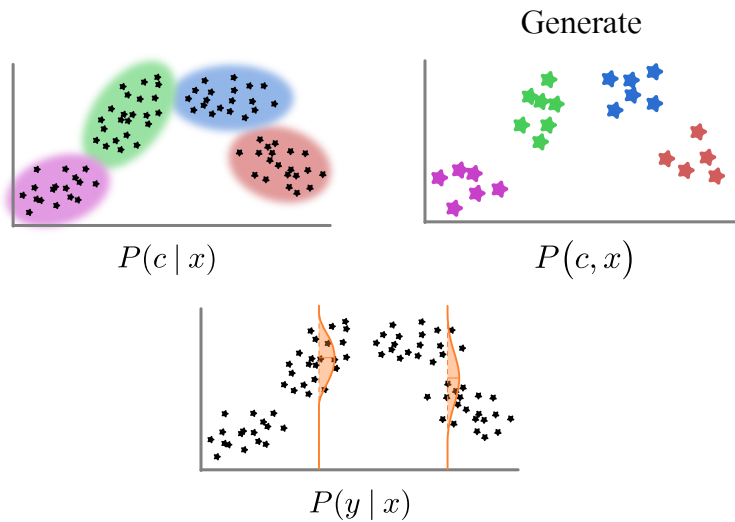


Figure 44: Generative models

Figure 44 illustrates the generative models. On the left we classify by computing the class-conditional probabilities. We are also able to determine areas with no information (centre region), as we know where the input distributions are limited. In a sense we can do outlier detection, and be warned if data is from a region for which the model has limited information. On the right we generate data using the joint distribution. At the bottom we condition on the input data, which provides a continuous distribution for regression purposes.

Full Causal Models



Up until now we have considered probability distributions for various problems and augmented them with graphical causal models. In this section we are going to expand these tools and make the first direct connection to artificial intelligence (AI) and machine learning (ML).

9.1 Generative Process Models

We now upgrade our generative models.

Generative Process Models: A generative model which models the *process* of the data and not only the distribution.

The difference between a generative model and a generative process model is a bit vague and weird, but in nutshell there can be many different generative models that all correctly models a problem, but there is only one generative process model. The generative process model is the *correct* generative model, which is correct even when we intervene and which can extrapolate to new systems. We can therefore use such a model for causality.

Figure 45 illustrates the usages of the generative process model and table 1 recaps the taxonomy of models.

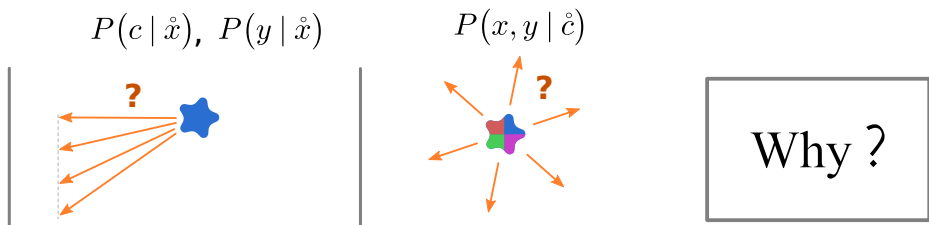


Figure 45: Generative process models. We can use these models to ask what happens when the input is intervened on (left), when the class is intervened on (center), as well as *why* a particular outcome was observed.

| Level | Model type | Modelled probability/density | Questions and usages |
|-------|--------------------|---|---|
| 1 | Basic | $\arg \max_c P(c x), \mathbb{E}[y x]$ | Single prediction. |
| 2 | Probabilistic | $P(c x), p(y x)$ | Prediction with uncertainty. |
| 3 | Generative | $P(c, x), p(y, x)$ | Generation and outlier warning. |
| 4 | Generative Process | $P^*(c, x), p^*(y, x)$ | What happens if we force x/c ? What if x/c had been different? Why did this happen? |

Table 1: Taxonomy of models. We use $*$ to indicate the "right" generative model, which is the generative process model.

9.2 Structural Equation Models

In causal literature we use a specific type of generative process model called a Structured Equation Model - SEM. An SEM is a set of ordered assignments of the following form

$$\begin{aligned} a &\leftarrow f_a() \\ b &\leftarrow f_b(a) \\ c &\leftarrow f_c(a, b) \\ &\vdots \end{aligned}$$

The model first determines a from some stochastic function $f_a()$. It then computes b using another stochastic function $f_b(a)$, which uses a as input. The formula for c uses a and b and so on. Note that while each variable has access to the previous ones, they do not *have* to use them (they may simply ignore some of them). The ordered set of assignments creates dependencies where each variable depends on earlier ones.

Let's take an example

$$\begin{aligned} A &\leftarrow N_A, & N_A &\sim \mathcal{N}(0, 1) \\ B &\leftarrow \frac{1}{2} \cdot A + N_B, & N_B &\sim \mathcal{N}(0, 1) \\ C &\leftarrow B + N_C, & N_C &\sim \mathcal{N}(0, 1) \\ D &\leftarrow B \cdot C + \frac{1}{3} \cdot N_D, & N_D &\sim \mathcal{N}(0, 1). \end{aligned}$$

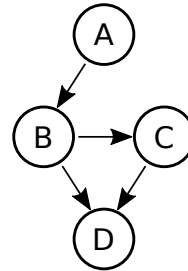


Figure 46:
Example graph.

We see that A depends on no other variable. B depends on A , C depends on B , and D depends on both B and C . Furthermore the stochasticity of the variables is represented by a stochastic variable for each (which in this example are all normally distributed).

For any SEM, we can graph the dependencies of the assignments. If the

SEM models a causal system, then this graph is the causal graph of the system - back to things we know about! If we graph the above example we get the graph in figure 46, where we see the expected dependencies.

A causal graph is a simplified SEM, where we know which parameters depend on which, but we do not know *how* they depend. With the SEM we know all the details.

Using the above SEM with a programming language with sampling libraries (such as Python) we can sample A - we just need to draw samples from a normal distribution. Similarly we can sample from B by first sampling from A and computing B with samples from N_B 's normal distribution. We can continue to do this and sample the whole process - we have a generative process model.

We can use the causal graph to do all the analyses we have previously discussed, but can also go a bit deeper. When asked how B depends on A we can actually specify *how much* B depends on A , because we can compute how much of B 's variance that comes from A (one third). Similarly we can say that two thirds of B 's variance is *independent* of A . SEM's gives a much more detailed view of the causal process.

Counterfactuals



10.1 A Counterfactual Question

We will now discuss a different type of causal question; "what if"-questions.

Counterfactual: A counterfactual question is one where we seek information about how the world would have been different given some specific change. For example; in a situation where X *did* happen, what if X had *not* happened?

For discussing counterfactuals we start with an example. At 23:40 on the 14th of April 1912, the RMS Titanic hit an iceberg which ultimately sank the ship after 2 hours and 40 minutes. The iceberg was spotted by Frederick Fleet, but only too late for the ship to navigate around it. Let us consider the following question; *if the RMS Titanic had started its voyage one day later, would it have made the crossing?*

At first we may get some information on the event. We know the date of the iceberg collision, and could probably get information about the weather and the preceding winter, perhaps even estimate the amount of ice in the waters based on the tales from the survivors. Let's gather all this information and put it into a variable U called utility information. Furthermore we know that the Titanic left from Southampton on the 10th of April and that it did sink; $D = 10$ and S . Let's first try to frame the problem mathematically. We want to determine the probability that the ship did not sink, given that it left the next day, while also conditioning

on the known information U as well as the fact that it left on the 10th and it did sink

$$P(\neg S \mid D = 11, U = u, D = 10, S). \quad (1)$$

Here we see the first obvious problem, which is what makes counterfactuals quite different. We are searching the probability of *not sinking* $\neg S$, given that we know that *the ship sank* S . Furthermore we are conditioning on the date being the 11th, $D = 11$, while also conditioning on the known fact that the ship did actually leave on the 10th, $D = 10$.

We first need to get rid of this confusion. Some variables hold information that we know. These are just like any other stochastic variable and we keep them as is. The other variables hold *counterfactual* information. They either hold the query that we wish to know about (in our case $\neg S$) or causal changes that we wish to investigate (in our case $D = 11$). We will here mark all counterfactual variables by a black dot \bullet , so the expression becomes

$$P(\neg \dot{S} \mid D \dot{\leftarrow} 11, U = u, D = 10, S). \quad (2)$$

A common interpretation of counterfactuals is that we consider a parallel universe, where everything is *exactly* the same, except for the counterfactual variables and everything affected by those variables. We therefore ask; in a parallel universe where U (all auxiliary information) is the same, what is the probability of $\neg \dot{S}$ given $D \dot{\leftarrow} 11$, when we know that in *our* universe we observed $D = 10$ and S .

10.2 Party Time!

The Titanic is an example of how important and relevant counterfactual questions can be. In order to learn about this topic though, we start with a simpler example; a party.

Bob is not too fond of parties but really likes Alice. If Alice is at the party, then Bob usually goes unless they are having a disagreement. Carl really likes partying, but Alice is his ex-girlfriend and sometimes he will avoid the party if she is there. Finally Bob and Carl cannot stand each other and if they are both at a party, they usually fight. Sometimes they even fight without being at a party.

We wrap this knowledge into the following causal model, where we have also specified the probabilities. A denotes Alice being at the party, B means Bob is at the party, C Carl and F means that there is a fight.

| value | if | $P(N_x)$ |
|--|-----------|----------|
| $A \leftarrow \begin{cases} 0 \\ 1 \end{cases}$ | $N_A = 0$ | 0.40 |
| | $N_A = 1$ | 0.60 |
| $B \leftarrow \begin{cases} 0 \\ A \\ 1 - A \end{cases}$ | $N_B = 0$ | 0.07 |
| | $N_B = 1$ | 0.90 |
| | $N_B = 2$ | 0.03 |
| $C \leftarrow \begin{cases} 0 \\ 1 - A \\ 1 \end{cases}$ | $N_C = 0$ | 0.05 |
| | $N_C = 1$ | 0.85 |
| | $N_C = 2$ | 0.10 |
| $F \leftarrow \begin{cases} 0 \\ C \cdot B \\ C \cdot B + (1 - C) \cdot (1 - B) \end{cases}$ | $N_F = 0$ | 0.05 |
| | $N_F = 1$ | 0.90 |
| | $N_F = 2$ | 0.05 |

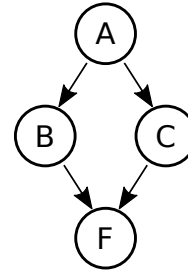


Figure 47:
Causal graph
for the party.

We have good intuitive explanations for each of the variables. N_A contains the variance of whether Alice attends the party. If $N_B = 0$ then Bob can't attend the party, if $N_B = 1$ then Bob will go if Alice goes, and if $N_B = 2$ then Bob will do the opposite of Alice, because they are having a

disagreement that night. With Carls attendance we see that Carl can't go to the party when $N_C = 0$, will go no matter what if $N_C = 2$, and might go but will avoid Alice if $N_C = 1$. N_F basically sets the mood of the party - either there is no fight, or Carl and Bob will fight if they are at the party, or they will fight if they are either both at the party or both somewhere else.

Before we start analysing this problem we note that we can expand on the causal graph in a useful way. In the graph on the right we have also made nodes for the noise/variance variables. Each noise variable is hidden because we cannot observe it. For example we can check whether Alice is at the party, but we do not know what specific factors made her come to the party. Furthermore we have modelled the problem in a way where all variance of each variable is collected in these variables, when conditioned on ancestors: if we know that Alice is at the party then all remaining variance of Bob's attendance is in N_B .

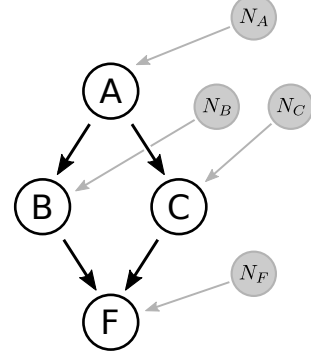


Figure 48: Causal graph for the party with noise/variance variables.

A Counterfactual Party

Last weekend there was a big party which Bob did not attend. We wish to determine the probability of a fight if Bob *had* gone

$$P(\dot{F} \mid \dot{B}, \neg B). \quad (3)$$

Here we have noted the only known information; Bob was not at the party $\neg B$. We are thinking of a parallel world where we make Bob go \dot{B} and see if there is a fight \dot{F} .

First we make a new causal graph where we duplicate all the nodes: one for the real universe and one for the counterfactual universe.

All observable nodes of the two universes have their own nodes, because we can make analyses where we inspect each side of the variables. Conversely we cannot observe any of the hidden variables, but we assume that everything we do not touch is identical for the two universes, so they only have one node each.

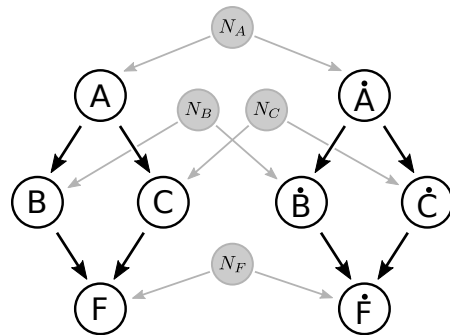


Figure 49: Causal graph for the party with counterfactual variables.

Now we can return to the query for $P(\dot{F} \mid \dot{B}, \neg B)$. We know that $B = 0$ in our universe. In the counterfactual we *intervene* to make $B \leftarrow 1$. Finally in our situation A , F and C are hidden. The causal graph that we are now analysing looks like the one in figure 50.

We have one constant node, one interventional node, which we forced ourselves, and the remaining are hidden nodes. N_B is no longer relevant as both its children are known/constant or intervened on.

We wish to determine the probability $P(\dot{F} \mid \dot{B}, \neg B)$, for which the variables A and C are marginalized. We therefore need sum the probability of F ,

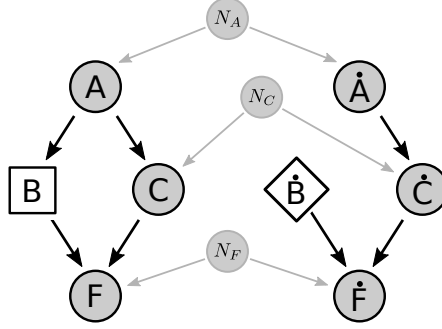


Figure 50: Counterfactual query graph for the party.

given we force \dot{B} , and conditioned on all possible values of A and C . The computation becomes

$$P(\dot{F} \mid \dot{B}, \neg B) = \sum_{a=0}^1 \sum_{c=0}^1 P(\dot{F} \mid \dot{B}, C=c, A=a) P(C=c, A=a \mid \neg B). \quad (4)$$

The fight depends on Carl, and Carl depends on Alice, so let's start with Alice. We compute the probability of Alice being at the party with inverse probabilities (Bayes Theorem)

$$P(A \mid \neg B) = \frac{P(\neg B \mid A) P(A)}{P(\neg B \mid A) P(A) + P(\neg B \mid \neg A) P(\neg A)} \quad (5)$$

$$= \frac{(0.07 + 0.03) \cdot 0.6}{(0.07 + 0.03) \cdot 0.6 + (0.07 + 0.90) \cdot 0.4} \approx 0.13, \quad (6)$$

$$P(\neg A \mid \neg B) \approx 1 - 0.093 \approx 0.87 \quad (7)$$

So Alice is most likely not at the party, because if she were then Bob would most likely also have been there. We have illustrated the flow of information in Figure 51a. We know the constant B and we know the distribution of N_A , which we can use to infer the distribution of A .

Now that we know the distribution of Alice, we can determine the distribution of Carl. In Figure 51b we use the information about A and

10.2. Party Time!

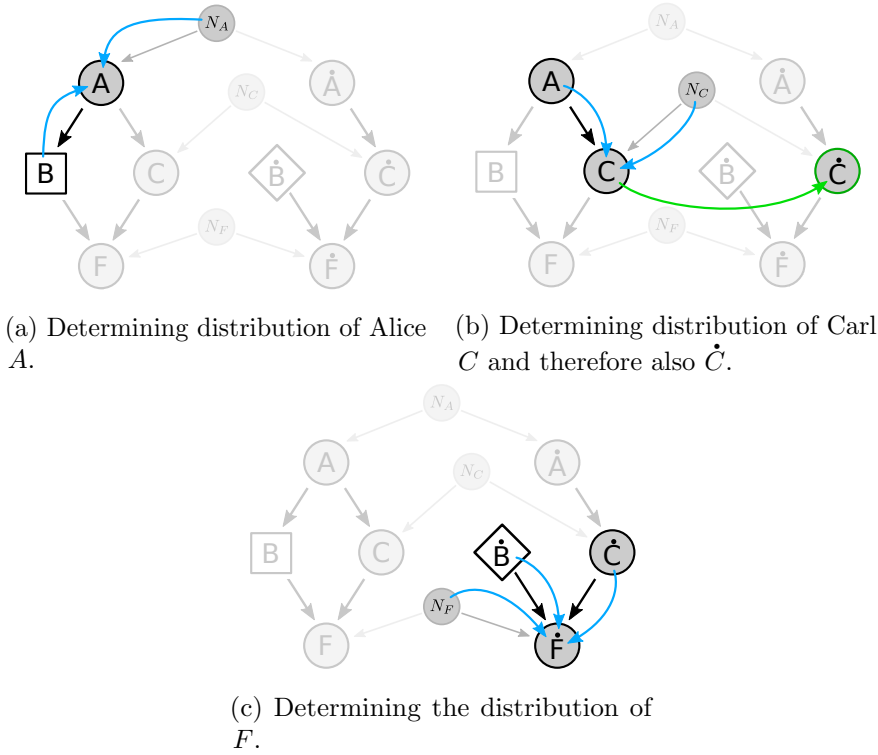


Figure 51: Flow of information for determining the counterfactual probability $P(\dot{F} \mid \dot{B}, \neg B)$.

distribution of N_C to infer the distribution of C . This is done by first getting the conditionals

$$P(C \mid A) = 0.10 \quad P(C \mid \neg A) = 0.95 \quad (8)$$

and then computing the probability of Carls attendance

$$P(C \mid \neg B) = \sum_{a=0}^1 P(C \mid A = a) P(A = a \mid \neg B) \quad (9)$$

$$\approx 0.10 \cdot 0.13 + 0.95 \cdot 0.87 \approx 0.84, \quad (10)$$

$$P(\neg C \mid \neg B) \approx 1 - 0.84 \approx 0.16. \quad (11)$$

We see that Carl is most likely at the party if Bob is not. Also, since we know that $C = \dot{C}$, we can transfer the knowledge about C to its counterfactual twin in the other universe.

We now know the distribution of \dot{C} , and we already knew the distribution of N_F and the forced \dot{B} . This allows us to finally consider the probability of a fight \dot{F} . Computing the counterfactual probability of a fight is seen in Figure 51c and is done by

$$P(\dot{F} \mid \dot{B}, \neg B) \quad (12)$$

$$= P(\dot{F} \mid \dot{B}, C) P(C \mid \neg B) + P(\dot{F} \mid \dot{B}, \neg C) P(\neg C \mid \neg B) \quad (13)$$

$$= 0.95 \cdot 0.84 + 0.00 \cdot 0.16 = 0.79. \quad (14)$$

We can therefore conclude, that if Bob was not at the party (which he wasn't), then there most likely would have been a fight between him and Carl, if we forced him to go.

10.3 Counterfactuals in General

We will now analyse the concept of counterfactual reasoning from a more generic point of view. Counterfactuals are all about hypothesising about the effect of making a specific change to a sample. In contrast, conditional and interventional distributions are about systems, and how they behave under conditioning and intervention. Counterfactuals are also about systems, but

Counterfactuals apply to specific samples.

If we gather all variables that we make counterfactual interventions on into X , and gather all variables that we are interested into Y , then a generic graph of the causal system is seen below¹³.

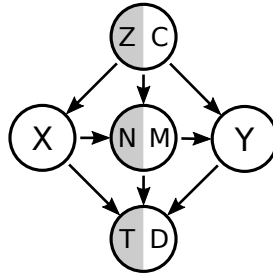


Figure 52: Generic causal graph investigating counterfactuals with X and Y .

We know that there is a causal relationships directed from X to Y . We have some mediators M between the two, some confounders C , and some descendants D . In counterfactuals it is quite crucial how much you know about the system. We have therefore split the confounders into visible confounders C and hidden confounders Z , split the mediators into visible M and hidden N , as well as split the descendants into visible D and hidden T . Furthermore note that it may be a direct link between X and Y and between Z/C and T/D , but we simply let these run through the mediator

¹³Again we assume no cycles.

node for simplicity.

Since we are doing counterfactuals we assume that we know the generative process model - so there is not ambiguity about cause and effect etc. and we now all relationships. Now we observe a sample $(x_0, y_0, c_0, m_0, d_0)$ and ask, what is the probability that we would have observed y_1 instead, given that we had intervened and forced x_1

$$P(\dot{Y} = y_1 \mid X \dot{\leftarrow} x_1, x_0, y_0, c_0, m_0, d_0). \quad (15)$$

First thing to note is that we keep all the known auxiliary information in place (c_0, m_0, d_0) and also condition on the true values of the counterfactual variables (x_0, y_0) . Knowing more about the system (for example moving variables from Z to C) will allow us to compute a more precise counterfactual probability, while losing information (moving variables from C to Z) will make the counterfactual probability less precise.

Hidden Variables

For predicting Y we need to figure out how the hidden variables might have behaved in the given situation. We thus inspect the distribution for Z , N and T . We know that they produced the sample $(x_0, y_0, c_0, m_0, d_0)$, and so their distributions should be conditioned on these

$$P(z, n, t \mid x_0, y_0, c_0, m_0, d_0) \quad (16)$$

which will involve forward probabilities and Bayes theorem for computing backwards probabilities. Also we are now performing a counterfactual intervention on \dot{X} which may change some of the hidden variables.

10.3. Counterfactuals in General

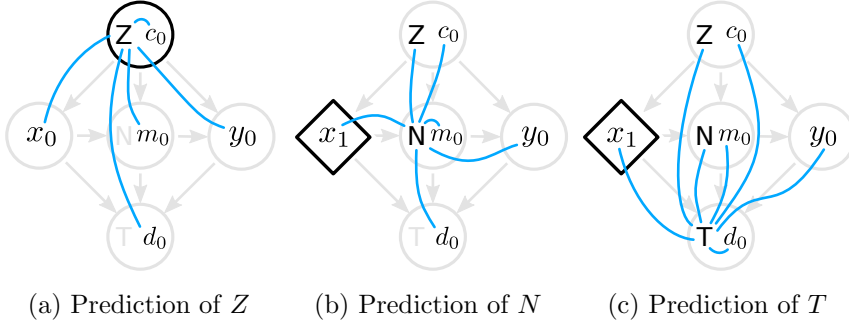


Figure 53: Information used for prediction of hidden variables.

We first determine the probability of Z - the highest, hidden node. This distribution is conditioned on $(x_0, y_0, c_0, m_0, d_0)$, because it help produce that outcome, which we illustrate in Figure 53a. It is not conditioned on \dot{x}_1 , because Z is not affected by X .

In Figure 53b we move down to N . Here we replace x_0 by \dot{x}_1 , because the value of N did not *cause* x_0 , but is itself *caused* by the new value \dot{x}_1 . The distribution of N is also conditioned on the value of its ancestor Z , whose distribution we have just determined, as well as all the other known information.

Finally in Figure 53c we consider the distribution of T . We condition the distribution of T on z, n, \dot{x} and k , because it is the descendant of all these variables. We let $k = (y_0, c_0, m_0, d_0)$ for reducing clutter and conclude that

$$P(z, n, t \mid \dot{x}_1, x_0, k) = \underbrace{P(t \mid \dot{x}_1, k, z, n)}_{\text{lowest node}} \underbrace{P(n \mid \dot{x}_1, k, z)}_{\text{center node}} \underbrace{P(z \mid x_0, k)}_{\text{highest node}}. \quad (17)$$

Prediction

If we know all variables, then the probability of y depends only on its ancestors (not D and T)

$$P(y \mid x, c, z, m, n, d, t) = P(y \mid x, c, z, m, n), \quad (18)$$

which can be computed directly from the generative process model.

We can now analyse the counterfactual probability in (15) by marginalizing over all unknown variables Z , N and T

$$P(\dot{Y} = y_1 \mid X \leftarrow x_1, x_0, y_0, c_0, m_0, d_0) \quad (19)$$

$$\begin{aligned} &= \sum_{znt} \underbrace{P(y_1 \mid \dot{x}_1, c_0, m_0, d_0, z, n, t)}_{\text{probability of } y_1 \text{ conditioned on all variables}} \underbrace{P(z, n, t \mid \dot{x}_1, x_0, y_0, c_0, m_0, d_0)}_{\text{probability of hidden variables}} \\ &= \sum_{znt} \underbrace{P(y_1 \mid \dot{x}_1, c_0, m_0, d_0, z, n, t)}_Y \underbrace{P(t \mid \dot{x}_1, k, z, n)}_T \underbrace{P(n \mid \dot{x}_1, k, z)}_N \underbrace{P(z \mid x_0, k)}_Z \quad [(17)] \\ &= \sum_{zn} \underbrace{P(y_1 \mid \dot{x}_1, c_0, m_0, z, n)}_Y \underbrace{P(n \mid \dot{x}_1, k, z)}_N \underbrace{P(z \mid x_0, k)}_Z \quad [(18)] \end{aligned}$$

$$= \sum_{zn} \underbrace{P(y_1 \mid \dot{x}_1, c_0, m_0, z, n)}_Y \underbrace{P(n \mid \dot{x}_1, y_0, c_0, m_0, d_0, z)}_N \underbrace{P(z \mid x_0, y_0, c_0, m_0, d_0)}_Z. \quad [k]$$

There are some important things to note here. First of all, the final computation of the probability of Y is not special. $P(y_1 \mid \dot{x}_1, c_0, m_0, z, n)$ is simply evaluating our model for y_1 conditioned on $(\dot{x}_1, c_0, m_0, z, n)$. Secondly the *hidden descendants* T are not used in any way and are irrelevant for the analysis. The only relevant descendants are the visible ones, because they provide information about the states of the other variables (so we included d_0 in k). Thirdly the ancestors of X uses x_0 , while the descendants of X uses \dot{x}_1 . This is because Z helped create x_0 , while the descendants are affected by the counterfactual intervention on X .

Variants of the Counterfactual Question

Say there is a fire in a building, which you have just left as part of the fire drill. You could ask "would the fire have been smaller if I had closed the basement window?" (do to less oxygen for the fire). It might be that someone else already closed it, in which case it would not have made a difference. If no one had the time to close that window, then it probably would have made a difference. This is a situation where X is hidden and we simply wish to know what the probability of Y is given that we *force* X .

Consider another situation where you just left the voting booth of an election. The results won't be released in several hours, but still you wonder whether it would have made a difference if you hadn't voted. This is a counterfactual question where you know X , but you don't know Y . Also you should always vote.

Thus it is possible to make counterfactuals where you don't know X and/or Y . In this case these variables won't appear in the conditionals of (19).

10.4 Counterfactuals and Interventions

In the previous example we knew $(x_0, y_0, c_0, m_0, d_0)$. Let's consider what happens if we know less information. First we eliminate x_0 and y_0 . These no longer appear in the conditionals in (19) and our probability become a bit more imprecise for the specific situation.

Now say we don't know any of the confounders (all of C is moved in to Z), which makes the probability more imprecise again. We can continue to move all of M in to N , and all of D into T . Now we do not know about any variables of the situation and the probability becomes

$$P(\dot{Y} = y_1 \mid X \leftarrow x_1) = \sum_{zn} P(y_1 \mid \dot{x}_1, z, n) P(n \mid \dot{x}_1, z) P(z) \quad (20)$$

$$= p(y_1 \mid X \overset{\circ}{\leftarrow} x_1), \quad (21)$$

which can be verified by inspecting (11). Thus the interventional distribution is the counterfactual distribution where we do not know anything! Or phrased differently

The interventional distribution is the expected counterfactual distribution over the observable population.

This makes intuitive sense: the interventional distribution evaluates the probability of an event given that we force on a variable, without knowing anything else about the system. A counterfactual can therefore be considered a conditional, interventional distribution, where we intervene while conditioning on a *specific sample*.

Perspectives on Causality

11.1 Why?

The word *why* is a bit overloaded. At the same time, it is the cornerstone of some of the most important scientific questions, and so important in causal analysis that Judea Pearl named his book on causality *The Book of Why*.

In order to specify what kinds of questions we will be concerned with, we first note the following three types of why-questions.

1. Which of these hypotheses made event X happen?
(why does X happen? - with hypotheses)
2. Why are we doing this?
3. Why does event X cause Y ? Why does X happen?
(without hypotheses)

The first question is also the simplest one. If a house has burned down, it is reasonable to ask why and investigate. This is a cause-and-effect question, where we use causal models to search for a likely cause for the effect. For example we know that faulty electrical systems is a common cause of fires, so if we find a device at the center of the fire, it may be a plausible explanation. The common methodology for solving this question is well understood using Bayesian approaches. Under multiple hypotheses for why there was a fire, we use probability theory to find the one that most precisely describes the outcome and evidence.

The second question is a confusing one, which commonly deteriorates some of the most heated debates. Usually because people do not agree what such a question means. It's questions like *why are we doing this?* and *why are we here?*. It concerns the *meaning* or *purpose* with an action, thing or person, which implies that someone else is giving it meaning or purpose. We are also going to ignore these types of questions.

This leaves the final question; why did X cause Y ? This is a causal question with great use and it is all about mediators and direct causal links.

Scurvy

Scurvy is a disease where early symptoms include weakness, feeling tired, sore arms and legs. If not cured the symptoms extend to gum disease, changes to hair, bleeding from the skin, poor wound healing and ultimately death from infection or bleeding. The disease plagued humanity for centuries and during the Age of Sail it was assumed that 50 percent of the sailors would die of scurvy on a given trip.

In 1753, a Scottish surgeon in the Royal Navy, James Lind, showed that scurvy could be successfully treated with citrus fruit. For a couple of decades few put his advice into practice, but by the 1800s the British Royal Navy routinely gave lemon juice to its sailors to avoid the disease.

While the discovery of curing scurvy with lemon juice was a tremendous and important achievement, it would still leave people perplexed, because it did not answer the question *why?* Why did a fruit cure a disease that had killed thousands if not millions of people over centuries? Why did people get scurvy?

These are causal questions, so let us use some of the tools we have learned about to analyse them. First of all why do people get scurvy? In this question we are interested in a causal variable *scurvy* S and we wish to know what *causes* it. In other words we want to find the ancestors of the variable

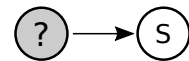


Figure 54: Causal graph for the question *why do we get scurvy?*

11.1. Why?

in whatever causal graph the variable operates in.

The second question also concerns scurvy S as well as the use of lemon juice L . We believe lemon juice prevents scurvy, but we do not know why. There could be different explanations. It could be that lemon juice is simply an essential part of human existence, in which case there is a direct link between lemon juice and scurvy. Alternatively it could be that lemon juice does something to our bodies or contains something which prevents scurvy, in which case there is some unknown mediator which we want to find.

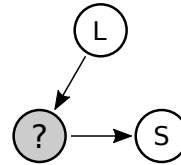
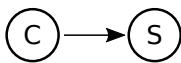
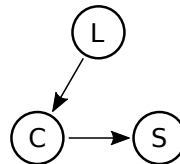


Figure 55: Causal graph for the question *why does lemon juice prevent scurvy?*

In 1912, Norwegian Axel Holst discovered that the cause of the cure in the citrus was vitamin C, and in 1928 Hungarian Albert von Szent-Györgyi was able to isolate vitamin C, which he was rewarded the Nobel prize for. This allows us to answer the causal questions! The reason why we get scurvy is because we lack vitamin C, which we need. The reason why lemon prevents scurvy is because it contains a lot of vitamin C which we can obtain.



(a) Answer to the question in figure 54



(b) Answer to the question in figure 55

Figure 56

This provides us with some answer for our causal questions about scurvy - but we might be more interested. For example why does the lack of vitamin C cause scurvy? Or why does it cause the phenomena which we label as scurvy? This can be seen as a further granulation of the question.

Now we could split the S -node into all the individual phenomena of scurvy and investigate any mediators between vitamin C and these new nodes. The new questions is about why we need vitamin C. What does vitamin C do? This way we can keep searching for more and more causal knowledge by increasing the detail of our causal graph, as illustrated in figure 57.

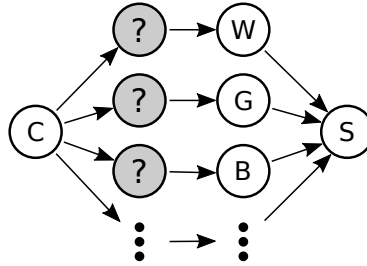


Figure 57: Expanded causal questioning of the scurvy problem. We are investigating weakness W , gum disease G , bleeding B and other phenomena of scurvy S . We wish to know why vitamin C affects (prevents) these phenomena.

But Really, Why?

Let us try to map some common why-questions to some of the terminology of causality.

If we want to know why something X happens then we are interested in the ancestors of X . If we find all ancestors of X , then we can say that X happens because the conditions of the ancestors are sufficient.

If we already know an ancestor Y of X , but we are unsure why the relationship exists, then we are searching for mediators between the two. If we find all mediators from Y to X , then we can say why Y causes X . In order to understand more and more about the process of Y and X we can increase the granularity of our causal model to find more and more mediators explaining each little bit of the process.

11.2. Explainability

Say a person makes decision X (assuming its a rational, voluntary decision). If we are pondering why the person did so we are searching for descendants. If we know all descendants of X , then we can answer the question of why a person might want to do X . We know all effects of X and some subset of those effects should be the rational reason for doing something.

We can also consider the last question retrospectively. Why would it have been a good idea to do X ? In this we are asking what the descendants of X would be if we counterfactually changed X .

11.2 Explainability

Say a banker needs to decide how good a loan they can make for a given person A . If person A is not satisfied with the conditions of the loan then the most common reaction would probably be to ask *why*. Why can I not get a good loan? Why is my credit rating lower than others? The banker will be able to answer these questions, because they have insight into the reasoning behind their decisions. Could for example be that person A has had a tough time paying back an earlier loan. Perhaps their monthly expenses are almost as big as their salary, making it irresponsible to provide a loan and thus more expenses.

In modern systems we are getting used to automated decision making provided either by human-made rule systems or artificial intelligence. In the latter case, explanations are becoming increasingly difficult to provide. Complicated, deep neural networks with millions of parameters are inherently difficult to understand and even if a professional AI-researcher managers to understand the workings of their model, how are they to explain the decision making process to laymen, who's expertise is not AI?

While we do not have an easy solution for this, we can make note of the properties of causal models. Causal models *do* provide explainability! If the model in an automated system is causal, then after a given decision we can start to analyse why the decision was made. What ancestors were important? From the example of the person taking a loan, the ancestors could be *previous loans* and *expenses/income balance*. Then we can dig

further into what the problem is with the two ancestors: the *previous loans* where not paid in time and the *expenses/income balance* was too close to 1. Thus in order to gain a loan person A should either convince the bank that paying back the loan will be different this time, that their expenses are being reduced, or that they will be earning a higher salary soon.

While producing causal models for all our automated systems is not an easy task, explainability is certainly a motivator for why we might want to solve that task.

11.3 Bias and Fairness

A hugely important and heavily discussed problem with automated systems is that if made negligently they can amplify or even create biased decision systems. For example in the case of providing loans to person A , one could ask; is the persons gender or ethnicity affecting their credit rating? A general goal for democratic societies is that we want equality between all citizens, regardless of sensitive properties such as gender, race, ethnicity, sexuality, religion etc. When designing automated decision systems, we therefore want to enforce some constraints on the system, enforcing our goals of equality and fairness.

Bias and fairness in artificial intelligence is a topic with a lot of attention and research, which we hope will see great progress in the next couple of years. An important stepping stone has been to mathematically define measures of fairness and bias in systems. This is not a trivial task and there are many approaches, but we here note that one of them involve causality.

Say we want to know whether the decision made for person A was discriminatory. If G is a good loan, we could ask the following counterfactual question; is

$$p(G \mid U, E) \neq p(\dot{G} \mid U, \dot{E}), \quad (1)$$

where U is a collection of used information and E is ethnicity. If the current probability of a good loan (left side) is equal to the counterfactual

11.4. Causality in Physics

probability of a good loan where we change the ethnicity of the person (right side), then it seems that the bank has not been discriminating. But the the probabilities are *not* equal, then ethnicity has an influence in their decision making of loans. We can therefore frame bias as decisions where we can counterfactually change the outcome by intervened on the sensitive variable that we are investigating.

Many methodologies are able to say whether there are biases in systems *on average*. While not having a bias on average is of course a good thing, but can there still be bias in individual cases? If there can, then the problem does not seem solved. One advantage of using the counterfactual approach is that it measures whether there was bias in a *specific case*.

11.4 Causality in Physics

This book discuss causality from a mathematical modelling point of view. If there exists a cause C and an effect E , then we can use our causal methods to do predictions, inference etc. Another question regarding causality is whether it applies to our universe; does the laws of physics include a law of causality? While I am in no way experienced in the field of physics, I will here attempt to describe some important concepts from physics regarding causality.

1. Definition

One thing to note is that physics does not provide much of a *definition* of causality, just like we did not in this book. Largely causality is defined similarly as here; Causality is the relationship between causes and effects.

2. Causality has historic roots in physics

Causality has deep roots in physics (and in other sciences), as we have often described physical systems as causes and effects. We accredit nuclear fusion in the sun with the light and warmth that makes life thrive on Earth. We say that the gravity and motion of the moon causes tides. We can

claim that a loud noise was caused by a tray of dishes falling to the floor. The natural sciences are very causal.

3. Causes for causes for causes ...

If a loud noise was caused by a falling tray of dishes, then we may say that the falling of the tray was caused by gravity. Furthermore we may use Einstein's theory of relativity to claim that the force of gravity is caused by the bending of space-time. But what causes matter to bend space-time? Can we keep making causes for causes? When does the causes of causes stop? We do not really know.

4. Can laws of physics attribute causal relationships?

There is also some dispute of whether the laws of physics actually explains causal relationships. Say we succeed in further developing our understanding of the laws of physics. We get to the point where we can predict virtually all physical phenomena using a specific set of equations. These equations may describe exactly where *every atom*¹⁴ is, how fast it moves etc. But still they may not say anything about what *causes* the equations. The equations just *work*, they always do, and we may never know what *causes* them to do so. This problem seems very related to the problem of working with an observational distribution, vs. an interventional distribution. We can do interventions within our world, while respecting the laws of physics, but per definition we can never intervene *on* the laws of physics. We can never *break* them to see what would happen. If we could break a law, then that law would not be universal and we would attempt replace it by a better law of physics. Thus our experience of the physical laws is observational and not interventional.

5. Does all things have a cause?

Here physics have managed to solve a problem from philosophy. A previous theory was that the universe is *deterministic*. This means we could potentially predict the future state of the universe if we knew everything about the universe as it is now. Due to the modern theories of quantum

¹⁴ or quark

11.4. Causality in Physics

physics we now know this is not the case. There are physical events that happen at random and they can not be predicted. Thus they don't really have a cause. One example is the radioactive decay of a single atom (an atom which is unstable and breaks into smaller atoms and other particles). We can make statistical models of how long we expect to wait before the atom will decay, but it is impossible to predict when a particular atom will decay, regardless of how long the atom has existed. While we have now sorted out that not all processes have a cause it does not make the whole discussion much simpler :)

6. Relation to time

One thing physics helps us with is to set up some requirements for causes and effects; a cause must precede an effect in *time*.

This requirement seems like common sense, but due to Einstein's theory of special relativity time is not as easy to operate with as we may perceive in our everyday lives.

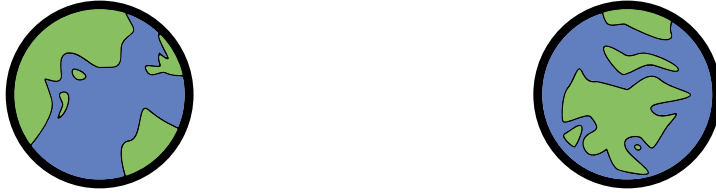
Say we are in a spaceship¹⁵. On our starboard side there is a particle *A* really far away. Similarly on our port side there is a particle *B*, which is equally far away. Say that all three objects are stationary in space and that both particles decay at the exact same "time". From our spaceship we will perceive exactly that; we will observe two particles decaying at the same time. Let's instead say that we moved really fast towards particle *A*. Due to special relativity time seems to move slower at particle *B* and faster at particle *A* from our perspective. Thus it seems like particle *A* decayed *before* particle *B*. If we were travelling the other way, then it would seem that particle *B* decayed *before* particle *A*! Thus for two different observers, events can seem to have opposite order in time depending of relative velocities. This makes the notion of time difficult. If a cause must happen before an effect, then what do we do if the order of events can be interchanged? For example if event *X* caused event *Y* is it possible for one observer to observe event *X* before *Y*?

Luckily physics have an answer to this problem. Consider first two planets

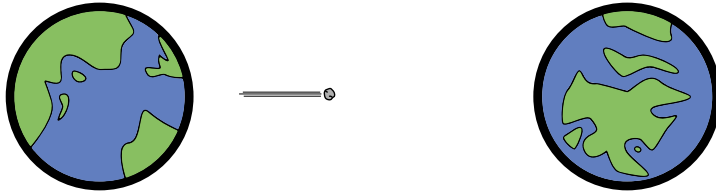
¹⁵ in space :)

11. PERSPECTIVES ON CAUSALITY

like below.



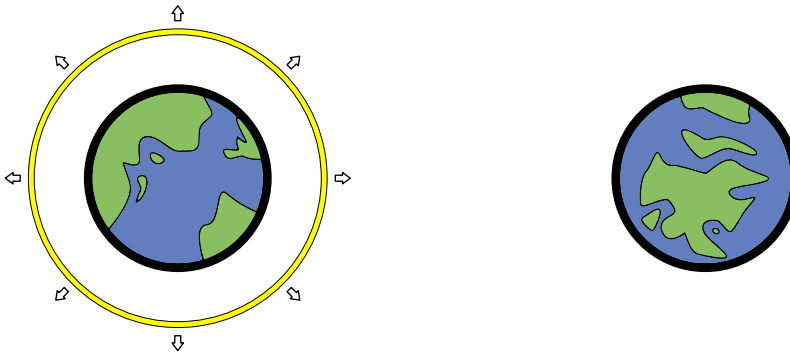
Now planet 1 on the left does not like planet 2 (surprisingly planet 2 really like planet 1, but I guess they are just nicer there). Therefore planet 1 throws a fairly big rock at planet 2.



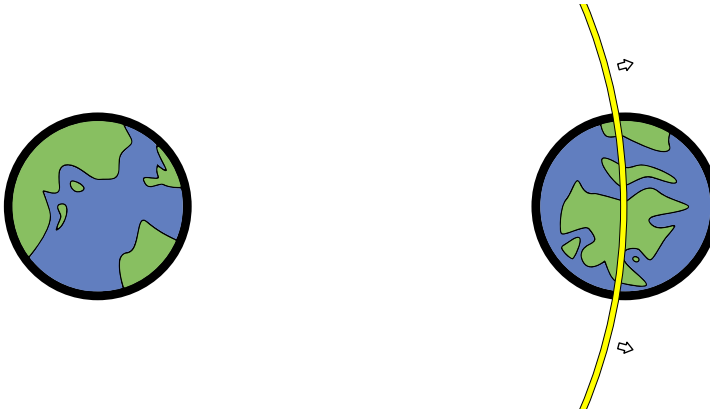
Planet 2 is certainly not going to like getting rocks thrown at them and it should cause a reaction. So planet 1 cause the rock to fly through space, and the rock later causes planet 2 to be surprised. Furthermore planet 2 will only be surprised when the rock gets there - thus there is a time delay between planet 1 throwing a rock and planet 2 being surprised. But what if planet 2 saw the rock coming? That would make them surprised earlier. This raises the question; how quickly can an event on planet 1 affect planet 2? Could planet 2 be affected instantly by an event on planet 1?

This is where the laws of physics help us. Say that instead of throwing a rock, planet 1 decides to send a message, "You guys are idiots", to planet 2 using light. They know that light travels at the fastest possible speed through space, so this must be the fastest way to annoy planet 2. The light travels outward from planet 1, creating a sphere sending the message in all directions.

11.4. Causality in Physics



At this moment the message has been sent, but planet 2 has not received it yet. Furthermore since the light travels faster than anything else in the universe it is impossible for planet 2 to realize that the message is coming beforehand. Even if they have astronauts in space, who sees the message first, they are unable to send any faster message to the planet to warn them. Only after the sphere of light reaches planet 2 is it possible for planet 2 to realize that a message was sent from planet 1.



Therefore it is impossible for planet 1 to do anything that affects planet 2 faster than the speed of light. This puts a restriction on causality: a cause and effect relationship can only occur at the speed of light.

In physics we describe this phenomenon with the concept of a *light-cone*. Say for a moment that the world is 2-dimensional. Figure 58a illustrates

11. PERSPECTIVES ON CAUSALITY

this world in orange, with an observer O in the center with a light-source. In the third dimension we now plot *time*. As time progresses (going up) light expands in circles around O . If O causes event A to happen, then event A must be within this *future light-cone* of O , because otherwise A would be so far away from O that light could not reach A in due time (and thus nothing else could either).

Similarly if event B causes something to happen at O , then B must be within the *past light-cone* (going down). If B is not within the past light-cone, then B is too far away to affect O at this time-instance.

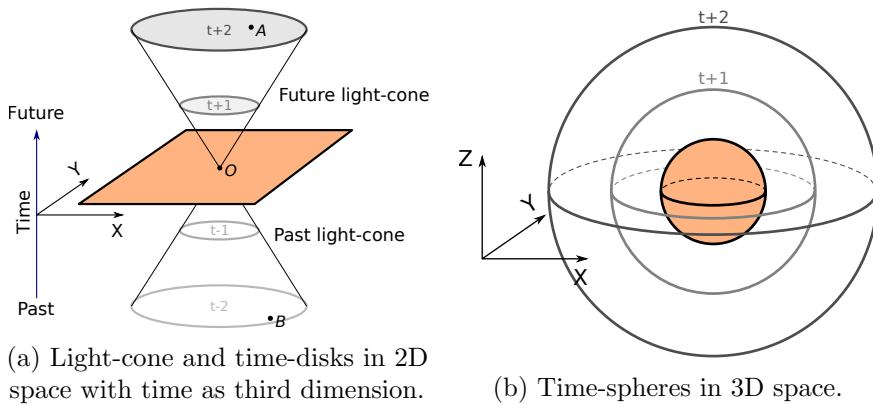


Figure 58

We can also describe a light-cone for 3-dimensional space, but the cone itself is tricky to draw because we have run out of dimensions. Instead we consider each time-instance. In 2D the light moved in circles away from the light source. In 3D light moves as a sphere. In figure 58b we have an observer in orange. If the observer sends out light, then this light will move as an expanding sphere (just like the message from planet 1). Anything that the observer affects must be within these spheres in order to be affected in due time.

While the constraint that causes and effects has to be within each others past and future light cones (must be within a certain distance to be affected in time), is a fairly loose restriction, it is one of the few absolute laws of causality that we can make.

11.5 Do-notation

This book has introduced you to some concepts of causality and use our own notation. There are great books and a lot of interesting research of causality out there, and they tend to use different notation.

Our Notation: In our notation we write \dot{x} to say that we intervene and force event x . We can specify a specific value which we enforce on x as $x \overset{\circ}{\leftarrow} 3$. We can also specify changes in probability distributions, for example used for randomized trials, $p_x \overset{\circ}{\leftarrow} \mathcal{N}(\mu, \sigma)$. Furthermore we can specify interventions which creates a forced constraint $y \overset{\circ}{\leftarrow} x$, in which y is set to follow x , but not the other way around. We use a very similar notation for counterfactuals, where \hat{x} denotes a counterfactual intervention of event x and $x \overset{\circ}{\leftarrow} 3$ denotes a counterfactual intervention of setting x to 3.

Pearl's Notation: Judea Pearl has been a key figure in the development of causal theory and coined the term do-notation. In Pearl's notation $do(x)$ indicates that we intervene and force event x , and $do(x = 3)$ indicates an intervention of setting x to a specific value. Pearl has often noted counterfactual variables with a tick or a hat, as in \hat{x} or x' , to note that the variable has been counterfactually changed.

Peters et al.: The book *Elements of Causal Inference* by Jonas Peters, Dominik Janzing, and Bernhard Schölkopf uses a very principled notation. They wish to underline the difference between conditioning and intervening, and maintain directionality in do-notation. $P_Y^{\mathfrak{C}; do(X:=3)}$ is the interventional distribution of variable Y , when we intervene to make X equal to 3 on the causal model \mathfrak{C} . First thing to note is that their notation includes the causal model \mathfrak{C} which we are assuming to produce the probability distribution. Furthermore $X := 3$ is directional allowing for expressions like $X := Y$ without any doubt of direction.

A counterfactual distribution is noted by $P_Y^{\mathfrak{C}|Z=2, X=1; do(X:=3)}$. Notice that $|$ still denotes conditioning, but $; do$ denotes interventions.

11. PERSPECTIVES ON CAUSALITY

The arguments for our notation is

- ① The notation is directional like causality is
- ② The notation aligns with assignment from computer science and intervention can be viewed as physical assignment
- ③ The notation allows for a clear and simple, link and distinction between intervention and counterfactual intervention
- ④ The notation avoids unnecessary bracketing and cluttering

Exercises

Python Code

The small codebase following this document is organized as follows

```
python/  
  document/  
  exercises/  
  project/  
  src/
```

The **src**-directory contains various methods and classes used for exercises and figures, but which should not be edited. In the **exercises**-directory there are scripts which can be run and edited to solve the exercises. The **document**-directory and **project**-directory has code related to figures in this document and to the project (page 155).

All python scripts should be run from the **python**-directory (the topmost directory), in order to access **src**-directory etc. For example say we wish to run the **switch_experiment.py**-file for exercise 1.4. We can do this by opening a terminal and moving to the **python**-directory.

On Windows we then call

```
      Run from here  
    ┌───────────┐  
SET PYTHONPATH=.  
python exercises\ex_1.4_switch\switch_experiment.py  
└──┬──┬──────────────────────────────────────────┘  
  run run this script  
python
```

while on Mac and Linux (bash) we call

```
PYTHONPATH=./ python exercises/ex_1.4_switch/switch_experiment.py  
└──┬──┬──────────────────────────────────────────┘  
  Run from here run this script  
python
```

If you use an IDE like Spider or PyCharm, you should simply ensure to add the **python**-directory as a source-location or as source/content-root.

1st Exercise

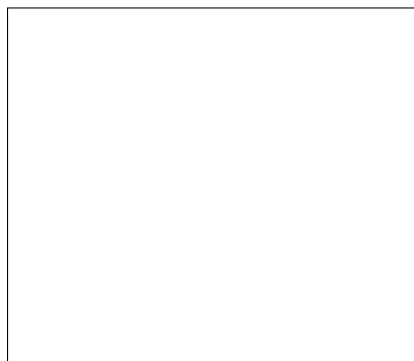
The exercises here concern the topics covered in sections 1-6.

1.1 Problems and Graphs

For each of the descriptions below, draw the related causal graph.

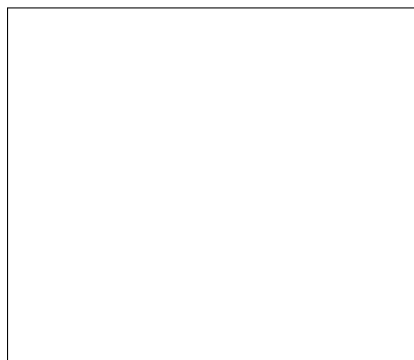
A Small Graph

J is a mediator between ancestor H and descendant G . P is a confounder of H and G , while Q is a child of J .



Ideal Gas Law

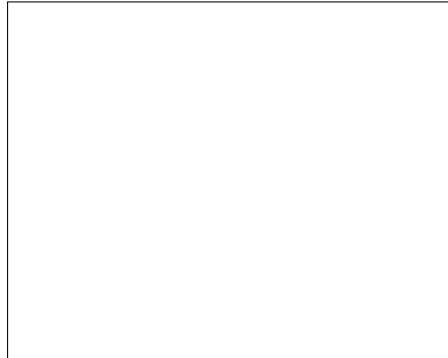
We have a nitrogen tank used for cooling with. We wish to know what its pressure is. We know various specs on the tank and can also do some measurements on it, but we don't know what we need. Do a Wikipedia search on the *Ideal Gas Law* and determine what factors we need to know about the tank, in order to determine the pressure inside the tank. Make a causal diagram illustrating this.



Relaxing is Dangerous?

A study from 2005 found that

... embarking on the Golden Years at age 55 doubled the risk for death before reaching age 65, compared with those who toiled beyond age 60.



That is; people who retired early seemed to die early¹⁶.

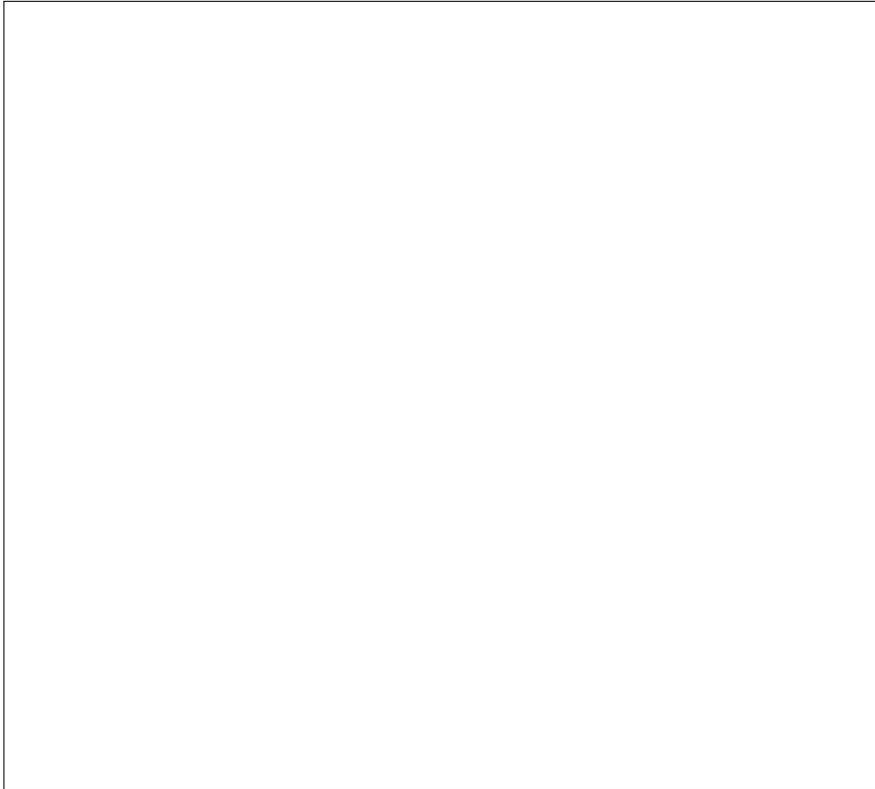
Make a causal graph of how you think the following variables interacts: death (D), current age (A), health-condition (H), job type (J), retirement age (R), personal economy (E). Assume everyone in the data already is retired.

¹⁶Luckily, a couple of later studies found the opposite to be true:
<https://www.nytimes.com/2018/01/29/upshot/early-retirement-longevity-health-wellness.html>

A Big Graph

In the following we enumerate the English alphabet (including W) in the normal order, starting with $a = 1$.

F is the parent of J , who is also a child of H . The 5th prime is parent of the 6th prime, who is a parent of the 7th prime, who is a parent of the first prime. All letters whose numbers are odd and greater than 2^4 are parents of the last letter in the alphabet. Letter 5^2 is the child of letter 2^3 . The vowels between letters "j" and "x" are children of J and parents of 2^4 . The only letter whose number is a divisor of 289 is a child of node $(2^4) - 1$.



1.2 A Paradox

Initial Analysis

We are going to analyse some numbers on the fall 1973 admissions to the University of California, Berkeley. We have a hypothesis that the university has a bias against the admission of women! Consider the table below and compute the percentages of admission.

| | Applicants | Admitted | Admitted % |
|-------|------------|----------|------------|
| Men | 2590 | 1269 | |
| Women | 1835 | 556 | |

Is there a bias?

Extended Analysis

We have some more data on the admission rates for each of 6 departments. We wish to find out which departments are most problematic. Fill out the percentages in table 1 to find any bias within the departments. Let's find all departments where the men/women admissions differs by more than 5%. That is, ignore small differences like 50% and 52%, but consider a difference like 50% and 55% a bias. Note the biases in the tick-boxes on the right of the table (bias against women/no-one/men).

| Depart. | Men | | | | Women | | | | Bias against W/-/M |
|---------|-------|--------|--------|---|-------|--------|--------|---|--|
| | Appl. | Admit. | Admit. | % | Appl. | Admit. | Admit. | % | |
| A | 825 | 511 | | | 108 | 88 | | | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> |
| B | 560 | 353 | | | 25 | 17 | | | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> |
| C | 325 | 120 | | | 593 | 202 | | | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> |
| D | 417 | 138 | | | 375 | 131 | | | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> |
| E | 191 | 53 | | | 393 | 94 | | | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> |
| F | 272 | 16 | | | 341 | 24 | | | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> |
| Total | 2590 | 1269 | | | 1835 | 556 | | | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> |

Table 1

Is there a bias? Who is the admissions biased against?

Causal Analysis

In order to figure out this problem we will analyse it using the causal tools we have learned. Hypothesis 1 was that there is a causal link from gender G to admission rate A - and no other variables are considered. In Hypothesis 2 we also consider department D . What could the causal links be between G , A and D ? Consider the graph for both Hypothesis 1 and 2 and draw them in Figure 71.



(a) Hypothesis 1



(b) Hypothesis 2

Figure 59

Can you use these graphs to determine whether there is a bias against women?

Simpson's Paradox

This is a case of *Simpson's Paradox*, which is perhaps the most famous statistical paradox. The paradox happens when a confounder (the department) inverts the analysed relationship (gender causes admission) in the aggregated data (total admissions), as opposed to within the subgroups. The true relationship is the one found in the subgroup, because the confounder no longer has influence.

| | Treatment A | Treatment B |
|--------------|-----------------|-----------------|
| Small stones | 93% (81 / 87) | 87% (234 / 270) |
| Large stones | 73% (192 / 263) | 69% (55 / 80) |
| Total | 78% (273 / 350) | 83% (289 / 350) |

Kidney stone experiment.

Another famous example is from a medical study comparing the success rates of two treatments for kidney stones, whose data is shown on the right. For the total, Treatment B seems best, but this is because Treatment B got more of the easy cases (small stones) and fewer of the difficult cases (large stones). In reality Treatment A is superior.

When we have information and data on the confounder we can eliminate the problem with conditioning. The problem really arises when we can't access the data or perhaps don't even know about the confounder. If we have no information about the size of the kidney stones, then the only possible conclusion is - incorrectly - that Treatment B is best for treating kidney stones. This is where random trials have their strengths. By randomizing who gets which treatment, A or B, we can ensure that no confounder has influenced the decision. This is the point of randomization in A/B testing.

The Berkeley and Kidney Stones examples are from

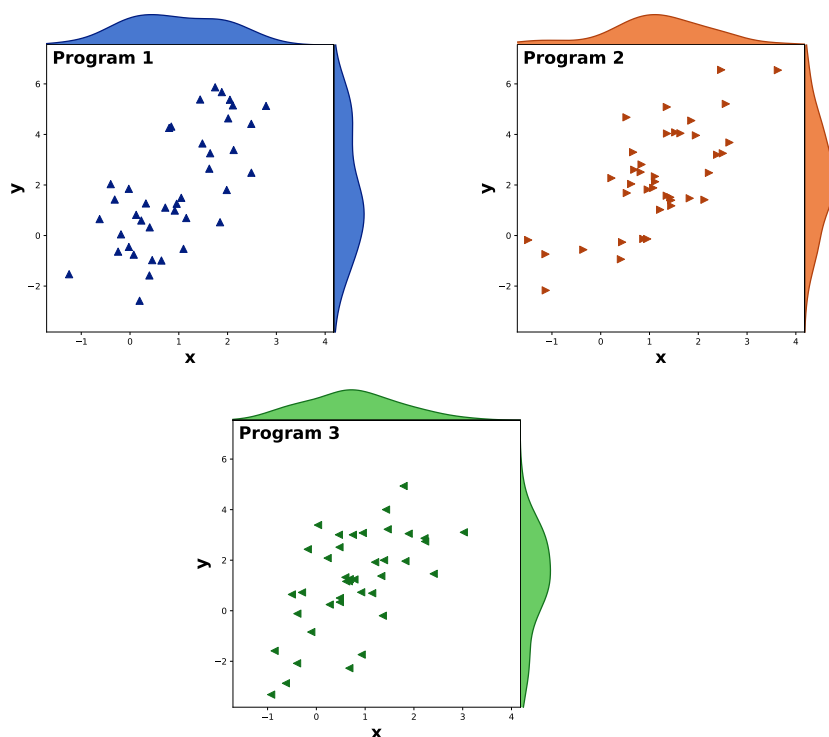
- P. J. Bickel, E. A. Hammel, and J. W. O'Connell. Sex Bias in Graduate Admissions: Data from Berkeley. *Science*, 187, 1975

- S A Julious and M A Mullee. Confounding and Simpson's paradox.
BMJ, 309, 1994

1.3 Intervention on Programs

Initial Analysis

In this exercise we will use the script `intervention_on_programs_1.py`. This script runs an experiment on three programs. Each of the programs produces 2D points like the one shown in the figure below, with density histograms along the axes



We want to know what the causal graph is for each program and how much they differ. The first thing we should consider is whether the distributions above are the same. Do you think they are?

It might be tough to say anything quantitatively about how equal these

distributions are just from inspection. In the script there are two settings `n_samples` and `fit_normal_distribution`. You can use these to get more information from the distribution. Try to figure out how many samples you need before you can see the parameters fitted. Are some of the distributions identical?

Inference of Causal Structures

We want to determine the causal structures of the programs, but first we will make some hypotheses. There are three basic causal structures between X and Y that can explain the data above. See if you can come up with these and sketch them here

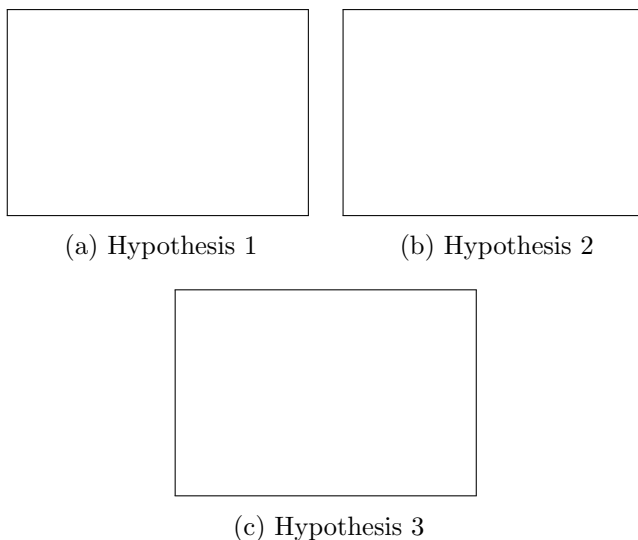


Figure 61

If we intervene on the programs they will act differently depending on their causal structure. Use the table below to check off the variable that you expect to change under intervention of X , or of Y , or under no intervention at all.

| | No interv. | Interv. on X | Interv. on Y |
|--------------|-------------------------|-------------------------|-------------------------|
| Change in: | $\boxed{x} : \boxed{y}$ | $\boxed{x} : \boxed{y}$ | $\boxed{x} : \boxed{y}$ |
| Hypothesis 1 | $\square : \square$ | $\square : \square$ | $\square : \square$ |
| Hypothesis 2 | $\square : \square$ | $\square : \square$ | $\square : \square$ |
| Hypothesis 3 | $\square : \square$ | $\square : \square$ | $\square : \square$ |

In there script there are two settings parameters: x and y . You can use these to set intervene on the programs. Try intervening on each of the variables and note below how the programs behave

| | No inter. | Inter. on X | Inter. on Y |
|------------|-------------------------|-------------------------|-------------------------|
| Change in: | $\boxed{x} : \boxed{y}$ | $\boxed{x} : \boxed{y}$ | $\boxed{x} : \boxed{y}$ |
| Program 1 | $\square : \square$ | $\square : \square$ | $\square : \square$ |
| Program 2 | $\square : \square$ | $\square : \square$ | $\square : \square$ |
| Program 3 | $\square : \square$ | $\square : \square$ | $\square : \square$ |

So which program corresponds to each of the hypotheses?

1.4 A New Switch

In this exercise we will again consider the switch problem. Let's assume that we have another switch-box. This time the lights are blue and purple and we don't know whether blue causes purple, whether purple causes blue or whether there is no causal relationship.

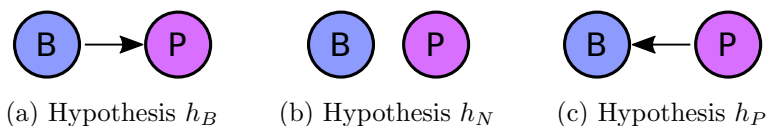


Figure 62

Furthermore we do not know any probabilities before hand. The space of hypotheses and outcome states is

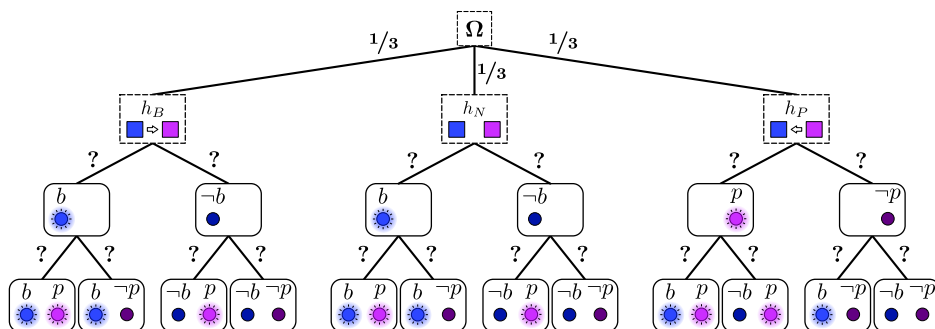


Figure 63: Advanced switch problem.

Using script `a_new_switch.py` we can get samples from the switch. We want to determine all marginal and conditional probabilities for analysing the problem. Use the script to determine these and fill them out below.

| | | | | | | |
|----------|---|---------------|---|---------------|---|--------------------|
| $P(b)$ | : | $P(\neg b)$ | : | $P(p)$ | : | $P(\neg p)$ |
| <hr/> | | | | | | |
| | : | | : | | : | |
| <hr/> | | | | | | |
| $P(p b)$ | : | $P(\neg p b)$ | : | $P(p \neg b)$ | : | $P(\neg p \neg b)$ |
| <hr/> | | | | | | |
| | : | | : | | : | |
| <hr/> | | | | | | |
| $P(b p)$ | : | $P(\neg b p)$ | : | $P(b \neg p)$ | : | $P(\neg b \neg p)$ |
| <hr/> | | | | | | |
| | : | | : | | : | |
| <hr/> | | | | | | |

Can we rule out any of the hypotheses?

We will be analysing the causal structure of the switch by intervention. In order to determine which hypothesis is correct we need to know what to expect from the intervention under each hypothesis. Fill out the table below according to what we expect to happen.

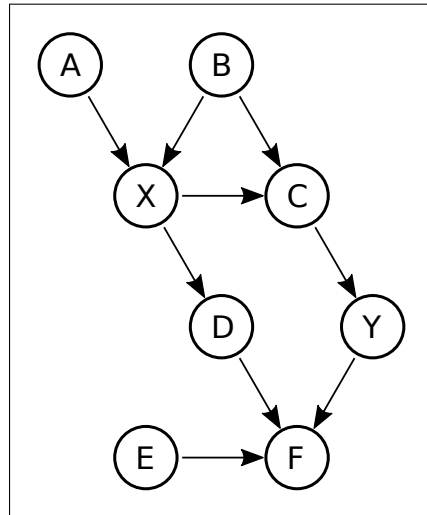
| | Intervention on Blue | | | | Intervention on Purple | | | |
|------------------|----------------------|---------------|---------------|--------------------|------------------------|---------------|---------------|--------------------|
| | $P(d p)$ | $P(\neg p b)$ | $P(p \neg b)$ | $P(\neg p \neg b)$ | $P(b p)$ | $P(\neg b p)$ | $P(b \neg p)$ | $P(\neg b \neg p)$ |
| Hypothesis h_B | | | | | | | | |
| Hypothesis h_N | | | | | | | | |
| Hypothesis h_P | | | | | | | | |

The method `sample_switch()` has two parameters for intervention: `intervene_blue` and `intervene_purple`. Use intervention to determine what hypothesis is correct.

1.5 Information Flow

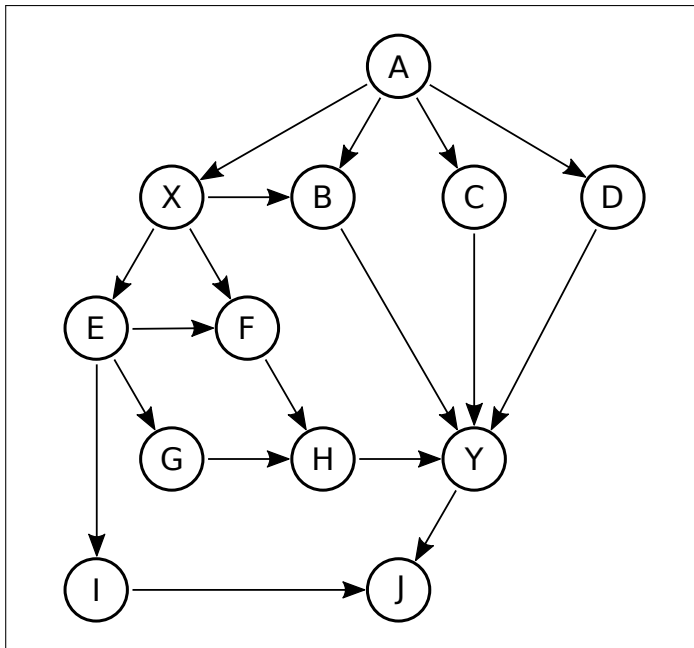
Flow Graph 1

On the right there is a causal graph. We wish to determine the causal relationship between X and Y . Draw all flow-paths between X and Y , and determine two experiments for determining the causal relationships.



Flow Graph 2

In the graph below we again want to determine the causal relationship between X and Y . Draw all information flows between the two nodes of interest.



How would you design experiments to determine the causal relationship?

2nd Exercise

The exercises here concern the topics covered in 7-10.

Some of the exercises spoil solutions to parts of 1st Exercise.

Gaussian Distribution

For some of these exercises we will use Gaussian/normal distributions. The following is a recap of a few properties of Gaussians.

A univariate Gaussian distribution is parameterized by

$$X \sim \mathcal{N}(\mu, \sigma^2), \quad \mu \in \mathbb{R}, \sigma \in \mathbb{R}^+, \quad (1)$$

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}.$$

The sum of two normally distributed, stochastic variables is again a normally distributed

$$X \sim \mathcal{N}(\mu_x, \sigma_x^2), \quad (2)$$

$$Y \sim \mathcal{N}(\mu_y, \sigma_y^2),$$

$$Z = aX + bY,$$

$$Z = aX - bY,$$

$$Z \sim \mathcal{N}(a\mu_x + b\mu_y, \sigma_z^2),$$

$$Z \sim \mathcal{N}(a\mu_x - b\mu_y, \sigma_z^2),$$

$$\sigma_z^2 = a^2\sigma_x^2 + b^2\sigma_y^2 + 2ab\rho\sigma_x\sigma_y, \quad \sigma_z^2 = a^2\sigma_x^2 + b^2\sigma_y^2 - 2ab\rho\sigma_x\sigma_y,$$

where ρ is correlation.

Categorical Distribution

We will also use a categorical distribution parameterized by

$$\begin{aligned} X &\sim \mathcal{C}(\mathbf{a}), \\ p(X = i) &= \mathbf{a}_i. \end{aligned} \tag{3}$$

For example

$$\begin{aligned} X &\sim \mathcal{C}([0.3, 0.4, 0.2, 0.1]) \\ p(X = 0) &= 0.3, \quad p(X = 1) = 0.4, \quad p(X = 2) = 0.2, \quad p(X = 3) = 0.1. \end{aligned} \tag{4}$$

Bernoulli Distribution

The Bernoulli distribution is a special case of the Categorical, with only two categories

$$\begin{aligned} X &\sim \mathcal{B}(0.8) \\ p(X = 0) &= 0.2, \quad p(X = 1) = 0.8. \end{aligned} \tag{5}$$

Exponential Distribution

The exponential distribution has the following form

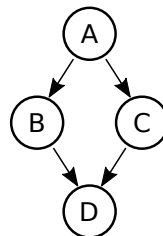
$$X \sim \mathcal{E}(a), \quad a \in \mathbb{R}^+ \tag{6}$$

$$p(x) = a \cdot e^{-a \cdot x}. \tag{7}$$

2.1 Signal Through Variables

We will here analyse the following model

$$\begin{aligned}
 A &\leftarrow N_A, & N_A &\sim \mathcal{N}(2, 2) \\
 B &\leftarrow \frac{1}{2} \cdot A + N_B, & N_B &\sim \mathcal{N}(-1/2, 2) \\
 C &\leftarrow A + N_C, & N_C &\sim \mathcal{N}(1, 1) \\
 D &\leftarrow \frac{2}{3} \cdot B - \frac{1}{3} \cdot C + N_D, & N_D &\sim \mathcal{N}(1, 1).
 \end{aligned}$$



Since all ancestors are normal distributions we know that D itself must also be normally distributed. Let's first consider the means of the four variables. Using (2) determine the means for the four variables

| | | | |
|-----------------|-----------------|-----------------|-----------------|
| $\mathbb{E}[A]$ | $\mathbb{E}[B]$ | $\mathbb{E}[C]$ | $\mathbb{E}[D]$ |
| | | | |

Now that we have those in place, we want to determine how much A controls the other variables. That is; how much of B 's, C 's and D 's variance is determined by A . Using (2) you can determine/compute the variance of A , B and C

| | | |
|-----------------|-----------------|-----------------|
| $\text{var}[A]$ | $\text{var}[B]$ | $\text{var}[C]$ |
| | | |

Determining the variance of D requires computing the correlation between B and C . Alternatively we can derive the final expression of D . Insert the definitions of A , B and C to determine the form of D as a function of N_A , N_B , N_C and N_D

What happened here and why is it interesting?

What is the variance of D ?

2.2 Programming Models

In this exercise we will try implementing a structural equation model of a causal problem in python. We will use the implementation to sample and run experiments on the model.

First of all let's define the model

$$A \leftarrow N_A, \quad N_A \sim \mathcal{E}(7).$$

$$B \leftarrow N_B, \quad N_B \sim \mathcal{N}(1, 4).$$

$$C \leftarrow N_C + 1, \quad N_C \sim \mathcal{C}([1/6, 2/6, 3/6]).$$

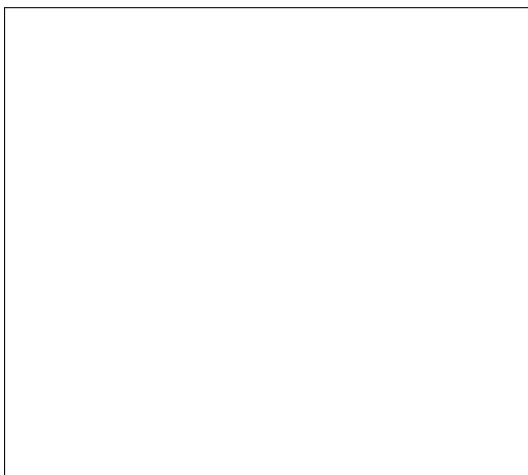
$$D \leftarrow A + B + 30 \cdot N_D, \quad N_D \sim \mathcal{B}(2/5).$$

$$E \leftarrow B + 20 \cdot C.$$

$$F \leftarrow \frac{3}{2} \cdot D + E.$$

We have used four different distributions: exponential distribution \mathcal{E} , normal distribution \mathcal{N} , categorical distribution \mathcal{C} and Bernoulli distribution \mathcal{B} .

As an initial analysis, draw the causal model below



Implementation

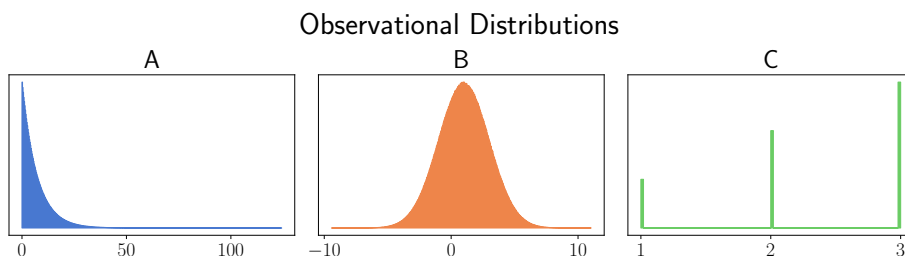
Implement the structural equation model in **python**. You do not need to calculate any probabilities, but simply be able to sample from the model. You can use the **numpy.random** library to find samples for all used distributions.

Observational Distribution

After implementing the model try sampling from the observational distribution. When we ran our programs we got the following means and variances for the first variables

$$\begin{array}{lll} \mathbb{E}[A] \approx 7.0 & \mathbb{E}[B] \approx 1.0 & \mathbb{E}[C] \approx 2.3 \\ \text{var}[A] \approx 49. & \text{var}[B] \approx 4.0 & \text{var}[C] \approx 0.55 \end{array}$$

Its also a good idea to get a visual of the distributions we are working with. Get a big number of samples from the model and plot histograms of the distributions. The first few distributions should look something like. Note that it does not matter if your distributions are more ragged than these - we used a lot of samples.



Interventional Distribution 1

Variable C is an interesting one, which makes some of the other distributions "jump around". Let's investigate what would happen if we controlled this variable. Investigate how the interventional distributions (histograms) of D , E and F look like when we intervene by $C \leftarrow 1$. Did you remove some of the jumping phenomena?

Interventional Distribution 2

This time we wish to reduce the size of E by halving it

$$\overset{\circ}{E} \leftarrow E/2. \quad (8)$$

While the above notation is easy to read a more correct notation may be to specify the probability distribution of the intervention. The probability of $\overset{\circ}{E} = e$ must be equal to the probability of $E = e \cdot 2$. Thus we can also write the intervention as

$$p_E(x \cdot 2) \overset{\circ}{\leftarrow} p_E. \quad (9)$$

You can though, use whichever notation you want.

Implement the above intervention and plot the histograms.

Interventional Distribution 3

Sometimes when working with complicated distributions, it can become relevant to make approximations. The exponential distribution of variable A is very simple, but let us assume that we would rather work with a normal distribution. We can fit a normal distribution to A by matching its mean and variance, which we already know. Make the intervention $p_A \leftarrow \mathcal{N}(\cdot)$, with suitable parameters, and determine the histograms of A , D , E and F . What seems to change?

2.3 Modelling Programs

For this exercise we return to the programs of exercise 1.3. The causal graphs which we inferred in that exercise are seen below. Furthermore the data looks Gaussian.

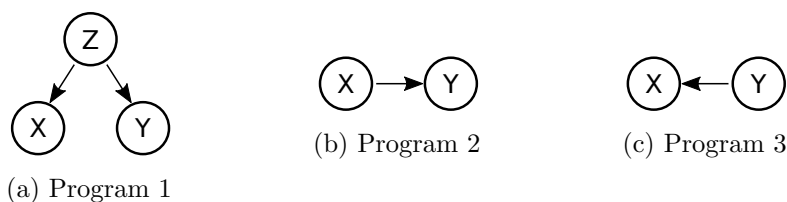
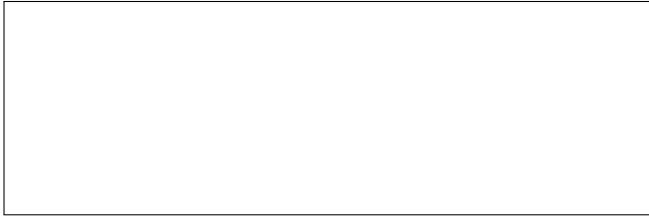


Figure 64

The Programs

We are going to assume that all the programs are normal distributions and sums of normal distributions - which are also normal distributions. First sketch out the program-structures - we've started by writing program 2. We call a_2 the *mixing coefficient* of program 2, as it determines how much X is "mixed into" Y . It does not matter what you name the mixing coefficients etc.



Program 1

$$\begin{array}{ll} X_2 \leftarrow N_{x2}, & N_{x2} \sim \mathcal{N}(\mu_{x2}, \sigma_{x2}) \\ Y_2 \leftarrow a_2 \cdot X_2 + N_{y2}, & N_{y2} \sim \mathcal{N}(\mu_{y2}, \sigma_{y2}) \end{array}$$

Program 2



Program 3

Let's start with the two smaller programs; 2 and 3. Can you come up with a simple strategy for determining the parameters of their normal distributions?

Plan some experiments (2-9) and run the three programs. Measure the means and variances of the data and store in the table below. Write what interventions you do (for example $X \leftarrow 3.14$) - if you do any.

| Experiment | Program 1 | | Program 2 | | Program 3 | | Intervention |
|------------------------------|-----------|-----|-----------|-----|-----------|-----|--------------|
| | X | Y | X | Y | X | Y | |
| 1 Mean Variance | | | | | | | |
| 2 Mean Variance | | | | | | | |
| 3 Mean Variance | | | | | | | |
| 4 Mean Variance | | | | | | | |
| 5 Mean Variance | | | | | | | |

Table 1: Experiments

| Experiment | Program 1 | | Program 2 | | Program 3 | | Intervention |
|------------|-----------|---|-----------|---|-----------|---|--------------|
| | X | Y | X | Y | X | Y | |
| 6 | Mean | | | | | | |
| | Variance | | | | | | |
| 7 | Mean | | | | | | |
| | Variance | | | | | | |
| 8 | Mean | | | | | | |
| | Variance | | | | | | |
| 9 | Mean | | | | | | |
| | Variance | | | | | | |
| 10 | Mean | | | | | | |
| | Variance | | | | | | |

Table 2: More Experiments

Can you determine the noise-distributions for program 2 and 3?
For example N_{x2} and N_{y2} ?

In order to completely determine programs 2 and 3, we need to determine the mixing coefficients like a_2 . Using (2), with the measured means and/or variances, derive the coefficients for programs 2 and 3.

In order to solve problem 1, we need to intervene on the confounder Z . In the programs, it is possible to do intervention on this variable by setting $_z$. Do an intervention experiment on the programs, where you intervene on Z and note the results in the data-table from earlier. Assume the distribution of Z is $Z \sim \mathcal{N}(0, 1)$. Can you determine the noise-distributions of X and Y , as well as the mixing coefficients?

2.4 Counterfactuals 1

Consider the following causal model

$$\begin{array}{ll} X \leftarrow N_x & N_x \sim \mathcal{B}\left(\frac{7}{10}\right) \\ Y \leftarrow N_y & N_y \sim \mathcal{C}\left(\left[\frac{1}{5}, \frac{3}{5}, \frac{1}{5}\right]\right) \\ Z \leftarrow X + Y - N_z & N_z \sim \mathcal{C}\left(\left[\frac{2}{5}, \frac{2}{5}, \frac{1}{5}\right]\right) \end{array}$$

All Known

Say we observe $x = 1, y = 1, z = 0$.

Determine the value of z , given that we counterfactually do $x \leftarrow 0$.

Unknown z

Now we observe $x = 1, y = 1$.

Determine the probability of $z = 1$, given that we counterfactually do $x \leftarrow 0$.

$$p(\dot{z} = 1 \mid x \leftarrow 0, x = 1, y = 1).$$

Unknown y

Now say we observe $x = 1, z = 0$.

We wish to determine the counterfactual distribution

$$p(\dot{z} \mid x \leftarrow 0, z = 0, x = 1).$$

In order to determine this probability, we need to determine the conditional distributions of N_z and N_y . We know what x is, so what is the constraint on the remaining variables influencing Z ?

2.5 Counterfactuals 2

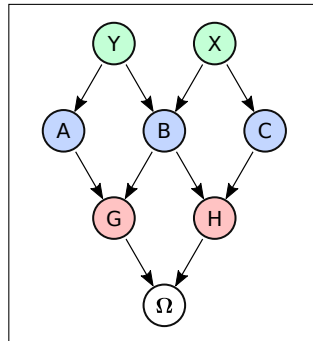
Consider the following causal model

$$N_y, N_x, N_\omega \sim \mathcal{C}\left(\left[\frac{1}{4}, \frac{1}{4}, \frac{2}{4}\right]\right)$$

$$N_a, N_b, N_c \sim \mathcal{B}\left(\frac{1}{3}\right)$$

$$N_g, N_h \sim \mathcal{B}\left(\frac{2}{3}\right)$$

$$\begin{aligned} Y &\leftarrow N_y & X &\leftarrow N_x \\ A &\leftarrow Y \cdot N_a & B &\leftarrow (Y + X) \cdot N_b \\ C &\leftarrow X \cdot (1 + N_c) \\ H &\leftarrow (B + C) \cdot N_h & G &\leftarrow A + B + N_g \\ \Omega &\leftarrow |N_\omega - 1| \cdot G + (N_\omega - 1) \cdot H \end{aligned}$$



Round 1

We want to determine

$$p(\dot{g} \mid y \leftarrow 0, y = 1, a = 1).$$

First write out G as an expression of noise-variables and counterfactual variables (G, Y are the counterfactual variables and N_a, N_b, N_x and N_g are the necessary noise-variables)

Can you determine the value of one of the noise variables already now?

Insert the known noise-variable and the counterfactual value (what we force on Y) to find a final expression for \dot{G}

Now we can make a table of values of N_x , N_b and N_g , together with the values for G and their probabilities

| N_x | N_b | N_g | G | $p(N_x, N_b, N_g)$ |
|-------|-------|-------|-----|--------------------|
| 0 | 0 | 0 | | |
| 1 | 0 | 0 | | |
| 2 | 0 | 0 | | |
| 0 | 1 | 0 | | |
| 1 | 1 | 0 | | |
| 2 | 1 | 0 | | |
| 0 | 0 | 1 | | |
| 1 | 0 | 1 | | |
| 2 | 0 | 1 | | |
| 0 | 1 | 1 | | |
| 1 | 1 | 1 | | |
| 2 | 1 | 1 | | |

What is the distribution of G ?

Round 2

We want to determine

$$p(\dot{g} \mid y \stackrel{\leftarrow}{\leftarrow} 0, c = 0, b > 0, g = 2).$$

Can we determine any noise-variables directly? Can you determine any relationship between noise-variables?

Now let us write G only as a function of noise-variables

Let us also write the counterfactual \dot{G} as a function of noise-variables and counterfactual variables

2.6 Counterfactuals 3

Round 3

Determine

$$p(\dot{\omega} \mid g = 2, h = 2, c < 2, b \leftarrow 0).$$

Project

Project in Causal Inference

3.1 Project Description

In this project we will infer a causal model and compete to see who can do it most efficiently!

Your teacher will have defined a causal model, which is running on a server. You can get samples from the causal model by simply emailing the server. You can ask for as many experiments and samples as you want, and make any combination of interventions, but there is a catch to greed: the winner will be the team that can correctly determine the causal structure using the least experiments and samples!

Details of experiments and competition are below. You need to ask your teacher for the following information:

Server email

Number of nodes
in causal system $N =$

Cost of sample $C_s =$

Cost of experiment $C_e =$

Cost of incorrect guess $C_i =$

3.2 Experiments

For each experiment you want to do, you have to send one email to the server. You write what experiment you want to do in the subject line (the body of the email is irrelevant). The format for communicating with the server is

`n_samples, <interventions>`,

where `n_samples` is the number of samples you want to get for this experiment and `<interventions>` is a comma-separated list of interventions you want to do (without the `<` and `>`).

Here's a couple of examples on the query-format

`10`

Get 10 observational samples.

`10,`

Also get 10 observational samples.

`5, X=1`

Get 5 interventional samples where we force $X \overset{\circ}{\leftarrow} 1$

`25, Z=3.14, X=1, Y=0`

Get 25 interventional samples where we force $Z \overset{\circ}{\leftarrow} 3.14$, $X \overset{\circ}{\leftarrow} 1$ and $Y \overset{\circ}{\leftarrow} 0$

After sending an email to the server it will perform the experiment and send the results back. The received email will have two files with the same data. One file will have the data in a human-readable format as in Table 1a. The other file will be in CSV-format and will have higher precision, but will not be as easy to inspect by people.

| | X | Y | Z | F | G | I |
|---|-----------|-----|-----|----------|-----------|----------|
| 0 | 1.028379 | 1.0 | 0.0 | 0.485250 | 1.276530 | 1.193008 |
| 1 | 2.670319 | 1.0 | 0.0 | 0.810593 | -0.301816 | 1.027937 |
| 2 | 5.294258 | 0.0 | 1.0 | 0.319836 | -3.224374 | 0.254132 |
| 3 | -3.494897 | 0.0 | 1.0 | 0.481506 | 0.640079 | 1.165772 |
| 4 | 2.440380 | 1.0 | 1.0 | 0.816298 | -4.420412 | 2.269149 |

(a) Readable data-format

Table 1: Example of file from experiment server.

3.3 Inferring the Graph

We want you to determine the causal graph of the model your teacher has set up. Your teacher should have informed you about how many nodes there are, but so far that is all we know. To get you started we should consider how many possible graphs you might be working with. The number of causal graphs (assuming no cycles) is the number of *labelled, acyclic, directed graphs* of a certain size. Robinson 1970 has shown that the number of these graphs possible for N nodes is

$$a_N = \sum_{k=1}^N (-1)^{k-1} \binom{N}{k} 2^{k(N-k)} a_{N-k}, \quad (1)$$

where $a_1 = 1$. How many possible graphs can you make with the N nodes provided in this project?¹⁷

Do you think it might be possible for you to try all combinations?

¹⁷If all else fails you can look up the A003024 sequence in OEIS.

If not, then you may have to figure out a more structured approach. A good start is to get some samples from the observational distribution and see what the data looks like. Through correlation you may be able to eliminate some possible causal structures. Then you should start thinking about what interventions would be best for determining the causal graph.

3.4 Making a Guess

You can make a guess of the causal graph by sending

```
guess: <edges>
```

to the server, there `<edges>` is a list of edges. The format of the list of edges is like you would initialize a list of tuples in python, where each tuple has two strings; first the name of the ancestor then the name of the descendant.

The following are examples of attempts at guessing the graph in Figure 28 on page 47

| | |
|--|-----------------|
| <code>guess: [('Y', 'Z'), ('X', 'Z')]</code> | correct guess |
| <code>guess: [('X', 'Z'), ('y', 'z')]</code> | correct guess |
| <code>guess: [('Y', 'Z')]</code> | incorrect guess |
| <code>guess: [('Y', 'Z'), ('Z', 'X')]</code> | incorrect guess |
| <code>guess: [('X', 'Z'), ('Y', 'Z'), ('Z', 'H')]</code> | incorrect guess |

Notice that the order of the edges and the casing of the names do not matter.

Also, every time you make an incorrect guess, the server will remember and decrease your performance.

3.5 Competition

The competition is as follows: we hope for all students to determine the causal graph of the problem. After they have guessed the correct graph, we compute a cost of

$$\begin{aligned} \text{cost} = & C_e \cdot (\# \text{ experiments}) + C_s \cdot (\# \text{ samples}) \\ & + C_i \cdot (\# \text{ incorrect guesses}). \end{aligned} \tag{2}$$

The winner is the team with the lowest cost.

If you make many small experiments, then $C_e \cdot (\# \text{ experiments})$ may get big and give you a bad score. Similarly if you make all your examples huge, then $C_s \cdot (\# \text{ samples})$ may give you a bad score. You therefore need to balance the two to get a good score. Also if you guess correct the first time, then you don't have to be bothered with C_i , but every time you guess wrong it increases the cost!

Finally note that you should prioritise your studies and your understanding of the course higher than your competition score :P.

Project in Causal Inference - Teacher

Read the project description made for the student before this document. Also it does not matter if your students read this document - there will not be any useful spoilers.

4.1 Project Scope

The idea of the project is that the students will attempt to guess a causal graph which you define. This is a very open problem, which can utilize many interesting tools of mathematics, computer science and machine learning, but can also be solved with fairly simple approaches in a less rigorous manner.

The simplest approach is probably to simply make histograms of the data and compare these under interventions. Relatively quickly you can determine the first links in the causal graph and take it from there.

More advanced approaches will include modelling the data (normal distributions etc.). This gives the students precise measures of when a variable has changed, as well as an idea about how many samples are necessary for reliably determining these changes. The methodology can potentially include tools like information theory, active learning and Bayesian optimization to determine the best experiments to perform to most efficiently determine the causal structure.

We recommend making a competition with prizes for the winners, but also

recommend basing any grades on a project report (or something similar), which is independent of the competition scores.

4.2 Server

We have attempted to make the server as easy as possible to use, but note that (depending on version of this document) it has had limited testing and may be improved in the future. In the supplied code base there is a `project`-directory. This directory has the following files at its root for you to use

1. `define_server.py`
2. `server_run.py`
3. `server_reset.py`
4. `server_inspect.py`
5. `plot_causal_graph.py`

The first file `define_server.py` is where you will define the problem which the students will be facing, as well set some settings for the server. The original version of the document is seen in Figures 69 and 70 on pages 170 and 171.

Gmail

Currently the system is set up for using a Gmail-account. You can (probably) use other email-providers, but you might have to rewrite some of the interactions with the server. It is likely that we will experiment with other email-providers as well, so it might be easier in future versions. If you already have a Gmail-account, then you should most definitely NOT use that account. We will be saving the password for the account as free-text on your computer and also remove some security settings.

For setting up Gmail

1. Make a Gmail-account dedicated for the course
2. Allow less-secure apps by
 - (a) Go to <https://myaccount.google.com>
 - (b) *Security*
 - (c) Find *Less secure app access*
 - (d) Turn on access
3. Write email-account name and password in `define_server.py` (more on that later)
4. Run competition
5. Delete account afterwards

Setup

For setting up the server you have to edit a couple of fields of the `ServerSettings`-class in `define_server.py`. Refer to Figures 69 and 70 to see lines and an example of setup (pages 170 and 171).

- Line 9: write the email of the account used for the server.
- Line 10: write the password of the account used for the server (this is why the email should only be used for this experiment).
- Lines 13, 14 and 15 are information about the email-provider. If you use Gmail keep as is.
- Line 18 notes how often the server will check the email (currently every 2 seconds).

Access

Not everyone should have access to the server. This is first of all to avoid the server crashing from spam etc., but also so that you know exactly which student got which data. In lines 20-26 we write down all emails that are allowed to participate in the competition, separated by line-breaks. In the example we have included the server's own email which we use for testing, but that is not necessary - just write the students emails.

When the students ask for data, the server will remember how much data they have received. This will be stored per-email. Thus if the students work in teams it is easiest to have a single student in each team with access to the server. This way you know exactly how much data each team has, because each team is uniquely identified by a single email address.

Also the responses from the server may sometimes end up in spam filters - make sure the students check that.

4.3 Causal Model

It's time to define the causal model for the students to analyse. We have also attempted to make this process as easy as we could. You define the model by creating the `_sample()`-method of the `ExperimentSystem`-class in line 39.

You define a causal variable by assigning some samples to its name as shown in lines 41, 42 and 43. In those lines we use three methods defined in the class for ease of use. All they do is create samples from a normal-distribution, a Bernoulli distribution (binary variable) and a categorical distribution. These samples are assigned to their respective nodes in the caused model - here we use the names `X`, `Y` and `Z`.

You can also use other samplers as long as they have exactly `n_samples` values. We illustrate this in line 47, where we use `numpy`'s beta-distribution sampler for setting variable `F`.

You can combine causal variables to make causal structures, by referring to the names. This is illustrated in line 50, where we make variable `G` a causal descendant of nodes `X` and `F`.

For ensuring that the system works as well, please do *not* make any temporary variables at all. Anything you save temporarily should be saved in a named causal variable. If you want to make hidden variables that the students can not see, name them something that starts with `_`, as shown in line 54. Variable `_H` is hidden and is in this example used as a hidden confounder for variable `I`.

When the students make experiments with the example system they receive data on `X`, `Y`, `Z`, `F`, `G` and `I`, while variable `_H` is hidden from them.

If you want to see the full data including all hidden variables, then you pass a query with a password. For example

```
5, X=1, password=Open_Sesame
```

Get 5 interventional samples
where we force $X \overset{\circ}{\leftarrow} 1$, and
use password `Open_Sesame`.

The password is set in line 31 and should of course be changed in case the students read this document :)

The server will return data-tables with samples, where the causal variables are per default arranged in the order of their definition. This gives the students a hint, because the first variables will be the ancestors and the last variables will be the descendants. In order to remove this hint you can define some other ordering as done in line 60.

We suggest computing a cost for the students as the sum of

$$\begin{aligned} \text{cost} = & C_e \cdot (\# \text{ experiments}) + C_s \cdot (\# \text{ samples}) \\ & + C_i \cdot (\# \text{ incorrect guesses}). \end{aligned}$$

In order to allow the server to compute the cost you should fill these in, in the lines starting at line 33.

You can plot the causal graph with the script `plot_causal_graph.py` to see if it matches your expectations.

Note that the system does not care about casing of the names, so do not name nodes the same name but with different case (for example `x` and `X`).

4.4 Running Server

The server can be run by running the script called `server_run.py`. While the server is running it will display an output similar to that of Figure 66.

In the example the server first checks the emails and concludes that there has been no queries. It then receives two queries from two students, which it successfully understood and answered. It then continues to wait for emails. The server writes out who it has received emails from and what their subject lines were.

If the server receives an email it does not understand, it will write a comment to the student about the format of the emails. One example of the server receiving a bad email is seen in Figure 67, where a student wrote an incorrect query.

```

2019-11-22 13:33:29 -> Checking emails
2019-11-22 13:33:31 -> No new emails.
2019-11-22 13:33:31 -> Sleeping 2s

2019-11-22 13:33:33 -> Checking emails
2019-11-22 13:33:36 -> Email 15: SUCCESS, [15, Y=0], [student_email_1@university.com]
2019-11-22 13:33:39 -> Email 16: SUCCESS, [20, X=1], [student_email_2@university.com]
2019-11-22 13:33:39 -> Sleeping 2s

2019-11-22 13:33:41 -> Checking emails
2019-11-22 13:33:43 -> No new emails.
2019-11-22 13:33:43 -> Sleeping 2s

2019-11-22 13:33:45 -> Checking emails
2019-11-22 13:33:47 -> No new emails.
2019-11-22 13:33:47 -> Sleeping 2s

```

Figure 66: Example of output from server when running.

```

2019-11-22 13:36:51 -> Checking emails
2019-11-22 13:36:53 -> Email 17: Cannot parse subject line,
2019-11-22 13:36:53 -> [X=1] [student_email_1@university.com] <--

```

Figure 67: Example of output from server which received incorrect query.

4.5 Results

The script `server_inspect.py` is used to see how everything is progressing. An example of the output of this program is seen in Figure 68.

First the server prints how many emails it has successfully understood. In the example the students have sent 2 emails with incorrect formatting.

The first table shows some utility information about the students

- how many emails they have sent
- how many samples and experiments they have made
- how many guesses and correct guesses they have made
- how many incorrectly formatted emails they have sent

The second table shows the competition data. It shows the number of experiments, samples and incorrect guesses for each student as well as their final score.

At the bottom we also see the costs of experiments, samples and incorrect guesses for reference.

Each row in the bottom table is only filled out once the student has guessed the correct graph and will remain static afterwards - even if the student keeps querying the system out of interest.

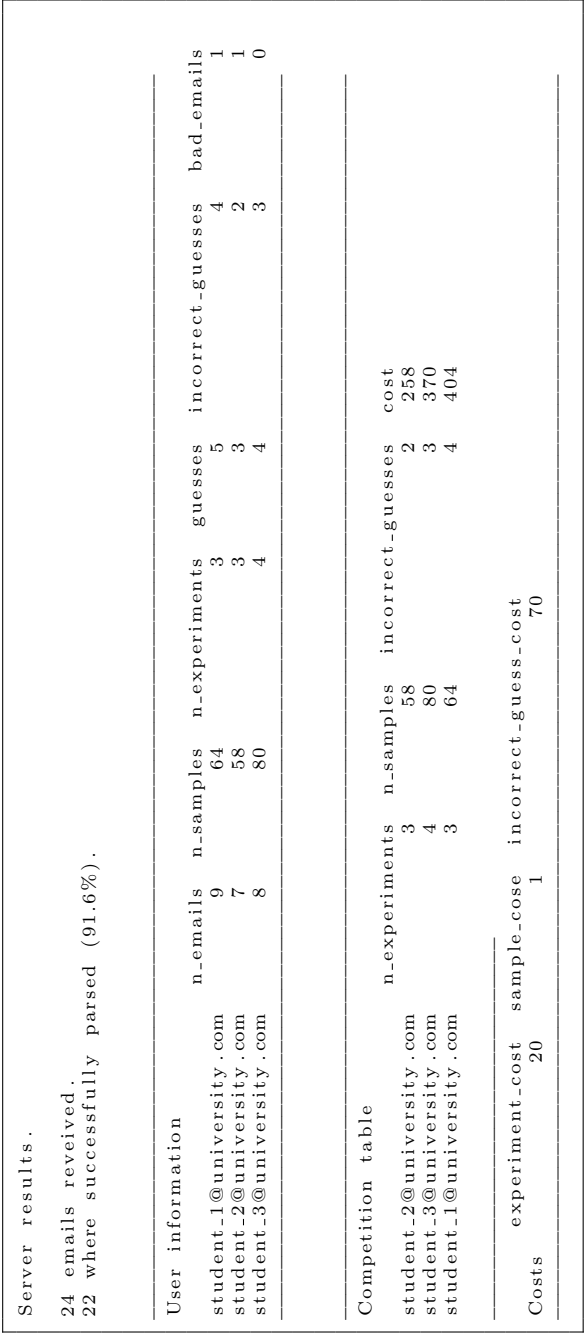


Figure 68: Example of output from server_inspect.py.

```
1 import numpy as np
2
3 from project.src.causal_system import CausalSystem
4
5
6 class ServerSettings:
7
8     # Email information -> you need to turn on "Less secure app access" on Gmail
9     username = "intervention.experiment@gmail.com" #
10    server_password = "hnJk9FYAWN03Er2p" #
11
12    # Email provider information
13    imap_host = "imap.gmail.com" #
14    smtp_host = "smtp.gmail.com" #
15    smtp_port = 587 #
16
17    # Other settings
18    check_email_delay = 2 #
19
20    # Emails with access to experiment #
21    allowed_emails = """
22    intervention.experiment@gmail.com
23    jepno@dtu.dk
24    student_email_1@university.com
25    student_email_2@university.com
26    """ #
```

Continued...

Figure 69: `define_server.py`

Continued...

```
29 # Define the causal system
30 class ExperimentSystem(CausalSystem): #
31     _project_password = "Open_Sesame" #
32
33     # Score settings for competition
34     experiment_cost = 20
35     sample_cost = 1
36     incorrect_guess_cost = 70
37
38     # Causal model
39     def _sample(self, n_samples): #
40         # Using predefined distributions
41         self["X"] = self.normal(mu=2, std=3) #
42         self["Y"] = self.binary(p_success=0.8) #
43         self["Z"] = self.categorical([7, 2, 1]) #
44
45         # You can also use numpy sampling
46         # - but you have to do it correctly by using _n_samples
47         self["F"] = np.random.beta(a=5, b=3, size=n_samples) #
48
49         # Combining is fine
50         self["G"] = self.normal(mu=0, std=1) * self["X"] + self["F"] #
51
52         # This variable is hidden (because of the "_")
53         # It will not be returned in the data unless you have the password
54         self["_H"] = self.normal(mu=0, std=1) #
55
56         # This variable has a hidden confounder
57         self["I"] = self["F"] + self["_H"] + self.normal(mu=0, std=0.2) #
58
59         # Ordering without hint of causal structure
60         self._ordering = ['I', 'Z', 'F', '_H', 'X', 'G', 'Y'] #
```

Figure 70: define_server.py

Exercises - With Solutions

1st Exercise

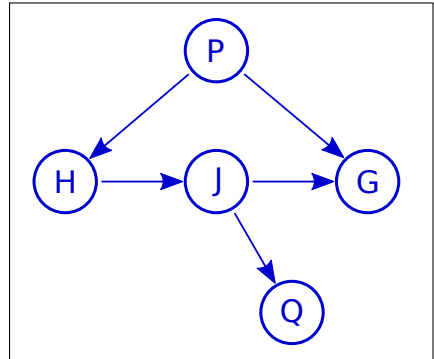
The exercises here concern the topics covered in sections 1-6.

1.1 Problems and Graphs

For each of the descriptions below, draw the related causal graph.

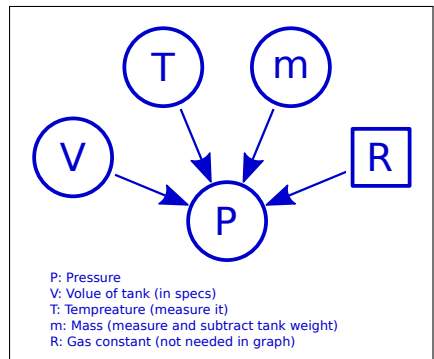
A Small Graph

J is a mediator between ancestor H and descendant G . P is a confounder of H and G , while Q is a child of J .



Ideal Gas Law

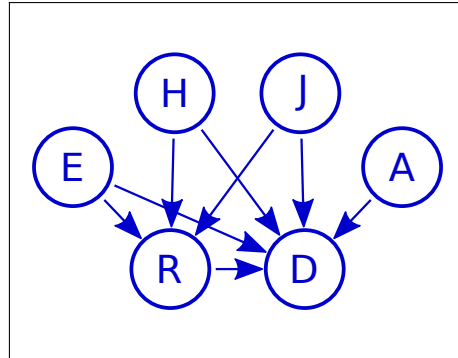
We have a nitrogen tank used for cooling with. We wish to know what its pressure is. We know various specs on the tank and can also do some measurements on it, but we don't know what we need. Do a Wikipedia search on the *Ideal Gas Law* and determine what factors we need to know about the tank, in order to determine the pressure inside the tank. Make a causal diagram illustrating this.



Relaxing is Dangerous?

A study from 2005 found that

... embarking on the Golden Years at age 55 doubled the risk for death before reaching age 65, compared with those who toiled beyond age 60.



That is; people who retired early seemed to die early¹⁸.

Make a causal graph of how you think the following variables interacts: death (D), current age (A), health-condition (H), job type (J), retirement age (R), personal economy (E). Assume everyone in the data already is retired.

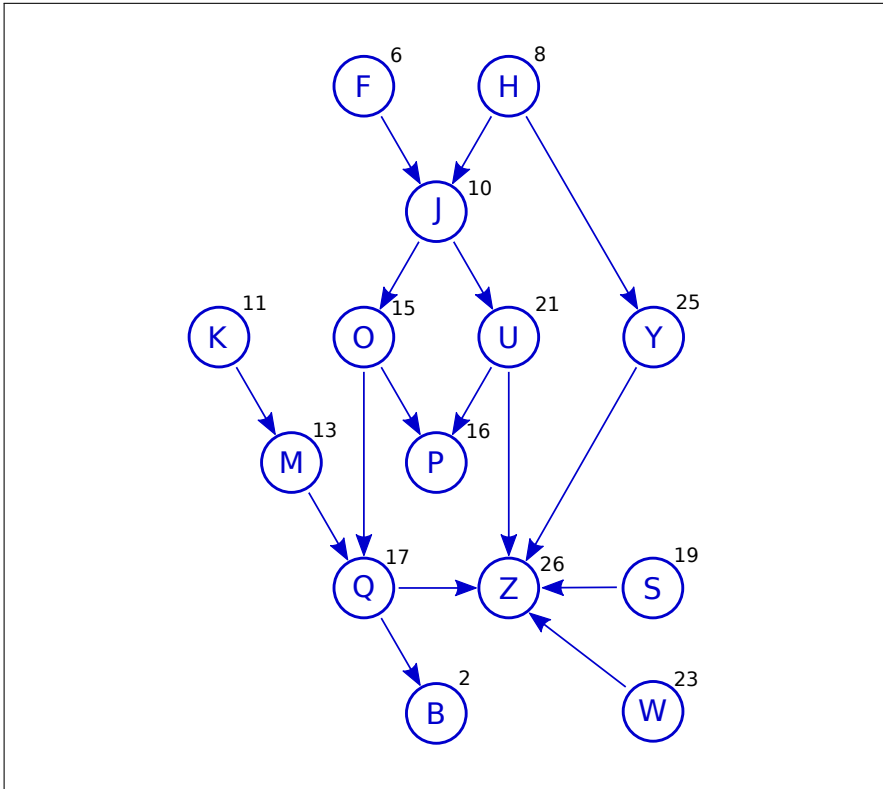
(E), (H) and (J) can all influence retirement age (but your current age post-retirement does not). All factors can influence risk of death. There are a LOT of confounders for this problem, so if the research was not done right, it may not be trustworthy.

¹⁸Luckily, a couple of later studies found the opposite to be true:
<https://www.nytimes.com/2018/01/29/upshot/early-retirement-longevity-health-wellness.html>

A Big Graph

In the following we enumerate the English alphabet (including W) in the normal order, starting with $a = 1$.

F is the parent of J , who is also a child of H . The 5th prime is parent of the 6th prime, who is a parent of the 7th prime, who is a parent of the first prime. All letters whose numbers are odd and greater than 2^4 are parents of the last letter in the alphabet. Letter 5^2 is the child of letter 2^3 . The vowels between letters "j" and "x" are children of J and parents of 2^4 . The only letter whose number is a divisor of 289 is a child of node $(2^4) - 1$.



1.2 A Paradox

Initial Analysis

We are going to analyse some numbers on the fall 1973 admissions to the University of California, Berkeley. We have a hypothesis that the university has a bias against the admission of women! Consider the table below and compute the percentages of admission.

| | Applicants | Admitted | Admitted % |
|-------|------------|----------|------------|
| Men | 2590 | 1269 | 49% |
| Women | 1835 | 556 | 30% |

Is there a bias?

It seems like there is a bias against women.

Extended Analysis

We have some more data on the admission rates for each of 6 departments. We wish to find out which departments are most problematic. Fill out the percentages in table 1 to find any bias within the departments. Let's find all departments where the men/women admissions differs by more that 5%. That is, ignore small differences like 50% and 52%, but consider a difference like 50% and 55% a bias. Note the biases in the tick-boxes on the right of the table (bias against women/no-one/men).

| Depart. | Men | | | | Women | | | | Bias against W/-/M |
|---------|-------|--------|----------|--|-------|--------|----------|--|--------------------------|
| | Appl. | Admit. | Admit. % | | Appl. | Admit. | Admit. % | | |
| A | 825 | 511 | 62% | | 108 | 88 | 82% | | □ □ ☒ |
| B | 560 | 353 | 63% | | 25 | 17 | 68% | | □ □ ☒ |
| C | 325 | 120 | 37% | | 593 | 202 | 34% | | □ ☒ □ |
| D | 417 | 138 | 33% | | 375 | 131 | 35% | | □ ☒ □ |
| E | 191 | 53 | 28% | | 393 | 94 | 24% | | □ ☒ □ |
| F | 272 | 16 | 6% | | 341 | 24 | 7% | | □ ☒ □ |
| Total | 2590 | 1269 | 49% | | 1835 | 556 | 30% | | ☒ □ □ |

Table 1

Is there a bias? Who is the admissions biased against?

If we look at the overall stats it seems like there is a bias against women, but if we look at the departments it seems like there is a bias against men.

Causal Analysis

In order to figure out this problem we will analyse it using the causal tools we have learned. Hypothesis 1 was that there is a causal link from gender G to admission rate A - and no other variables are considered. In Hypothesis 2 we also consider department D . What could the causal links be between G , A and D ? Consider the graph for both Hypothesis 1 and 2 and draw them if Figure 71.

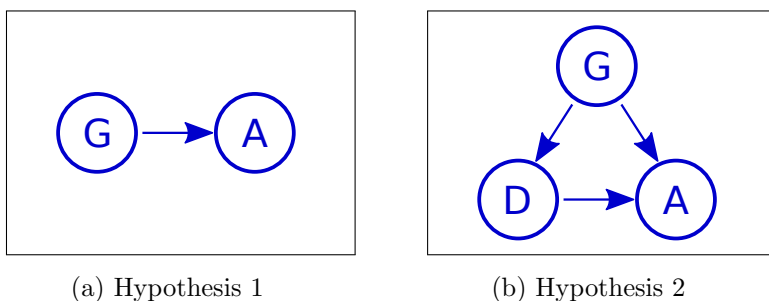


Figure 71

The gender could cause a bias in both hypotheses, but the gender could also cause choice of department which could affect admission rate.

Can you use these graphs to determine whether there is a bias against women?

If hypothesis 1 is correct, then conditioning on the department should not make a difference - but it does. Hypothesis 2 therefore seems more likely. If we want to assess the magnitude of the direct causal link between gender and admission rate under Hypothesis 2, then we have to condition on department (just like we did). When conditioning on department there actually seem to be a small bias against men.



Probabilistic Analysis

The goal of the above exercise is merely to become familiar with Simpson's paradox, know why it occurs and introduce simple tools for handling it. We decided on whether there was a bias or not, based on a qualitative analysis of the conditional admissions rates. It is possible to more formally assess the evidence for and against bias. This requires some fairly advanced Bayesian techniques and is thus only included here for completeness.

Initial analysis

In order to calculate evidence for and against a bias we can look at the problem from a model selection point of view. We start by assuming that we have two equally probable hypotheses; bias H_{bias} or no bias $H_{\neg\text{bias}}$, so that

$$P_0(H_{\neg\text{bias}}) = P_0(H_{\text{bias}}) = \frac{1}{2},$$

where P_0 denotes prior belief.

We will further divide the bias hypothesis into two equally probable halves; bias against women H_{bias_w} , and bias against men H_{bias_m}

$$\begin{aligned} P_0(H_{\text{bias}_w}) + P_0(H_{\text{bias}_m}) &= P_0(H_{\text{bias}}) \\ P_0(H_{\text{bias}_w}) &= P_0(H_{\text{bias}_m}) = \frac{1}{4}. \end{aligned}$$

If the probability of each student is constant, then the number of admissions will be distributed according to a binomial distribution, with some admission rate parameter ρ . We assume that we have a ρ_m for men and a ρ_w for women, and that the prior on the admissions rates is uniform (all admission rates are equally probable).

To fulfil our requirements for the prior we have

$$P_0(H_{\text{bias}_m}) = P_0(\rho_m > \rho_w) = \int_0^1 \int_0^{\rho_m} p_0(\rho_m, \rho_w) \, d\rho_w \, d\rho_m = \frac{1}{4}.$$

Since the prior on admission rates are uniform we have

$$p_0(\rho_m, \rho_w) = \frac{1}{2}.$$

Now given two sets of observations one for men and one for women, we further have

$$P(k_m, n_m | \rho_m) = \binom{n_m}{k_m} \rho_m^{k_m} (1 - \rho_m)^{n_m - k_m}$$

$$P(k_w, n_w | \rho_w) = \binom{n_w}{k_w} \rho_w^{k_w} (1 - \rho_w)^{n_w - k_w}$$

with joint probability:

$$P(k_m, n_m, k_w, n_w | \rho_m, \rho_w) = P(k_m, n_m | \rho_m) P(k_w, n_w | \rho_w)$$

We want to assess the probability of each of the three hypotheses using Bayes rule. For the no-bias hypothesis we have

$$P(H_{\neg \text{bias}} | D, \rho) = \frac{P(D | H_{\neg \text{bias}}, \rho) P_0(H_{\neg \text{bias}} | \rho)}{P(D | \rho)}$$

where we used the D to represent the data $\{k_m, n_m, k_w, n_w\}$, and ρ denotes both ρ_m and ρ_w . The denominator is given by a sum across the three hypotheses. As we do not know ρ we marginalize with respect to this parameter:

$$P(D | H_{\neg \text{bias}}) = \int_0^1 P(D | H_{\neg \text{bias}}, \rho) P_0(H_{\neg \text{bias}} | \rho) d\rho$$

$$P(D | H_{\text{bias}_m}) = \int_0^1 \int_0^{\rho_m} P(D | \rho_m, \rho_w) P_0(H_{\text{bias}_m} | \rho) d\rho_w d\rho_m$$

$$P(D | H_{\text{bias}_w}) = \int_0^1 \int_0^{\rho_w} P(D | \rho_m, \rho_w) P_0(H_{\text{bias}_w} | \rho) d\rho_m d\rho_w$$

Consequently we have

$$P(H_{\neg \text{bias}} | k_m, n_m, k_w, n_w) = \frac{P(D | H_{\neg \text{bias}})}{P(D | H_{\neg \text{bias}}) + P(D | H_{\text{bias}_m}) + P(D | H_{\text{bias}_w})}$$

As the integrals across ρ are complicated it is most straightforward to approximate them with a sum. Using this formulation we find very strong evidence for a bias against women in the aggregated data $P(H_{\text{bias}_m}) > 0.999$.

Conditioned by department

To evaluate if there is evidence for division into departments we formulate a total of six hypotheses, three where data is aggregated across departments (meaning that the ρ will be shared across departments); no bias, bias against women and bias against men. In addition, we then have three hypotheses where we condition on department, that is we allow different ρ_i for each department. Consistent with the previous analysis, we will here set prior probabilities to $\frac{1}{4}$ for the two no bias hypotheses and $\frac{1}{8}$ for the four bias hypotheses. We can then write the likelihood for each of the hypotheses, as follows

$$\begin{aligned}
 P\left(H_{\neg \text{bias}}^{\text{agg}}\right) &= \int_0^1 \prod_i P(D_{m,i}|\rho)P(D_{w,i}|\rho) \, d\rho \\
 P\left(H_{\text{bias}_m}^{\text{agg}}\right) &= \int_0^1 \int_0^{\rho_m} \prod_i P(D_{m,i}|\rho_m)P(D_{w,i}|\rho_w) \, d\rho_w \, d\rho_m \\
 P\left(H_{\text{bias}_w}^{\text{agg}}\right) &= \int_0^1 \int_0^{\rho_w} \prod_i P(D_{m,i}|\rho_m)P(D_{w,i}|\rho_w) \, d\rho_m \, d\rho_w \\
 P\left(H_{\neg \text{bias}}^{\text{dep}}\right) &= \prod_i \int_0^1 P(D_{m,i}|\rho_i)P(D_{w,i}|\rho_i) \, d\rho_i \\
 P\left(H_{\text{bias}_m}^{\text{dep}}\right) &= \prod_i \int_0^1 \int_0^{\rho_{m,i}} P(D_{m,i}|\rho_{m,i})P(D_{w,i}|\rho_{w,i}) \, d\rho_{w,i} \, d\rho_{m,i} \\
 P\left(H_{\text{bias}_w}^{\text{dep}}\right) &= \prod_i \int_0^1 \int_0^{\rho_{w,i}} P(D_{m,i}|\rho_{m,i})P(D_{w,i}|\rho_{w,i}) \, d\rho_{m,i} \, d\rho_{w,i}
 \end{aligned}$$

Here, $D_{m,i}$ and $D_{w,i}$ denote the data recorded for men and women respectively in department i . As an important note there now is an important distinction between modelling bias against men and women separately, doing it this way (separate) the bias would have to be in the same direction across departments, whereas if bias is modelled collectively it could be in different directions across departments. To evaluate the evidence for each of the six hypotheses we simply weight these each of these integrals with the prior probabilities and normalize. This is done in the accompanying python file where you also find this analytically evaluated when collapsing the two bias conditions at the end of the script. We here find

compelling evidence for conditioning on department and having no bias $P(H_{\neg \text{bias}}^{\text{dep}}) \approx 0.99$, importantly here the probabilities for bias actually point more towards a bias against men than women (even though there is really no evidence for either based on this analysis).

Simpson's Paradox

This is a case of *Simpson's Paradox*, which is perhaps the most famous statistical paradox. The paradox happens when a confounder (the department) inverts the analysed relationship (gender causes admission) in the aggregated data (total admissions), as opposed to within the subgroups. The true relationship is the one found in the subgroup, because the confounder no longer has influence.

| | Treatment A | Treatment B |
|--------------|-----------------|-----------------|
| Small stones | 93% (81 / 87) | 87% (234 / 270) |
| Large stones | 73% (192 / 263) | 69% (55 / 80) |
| Total | 78% (273 / 350) | 83% (289 / 350) |

Kidney stone experiment.

Another famous example is from a medical study comparing the success rates of two treatments for kidney stones, whose data is shown on the right. For the total, Treatment B seems best, but this is because Treatment B got more of the easy cases (small stones) and fewer of the difficult cases (large stones). In reality Treatment A is superior.

When we have information and data on the confounder we can eliminate the problem with conditioning. The problem really arises when we can't access the data or perhaps don't even know about the confounder. If we have no information about the size of the kidney stones, then the only possible conclusion is - incorrectly - that Treatment B is best for treating kidney stones. This is where random trials have their strengths. By randomizing who gets which treatment, A or B, we can ensure that no confounder has influenced the decision. This is the point of randomization in A/B testing.

The Berkeley and Kidney Stones examples are from

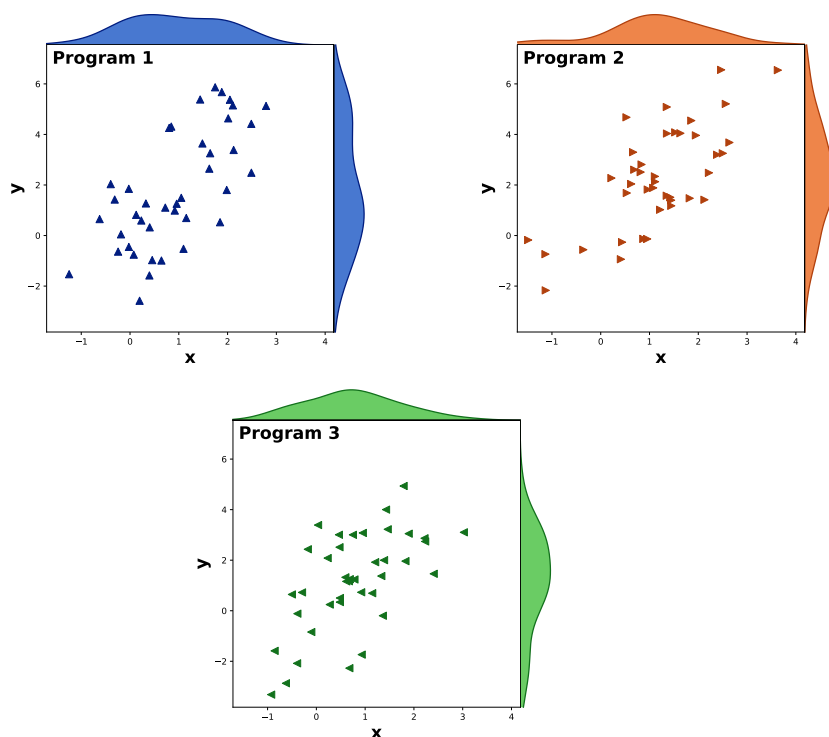
- P. J. Bickel, E. A. Hammel, and J. W. O'Connell. Sex Bias in Graduate Admissions: Data from Berkeley. *Science*, 187, 1975

- S A Julious and M A Mullee. Confounding and Simpson's paradox.
BMJ, 309, 1994

1.3 Intervention on Programs

Initial Analysis

In this exercise we will use the script `intervention_on_programs_1.py`. This script runs an experiment on three programs. Each of the programs produces 2D points like the one shown in the figure below, with density histograms along the axes



We want to know what the causal graph is for each program and how much they differ. The first thing we should consider is whether the distributions above are the same. Do you think they are?

Maybe - difficult to say with this few points

It might be tough to say anything quantitatively about how equal these

distributions are just from inspection. In the script there are two settings `n_samples` and `fit_normal_distribution`. You can use these to get more information from the distribution. Try to figure out how many samples you need before you can see the parameters fitted. Are some of the distributions identical?

I ran about 500 points. The distributions seem quite similar. If I run 50.000 points then they seem really identical.

Inference of Causal Structures

We want to determine the causal structures of the programs, but first we will make some hypotheses. There are three basic causal structures between X and Y that can explain the data above. See if you can come up with these and sketch them here

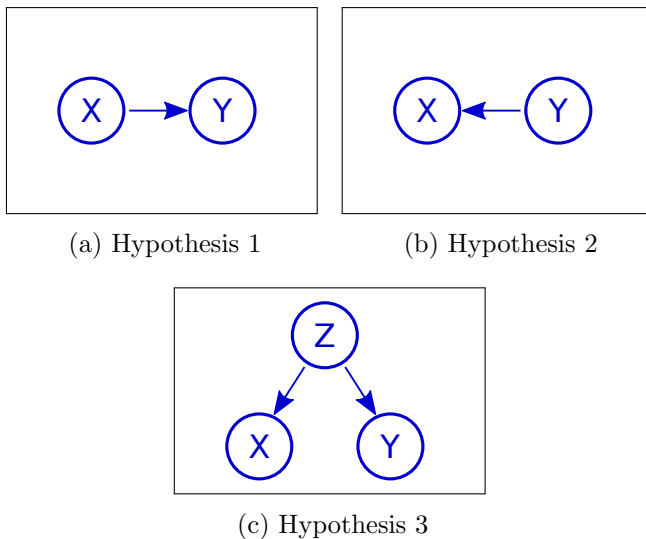


Figure 73

If we intervene on the programs they will act differently depending on their causal structure. Use the table below to check off the variable that you expect to change under intervention of X , or of Y , or under no intervention

at all.

| | No interv. | Interv. on X | Interv. on Y |
|--------------|-------------------------|-------------------------|-------------------------|
| Change in: | $\boxed{X} : \boxed{Y}$ | $\boxed{X} : \boxed{Y}$ | $\boxed{X} : \boxed{Y}$ |
| Hypothesis 1 | $\square : \square$ | $\boxtimes : \boxtimes$ | $\square : \boxtimes$ |
| Hypothesis 2 | $\square : \square$ | $\boxtimes : \square$ | $\boxtimes : \boxtimes$ |
| Hypothesis 3 | $\square : \square$ | $\boxtimes : \square$ | $\square : \boxtimes$ |

The order of the answers of cause change depending on order of the hypotheses.

In there script there are two settings parameters: x and y . You can use these to set intervene on the programs. Try intervening on each of the variables and note below how the programs behave

| | No inter. | Inter. on X | Inter. on Y |
|------------|-------------------------|-------------------------|-------------------------|
| Change in: | $\boxed{X} : \boxed{Y}$ | $\boxed{X} : \boxed{Y}$ | $\boxed{X} : \boxed{Y}$ |
| Program 1 | $\square : \square$ | $\boxtimes : \square$ | $\square : \boxtimes$ |
| Program 2 | $\square : \square$ | $\boxtimes : \boxtimes$ | $\square : \boxtimes$ |
| Program 3 | $\square : \square$ | $\boxtimes : \square$ | $\boxtimes : \boxtimes$ |

So which program corresponds to each of the hypotheses?

Using our ordering we find that Program 1 is Hypothesis 3, Program 2 is Hypothesis 1 and Program 3 is Hypothesis 2.

1.4 A New Switch

In this exercise we will again consider the switch problem. Let's assume that we have another switch-box. This time the lights are blue and purple and we don't know whether blue causes purple, whether purple causes blue or whether there is no causal relationship.

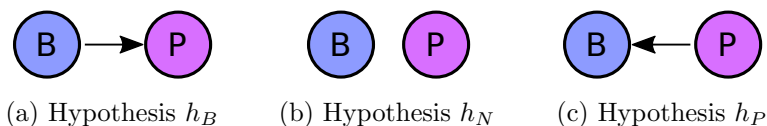


Figure 74

Furthermore we do not know any probabilities before hand. The space of hypotheses and outcome states is

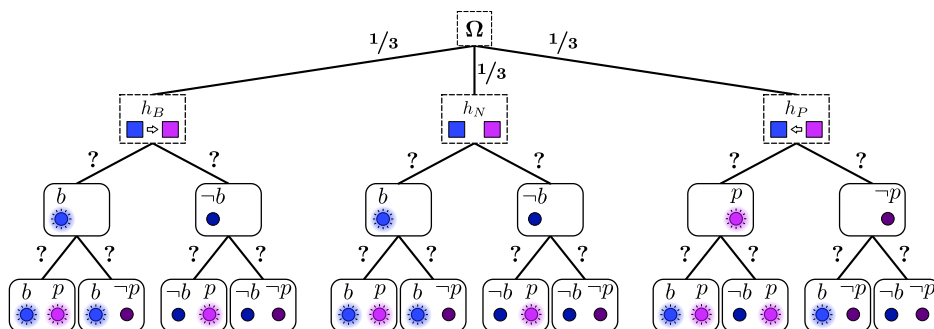


Figure 75: Advanced switch problem.

Using script `a_new_switch.py` we can get samples from the switch. We want to determine all marginal and conditional probabilities for analysing the problem. Use the script to determine these and fill them out below.

| | | | | | | |
|--------|---|-------------|---|--------|---|-------------|
| $P(b)$ | : | $P(\neg b)$ | : | $P(p)$ | : | $P(\neg p)$ |
| 0.50 | : | 0.50 | : | 0.50 | : | 0.50 |

| | | | | | | |
|------------|---|-----------------|---|-----------------|---|----------------------|
| $P(p b)$ | : | $P(\neg p b)$ | : | $P(p \neg b)$ | : | $P(\neg p \neg b)$ |
| 0.65 | : | 0.35 | : | 0.35 | : | 0.65 |

| | | | | | | |
|------------|---|-----------------|---|-----------------|---|----------------------|
| $P(b p)$ | : | $P(\neg b p)$ | : | $P(b \neg p)$ | : | $P(\neg b \neg p)$ |
| 0.65 | : | 0.35 | : | 0.35 | : | 0.65 |

Can we rule out any of the hypotheses?

We can rule out h_N , because the variables are clearly not independent:

$$P(p | b) \neq P(p), \quad P(b | p) \neq P(b), \quad P(b, p) \neq P(b) P(p).$$

We will be analysing the causal structure of the switch by intervention. In order to determine which hypothesis is correct we need to know what to expect from the intervention under each hypothesis. Fill out the table below according to what we expect to happen.

| | Intervention on Blue | | | | Intervention on Purple | | | |
|------------------|----------------------|---------------|---------------|--------------------|------------------------|---------------|---------------|--------------------|
| | $P(d b)$ | $P(\neg d b)$ | $P(d \neg b)$ | $P(\neg d \neg b)$ | $P(b p)$ | $P(\neg b p)$ | $P(b \neg p)$ | $P(\neg b \neg p)$ |
| Hypothesis h_B | 0.65 | 0.35 | 0.35 | 0.65 | 0.50 | 0.50 | 0.50 | 0.50 |
| Hypothesis h_N | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| Hypothesis h_P | 0.50 | 0.50 | 0.50 | 0.50 | 0.35 | 0.65 | 0.65 | 0.35 |

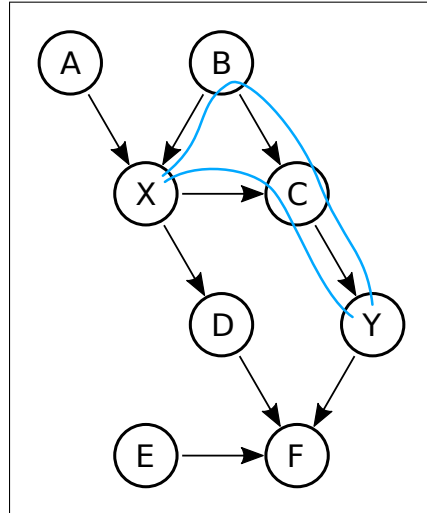
The method `sample_switch()` has two parameters for intervention: `intervene_blue` and `intervene_purple`. Use intervention to determine what hypothesis is correct.
Hypothesis h_P is correct.

1.5 Information Flow

Flow Graph 1

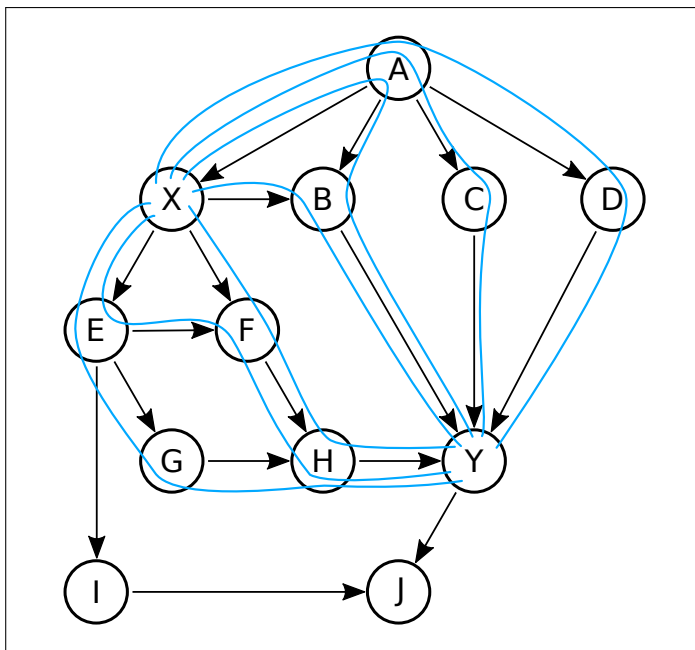
On the right there is a causal graph. We wish to determine the causal relationship between X and Y . Draw all flow-paths between X and Y , and determine two experiments for determining the causal relationships.

We want to estimate the causal path $X \rightarrow C \rightarrow Y$ while ignoring the flow of information of the path $X \leftarrow B \rightarrow C \rightarrow Y$. In order to do this we can make an observational experiment and condition on B , or we could make an interventional experiment where we intervene on X .



Flow Graph 2

In the graph below we again want to determine the causal relationship between X and Y . Draw all information flows between the two nodes of interest.



How would you design experiments to determine the causal relationship?

We can again intervene on X if possible.

We can also condition on B , C and D , but the better solution is to only condition on A .

2nd Exercise

The exercises here concern the topics covered in 7-10.
Some of the exercises spoil solutions to parts of 1st Exercise.

Gaussian Distribution

For some of these exercises we will use Gaussian/normal distributions. The following is a recap of a few properties of Gaussians.

A univariate Gaussian distribution is parameterized by

$$\begin{aligned} X &\sim \mathcal{N}(\mu, \sigma^2), & \mu &\in \mathbb{R}, \sigma \in \mathbb{R}^+, \\ p(x) &= \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}. \end{aligned} \tag{1}$$

The sum of two normally distributed, stochastic variables is again a normally distributed

$$\begin{aligned} X &\sim \mathcal{N}(\mu_x, \sigma_x^2), \\ Y &\sim \mathcal{N}(\mu_y, \sigma_y^2), \end{aligned} \tag{2}$$

$$\begin{aligned} Z &= aX + bY, & Z &= aX - bY, \\ Z &\sim \mathcal{N}(a\mu_x + b\mu_y, \sigma_z^2), & Z &\sim \mathcal{N}(a\mu_x - b\mu_y, \sigma_z^2), \\ \sigma_z^2 &= a^2\sigma_x^2 + b^2\sigma_y^2 + 2ab\rho\sigma_x\sigma_y, & \sigma_z^2 &= a^2\sigma_x^2 + b^2\sigma_y^2 - 2ab\rho\sigma_x\sigma_y, \end{aligned}$$

where ρ is correlation.

Categorical Distribution

We will also use a categorical distribution parameterized by

$$\begin{aligned} X &\sim \mathcal{C}(\mathbf{a}), \\ p(X = i) &= \mathbf{a}_i. \end{aligned} \tag{3}$$

For example

$$\begin{aligned} X &\sim \mathcal{C}([0.3, 0.4, 0.2, 0.1]) \\ p(X = 0) &= 0.3, \quad p(X = 1) = 0.4, \quad p(X = 2) = 0.2, \quad p(X = 3) = 0.1. \end{aligned} \tag{4}$$

Bernoulli Distribution

The Bernoulli distribution is a special case of the Categorical, with only two categories

$$\begin{aligned} X &\sim \mathcal{B}(0.8) \\ p(X = 0) &= 0.2, \quad p(X = 1) = 0.8. \end{aligned} \tag{5}$$

Exponential Distribution

The exponential distribution has the following form

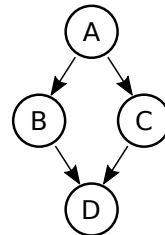
$$X \sim \mathcal{E}(a), \quad a \in \mathbb{R}^+ \tag{6}$$

$$p(x) = a \cdot e^{-a \cdot x}. \tag{7}$$

2.1 Signal Through Variables

We will here analyse the following model

$$\begin{aligned}
 A &\leftarrow N_A, & N_A &\sim \mathcal{N}(2, 2) \\
 B &\leftarrow \frac{1}{2} \cdot A + N_B, & N_B &\sim \mathcal{N}\left(-1/2, 2\right) \\
 C &\leftarrow A + N_C, & N_C &\sim \mathcal{N}(1, 1) \\
 D &\leftarrow \frac{2}{3} \cdot B - \frac{1}{3} \cdot C + N_D, & N_D &\sim \mathcal{N}(1, 1).
 \end{aligned}$$



Since all ancestors are normal distributions we know that D itself must also be normally distributed. Let's first consider the means of the four variables. Using (2) determine the means for the four variables

| $\mathbb{E}[A]$ | $\mathbb{E}[B]$ | $\mathbb{E}[C]$ | $\mathbb{E}[D]$ |
|-----------------|---|-----------------|---|
| 2 | $\frac{1}{2} \cdot 2 - \frac{1}{2} = \frac{1}{2}$ | $2 + 1 = 3$ | $\frac{2}{3} \cdot \frac{1}{2} - \frac{1}{3} \cdot 3 + 1 = \frac{1}{3}$ |

Now that we have those in place, we want to determine how much A controls the other variables. That is; how much of B 's, C 's and D 's variance is determined by A . Using (2) you can determine/compute the variance of A , B and C

| $\text{var}[A]$ | $\text{var}[B]$ | $\text{var}[C]$ |
|-----------------|--|-----------------|
| 2 | $\left(\frac{1}{2}\right)^2 \cdot 2 + 2 = \frac{5}{2}$ | $2 + 1 = 3$ |

Determining the variance of D requires computing the correlation between B and C . Alternatively we can derive the final expression of D . Insert the definitions of A , B and C to determine the form of D as a function of N_A , N_B , N_C and N_D

$$\begin{aligned}
 D &\leftarrow \frac{2}{3} \cdot B - \frac{1}{3} \cdot C + N_D \\
 &= \frac{2}{3} \cdot \underbrace{\left(\frac{1}{2} \cdot A + N_B\right)}_B - \frac{1}{3} \cdot \underbrace{(A + N_C)}_C + N_D \\
 &= \frac{1}{3} \cdot A + \frac{2}{3} \cdot N_B - \frac{1}{3} \cdot A - \frac{1}{3} \cdot N_C + N_D \\
 &= \frac{2}{3} \cdot N_B - \frac{1}{3} \cdot N_C + N_D.
 \end{aligned}$$

What happened here and why is it interesting?

The signal of A through B and C exactly cancels out, making D completely independent of A !

From the graph we would expect D to be caused by A , but it turns out it is not.

While it may rarely work out this precisely in practice, there are situations where signals can cancel each other out, so that expected causes has almost no influence of effects.

What is the variance of D ?

$$\text{var}[D] = \left(\frac{2}{3}\right)^2 \cdot 2 + \left(-\frac{1}{3}\right)^2 \cdot 1 + 1 = \frac{8}{9} + \frac{1}{9} + 1 = 2.$$

2.2 Programming Models

In this exercise we will try implementing a structural equation model of a causal problem in python. We will use the implementation to sample and run experiments on the model.

First of all let's define the model

$$A \leftarrow N_A, \quad N_A \sim \mathcal{E}(7).$$

$$B \leftarrow N_B, \quad N_B \sim \mathcal{N}(1, 4).$$

$$C \leftarrow N_C + 1, \quad N_C \sim \mathcal{C}([1/6, 2/6, 3/6]).$$

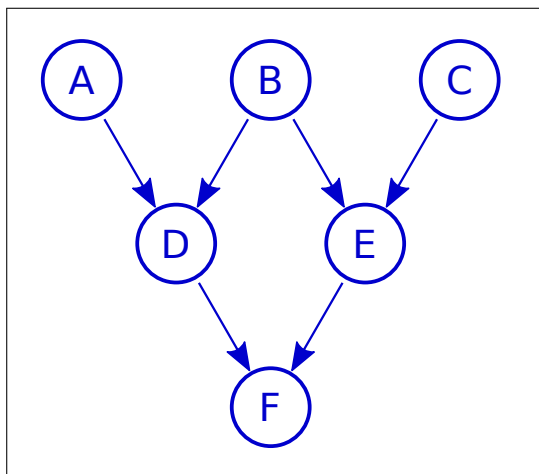
$$D \leftarrow A + B + 30 \cdot N_D, \quad N_D \sim \mathcal{B}(2/5).$$

$$E \leftarrow B + 20 \cdot C.$$

$$F \leftarrow \frac{3}{2} \cdot D + E.$$

We have used four different distributions: exponential distribution \mathcal{E} , normal distribution \mathcal{N} , categorical distribution \mathcal{C} and Bernoulli distribution \mathcal{B} .

As an initial analysis, draw the causal model below



Implementation

Implement the structural equation model in `python`. You do not need to calculate any probabilities, but simply be able to sample from the model. You can use the `numpy.random` library to find samples for all used distributions.

Observational Distribution

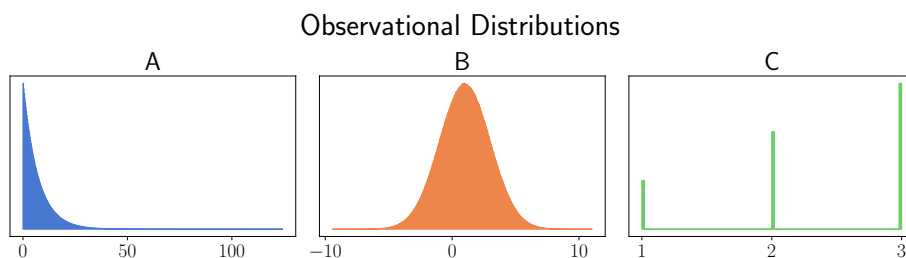
After implementing the model try sampling from the observational distribution. When we ran our programs we got the following means and variances for the first variables

$$\begin{array}{lll} \mathbb{E}[A] \approx 7.0 & \mathbb{E}[B] \approx 1.0 & \mathbb{E}[C] \approx 2.3 \\ \text{var}[A] \approx 49. & \text{var}[B] \approx 4.0 & \text{var}[C] \approx 0.55 \end{array}$$

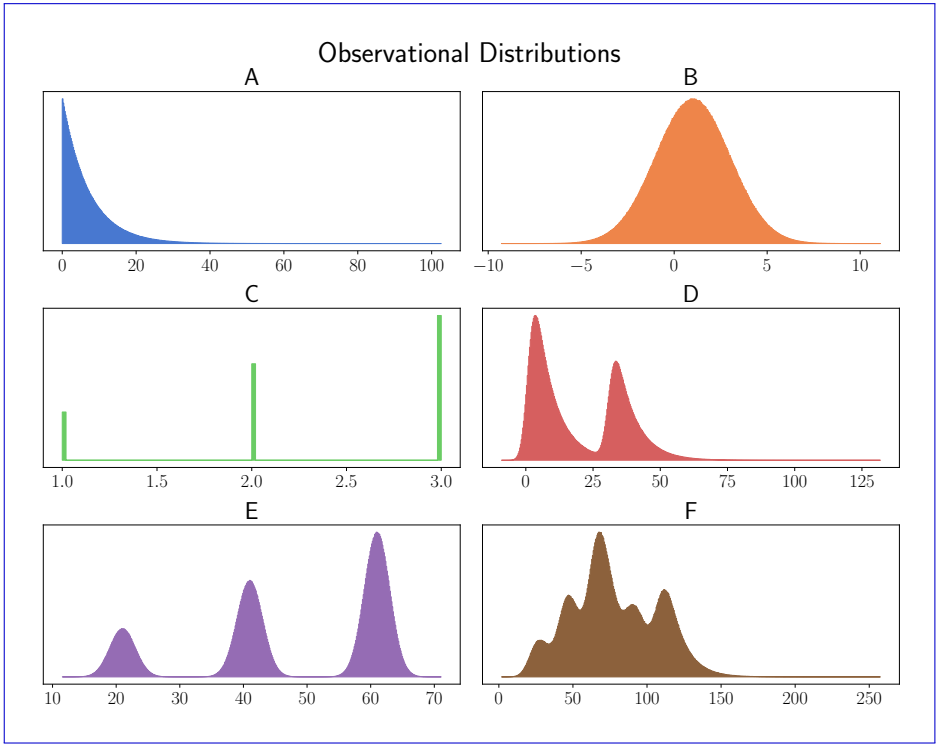
For the remaining variables we got

$$\begin{array}{lll} \mathbb{E}[D] \approx 20. & \mathbb{E}[E] \approx 48. & \mathbb{E}[F] \approx 78. \\ \text{var}[D] \approx 270. & \text{var}[E] \approx 230. & \text{var}[F] \approx 840. \end{array}$$

Its also a good idea to get a visual of the distributions we are working with. Get a big number of samples from the model and plot histograms of the distributions. The first few distributions should look something like. Note that it does not matter if your distributions are more ragged than these - we used a lot of samples.



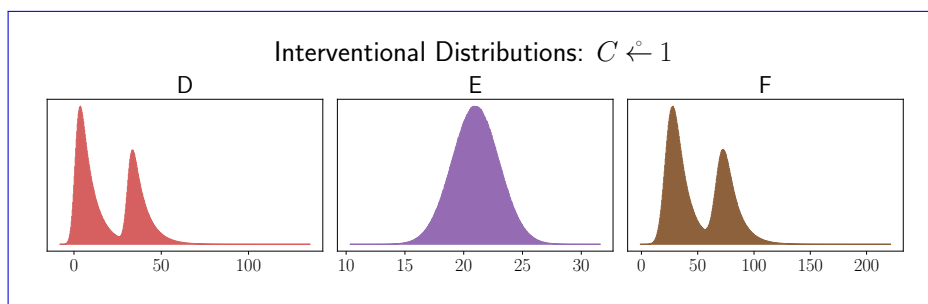
The full observational distributions looks like below.



Interventional Distribution 1

Variable C is an interesting one, which makes some of the other distributions "jump around". Let's investigate what would happen if we controlled this variable. Investigate how the interventional distributions (histograms) of D , E and F look like when we intervene by $C \leftarrow 1$. Did you remove some of the jumping phenomena?

We find the interventional distributions to be:



We have removed the "jumps" from variable E , but still have the jumps in D , because of its own noise-variable N_D .

Interventional Distribution 2

This time we wish to reduce the size of E by halving it

$$\overset{\circ}{E} \leftarrow E/2. \quad (8)$$

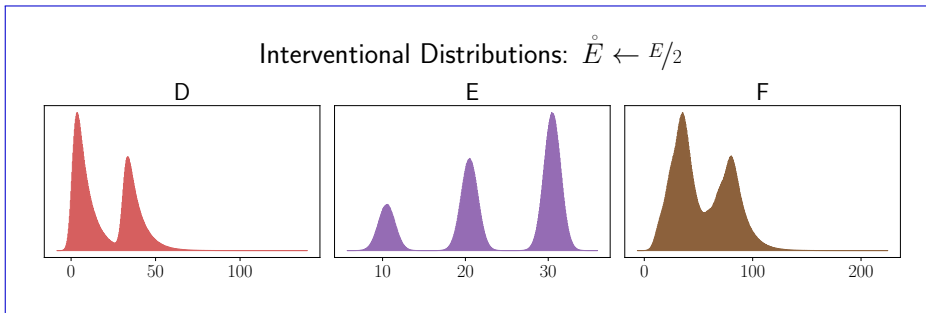
While the above notation is easy to read a more correct notation may be to specify the probability distribution of the intervention. The probability of $\overset{\circ}{E} = e$ must be equal to the probability of $E = e \cdot 2$. Thus we can also write the intervention as

$$p_E(x \cdot 2) \overset{\circ}{\leftarrow} p_E. \quad (9)$$

You can though, use whichever notation you want.

Implement the above intervention and plot the histograms.

We find the interventional distributions to be:

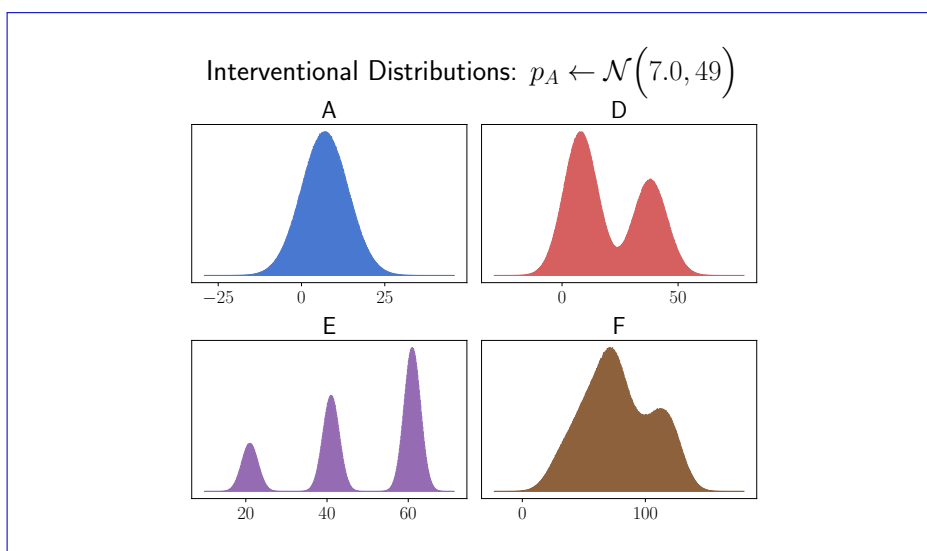


Notice that E look completely identical, but the x-axis has changed.

Interventional Distribution 3

Sometimes when working with complicated distributions, it can become relevant to make approximations. The exponential distribution of variable A is very simple, but let us assume that we would rather work with a normal distribution. We can fit a normal distribution to A by matching its mean and variance, which we already know. Make the intervention $p_A \leftarrow \mathcal{N}(\cdot)$, with suitable parameters, and determine the histograms of A , D , E and F . What seems to change?

We find the interventional distributions to be:



The D -variable now has two normal distributions instead of two exponential distributions. This seems to smooth out F , which becomes more of a blob.

2.3 Modelling Programs

For this exercise we return to the programs of exercise 1.3. The causal graphs which we inferred in that exercise are seen below. Furthermore the data looks Gaussian.

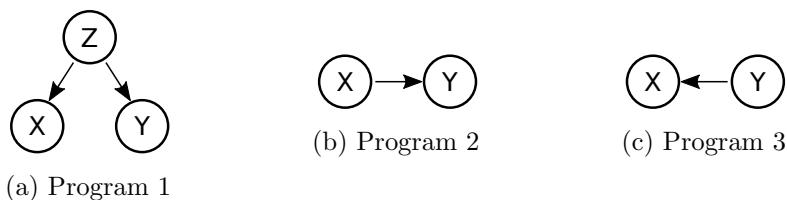


Figure 76

The Programs

We are going to assume that all the programs are normal distributions and sums of normal distributions - which are also normal distributions. First sketch out the program-structures - we've started by writing program 2. We call a_2 the *mixing coefficient* of program 2, as it determines how much X is "mixed into" Y . It does not matter what you name the mixing coefficients etc.

$$\begin{array}{ll}
Z_1 \leftarrow N_{z1}, & N_{z1} \sim \mathcal{N}(\mu_{z1}, \sigma_{z1}) \\
X_1 \leftarrow a_{1x} \cdot Z_1 + N_{x1}, & N_{x1} \sim \mathcal{N}(\mu_{x1}, \sigma_{x1}) \\
Y_1 \leftarrow a_{1y} \cdot Z_1 + N_{y1}, & N_{y1} \sim \mathcal{N}(\mu_{y1}, \sigma_{y1})
\end{array}$$

Program 1

$$\begin{array}{ll}
X_2 \leftarrow N_{x2}, & N_{x2} \sim \mathcal{N}(\mu_{x2}, \sigma_{x2}) \\
Y_2 \leftarrow a_2 \cdot X_2 + N_{y2}, & N_{y2} \sim \mathcal{N}(\mu_{y2}, \sigma_{y2})
\end{array}$$

Program 2

$$\begin{array}{ll}
Y_3 \leftarrow N_{y3}, & N_{y3} \sim \mathcal{N}(\mu_{y3}, \sigma_{y3}) \\
X_3 \leftarrow a_3 \cdot Y_3 + N_{x3}, & N_{x3} \sim \mathcal{N}(\mu_{x3}, \sigma_{x3})
\end{array}$$

Program 3

Let's start with the two smaller programs; 2 and 3. Can you come up with a simple strategy for determining the parameters of their normal distributions?

For program 2 we can measure X 's distribution directly, and if we intervene and set $X \leftarrow 0$ we can measure Y 's distribution directly. For program 3 it's the other way around.

Plan some experiments (2-9) and run the three programs. Measure the means and variances of the data and store in the table below. Write what interventions you do (for example $X \leftarrow 3.14$) - if you do any.

| Experiment | Program 1 | | Program 2 | | Program 3 | | Intervention |
|------------|-----------|-----|-----------|-----|-----------|---------|--------------|
| | X | Y | X | Y | X | Y | |
| 1 | Mean | 1.0 | 2.0 | 1.0 | 2.0 | 1.0 | 2.0 |
| | Variance | 1.0 | 4.0 | 1.0 | 4.0 | 1.0 | 4.0 |
| 2 | Mean | 0.0 | 2.0 | 0.0 | [0.6] | 0.0 | 2.0 |
| | Variance | 0.0 | 4.0 | 0.0 | [2.0] | 0.0 | 4.0 |
| 3 | Mean | 1.0 | 0.0 | 1.0 | 0.0 | [0.3] | 0.0 |
| | Variance | 1.0 | 0.0 | 1.0 | 0.0 | [0.5] | 0.0 |
| 4 | Mean | | | | | | |
| | Variance | | | | | | |
| 5 | Mean | | | | | | |
| | Variance | | | | | | |

Table 1: Experiments

| | | Program 1 | | Program 2 | | Program 3 | |
|------------|----------|-----------|---------|-----------|-----|-----------|-----|
| Experiment | | X | Y | X | Y | X | Y |
| 6 | Mean | | | | | | |
| | Variance | | | | | | |
| 7 | Mean | | | | | | |
| | Variance | | | | | | |
| 8 | Mean | | | | | | |
| | Variance | | | | | | |
| 9 | Mean | [0.4] | [0.8] | 1.0 | 2.0 | 1.0 | 2.0 |
| | Variance | [0.3] | [1.2] | 1.0 | 4.0 | 1.0 | 4.0 |
| 10 | Mean | [1.6] | [3.2] | 1.0 | 2.0 | 1.0 | 2.0 |
| | Variance | [0.3] | [1.2] | 1.0 | 4.0 | 1.0 | 4.0 |

Table 2: More Experiments

The numbers of the table are just suggestions. The last rows are for determining program 1, which needs additional interventions. We have marked the places where an intervention on a variable has changed *another* variable - which indicates a causal relationship.

Can you determine the noise-distributions for program 2 and 3?

For example N_{x2} and N_{y2} ?

Yes we can! :D

From observational data we can determine the distributions of N_{x2} and N_{y3} .

From the interventional data where $X \overset{\circ}{\leftarrow} 0$ we can determine N_{y2} .

From the data where $Y \overset{\circ}{\leftarrow} 0$ we can determine N_{x3} .

$$\begin{aligned} N_{x2} &\sim \mathcal{N}(1, 1) & N_{y2} &\sim \mathcal{N}(0.6, 2) \\ N_{x3} &\sim \mathcal{N}(0.3, 0.5) & N_{y3} &\sim \mathcal{N}(2, 4) \end{aligned}$$

In order to completely determine programs 2 and 3, we need to determine the mixing coefficients like a_2 . Using (2), with the measured means and/or variances, derive the coefficients for programs 2 and 3.

We know what the observational means/variances are, and what the means/variances of the noise-distributions are. From this we can determine the coefficients of the programs by

$$\begin{aligned} \mathbb{E}[Y_2] &= a_2 \cdot \mu_{x2} + \mu_{y2} = a_2 \cdot 1 + 0.6 = 2 \quad \Leftrightarrow \quad a_2 = \frac{2 - 0.6}{1} = 1.4, \\ \text{var}[Y_2] &= a_2^2 \cdot \sigma_{x2}^2 + \sigma_{y2}^2 = a_2^2 \cdot 1 + 2 = 4 \quad \Leftrightarrow \quad a_2 = \sqrt{\frac{4 - 2}{1}} \approx 1.41, \\ \mathbb{E}[X_3] &= a_3 \cdot \mu_{y3} + \mu_{x3} = a_3 \cdot 2 + 0.3 = 1 \quad \Leftrightarrow \quad a_3 = \frac{1 - 0.3}{2} = 0.35, \\ \text{var}[X_3] &= a_3^2 \cdot \sigma_{y3}^2 + \sigma_{x3}^2 = a_3^2 \cdot 4 + 0.5 = 1 \quad \Leftrightarrow \quad a_3 = \sqrt{\frac{1 - 0.5}{4}} \approx 0.35. \end{aligned}$$

Thus we can use either the means or the variances to determine the mixing coefficients - or both to verify our results.

In order to solve problem 1, we need to intervene on the confounder Z . In the programs, it is possible to do intervention on this variable by setting $_z$. Do an intervention experiment on the programs, where you intervene on Z and note the results in the data-table from earlier. Assume the distribution

of Z is $Z \sim \mathcal{N}(0, 1)$. Can you determine the noise-distributions of X and Y , as well as the mixing coefficients?

Yes we can.

If we make an experiment where $Z \overset{\circ}{\leftarrow} 0$, then we can read off the noise-distributions directly as

$$\begin{aligned} N_{x1} &\sim \mathcal{N}(0.4, 0.3) \\ N_{y1} &\sim \mathcal{N}(0.8, 1.2). \end{aligned}$$

From the interventional distribution where $Z \overset{\circ}{\leftarrow} 2$ we have

$$\begin{aligned} \mathbb{E}[X_1] &= a_{x1} \cdot \overset{\circ}{\mu}_{z1} + \mu_{x1} = a_{x1} \cdot 2.0 + 0.4 = 1.6 &\Leftrightarrow& a_{x1} = \frac{1.6 - 0.4}{2.0} = 0.6, \\ \text{var}[X_1] &= a_{x1}^2 \cdot \overset{\circ}{\sigma}_{z1}^2 + \sigma_{x1}^2 = a_{x1}^2 \cdot 0.0 + 0.3 = 0.3 &\Leftrightarrow& 0.3 = 0.3, \end{aligned}$$

$$\begin{aligned} \mathbb{E}[Y_1] &= a_{y1} \cdot \overset{\circ}{\mu}_{z1} + \mu_{y1} = a_{y1} \cdot 2.0 + 0.8 = 3.2 &\Leftrightarrow& a_{y1} = \frac{3.2 - 0.8}{2.0} = 1.2, \\ \text{var}[Y_1] &= a_{y1}^2 \cdot \overset{\circ}{\sigma}_{z1}^2 + \sigma_{y1}^2 = a_{y1}^2 \cdot 0.0 + 1.2 = 1.2 &\Leftrightarrow& 1.2 = 1.2 \end{aligned}$$

Knowing the mixing coefficients we can now use the observational means and variances to find

$$\begin{aligned} \mathbb{E}[X_1] &= a_{x1} \cdot \mu_{z1} + \mu_{x1} = 0.6 \cdot \mu_{z1} + 0.4 = 1.0 &\Leftrightarrow& \mu_{z1} = \frac{1.0 - 0.4}{0.6} = 1.0, \\ \text{var}[X_1] &= a_{x1}^2 \cdot \sigma_{z1}^2 + \sigma_{x1}^2 = 0.6^2 \cdot \sigma_{z1}^2 + 0.3 = 1.0 &\Leftrightarrow& \sigma_{z1}^2 = \frac{1.0 - 0.3}{0.6^2} = 1.95, \end{aligned}$$

$$\begin{aligned} \mathbb{E}[Y_1] &= a_{y1} \cdot \mu_{z1} + \mu_{y1} = 1.2 \cdot \mu_{z1} + 0.8 = 2.0 &\Leftrightarrow& \mu_{z1} = \frac{2.0 - 0.8}{1.2} = 1.0, \\ \text{var}[Y_1] &= a_{y1}^2 \cdot \sigma_{z1}^2 + \sigma_{y1}^2 = 1.2^2 \cdot \sigma_{z1}^2 + 1.2 = 4.0 &\Leftrightarrow& \sigma_{z1}^2 = \frac{4.0 - 1.2}{1.2^2} = 1.95 \end{aligned}$$

The true values used in the program is $\mu_{z1} = 1.0$ and $\sigma_{z1}^2 = 2.0$, which can be found if we use higher precision.

2.4 Counterfactuals 1

Consider the following causal model

$$\begin{array}{ll} X \leftarrow N_x & N_x \sim \mathcal{B}\left(\frac{7}{10}\right) \\ Y \leftarrow N_y & N_y \sim \mathcal{C}\left(\left[\frac{1}{5}, \frac{3}{5}, \frac{1}{5}\right]\right) \\ Z \leftarrow X + Y - N_z & N_z \sim \mathcal{C}\left(\left[\frac{2}{5}, \frac{2}{5}, \frac{1}{5}\right]\right) \end{array}$$

All Known

Say we observe $x = 1, y = 1, z = 0$.

Determine the value of z , given that we counterfactually do $x \leftarrow 0$.

First we determine all the noise-variables we can.

From the values of x and y we have

$$X = N_x = 1, \quad Y = N_y = 1.$$

From this we find that

$$Z = X + Y - N_z = N_x + N_y - N_z = 1 + 1 - N_z = 0 \quad \Leftrightarrow \quad N_z = 2.$$

So the counterfactual z will be

$$z = \dot{N}_x + N_y - N_z = 0 + 1 - 2 = -1.$$

This is hardly a counterfactual question; we know all information.

Unknown z

Now we observe $x = 1, y = 1$.

Determine the probability of $z = 1$, given that we counterfactually do $x \leftarrow 0$.

$$p(\dot{z} = 1 \mid x \leftarrow 0, x = 1, y = 1).$$

We know all variables determining z , except for one; the noise-variable N_z . We can marginalize over this variable in order to determine the probability of $z = 1$

$$\begin{aligned}
 p(\dot{z} = 1 \mid x \leftarrow 0, x = 1, y = 1) \\
 &= \sum_{i=0}^2 p(z = 1 \mid x = 0, y = 1, N_z = i) \cdot p(N_z = i) \\
 &= p(z = 1 \mid x = 0, y = 1, N_z = 0) \cdot p(N_z = 0) = 1 \cdot \frac{2}{5} = \frac{2}{5}.
 \end{aligned}$$

This also, is not really a counterfactual question; we are simply computing the forward probabilities of z .

Unknown y

Now say we observe $x = 1, z = 0$.

We wish to determine the counterfactual distribution

$$p(\dot{z} \mid x \leftarrow 0, z = 0, x = 1).$$

In order to determine this probability, we need to determine the conditional distributions of N_z and N_y . We know what x is, so what is the constraint on the remaining variables influencing Z ?

We observed

$$\begin{aligned}
 z &= x + Y - N_z \\
 0 &= 1 + Y - N_z \\
 Y - N_z &= -1.
 \end{aligned}$$

Since Y and N_z are not affected by the counterfactual intervention on x they remain constant. Thus by changing x we simply have

$$\dot{z} = \dot{x} + Y - N_z = 0 - 1 = -1. \quad (10)$$

Therefore the counterfactual distribution of \dot{z} is simply a constant of -1 .

2.5 Counterfactuals 2

Consider the following causal model

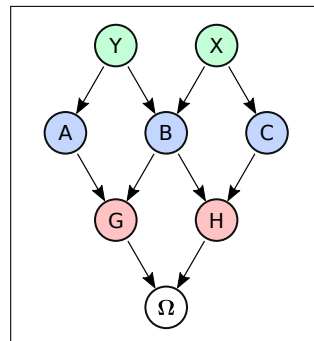
$$N_y, N_x, N_\omega \sim \mathcal{C}\left(\left[\frac{1}{4}, \frac{1}{4}, \frac{2}{4}\right]\right)$$

$$N_a, N_b, N_c \sim \mathcal{B}\left(\frac{1}{3}\right)$$

$$N_g, N_h \sim \mathcal{B}\left(\frac{2}{3}\right)$$

$$\begin{aligned} Y &\leftarrow N_y & X &\leftarrow N_x \\ A &\leftarrow Y \cdot N_a & B &\leftarrow (Y + X) \cdot N_b \\ C &\leftarrow X \cdot (1 + N_c) \\ H &\leftarrow (B + C) \cdot N_h & G &\leftarrow A + B + N_g \end{aligned}$$

$$\Omega \leftarrow |N_\omega - 1| \cdot G + (N_\omega - 1) \cdot H$$



Round 1

We want to determine

$$p(\dot{g} \mid y \leftarrow 0, y = 1, a = 1).$$

First write out G as an expression of noise-variables and counterfactual variables (G, Y are the counterfactual variables and N_a, N_b, N_x and N_g are the necessary noise-variables)

$$\begin{aligned} G \leftarrow A + B + N_g &= Y \cdot N_a + (Y + X) \cdot N_b + N_g \\ &= Y \cdot N_a + (Y + N_x) \cdot N_b + N_g. \end{aligned}$$

Can you determine the value of one of the noise variables already now?

Yes! We can determine N_a . We know that $a = 1$ and this can only happen if $Y = 1$ (which it was) and $N_a = 1$.

Insert the known noise-variable and the counterfactual value (what we force on Y) to find a final expression for \dot{G}

For the counterfactual case we have

$$\begin{aligned} \dot{G} \leftarrow \dot{y} \cdot N_a + (\dot{y} + N_x) \cdot N_b + N_g \\ = 0 \cdot 1 + (0 + N_x) \cdot N_b + N_g = N_x \cdot N_b + N_g. \end{aligned}$$

Now we can make a table of values of N_x , N_b and N_g , together with the values for G and their probabilities

| N_x | N_b | N_g | G | $p(N_x, N_b, N_g)$ | $= p(N_x) \cdot p(N_b) \cdot p(N_g)$ |
|-------|-------|-------|-----|---|--------------------------------------|
| 0 | 0 | 0 | 0 | $1/4 \cdot 2/3 \cdot 1/3 = 2/36 = 1/18$ | |
| 1 | 0 | 0 | 0 | $1/4 \cdot 2/3 \cdot 1/3 = 2/36 = 1/18$ | |
| 2 | 0 | 0 | 0 | $2/4 \cdot 2/3 \cdot 1/3 = 4/36 = 1/9$ | |
| 0 | 1 | 0 | 0 | $1/4 \cdot 1/3 \cdot 1/3 = 1/36 = 1/36$ | |
| 1 | 1 | 0 | 1 | $1/4 \cdot 1/3 \cdot 1/3 = 1/36 = 1/36$ | |
| 2 | 1 | 0 | 2 | $2/4 \cdot 1/3 \cdot 1/3 = 2/36 = 1/18$ | |
| 0 | 0 | 1 | 1 | $1/4 \cdot 2/3 \cdot 2/3 = 4/36 = 1/9$ | |
| 1 | 0 | 1 | 1 | $1/4 \cdot 2/3 \cdot 2/3 = 4/36 = 1/9$ | |
| 2 | 0 | 1 | 1 | $2/4 \cdot 2/3 \cdot 2/3 = 8/36 = 2/9$ | |
| 0 | 1 | 1 | 1 | $1/4 \cdot 1/3 \cdot 2/3 = 2/36 = 1/18$ | |
| 1 | 1 | 1 | 2 | $1/4 \cdot 1/3 \cdot 2/3 = 2/36 = 1/18$ | |
| 2 | 1 | 1 | 3 | $2/4 \cdot 1/3 \cdot 2/3 = 4/36 = 1/9$ | |
| | | | | $36/36 = 1$ | |

What is the distribution of G ?

$$G \sim \mathcal{C} \left(\left[9/36, 19/36, 4/36, 4/36 \right] \right) = \mathcal{C} \left(\left[1/4, 19/36, 1/9, 1/9 \right] \right).$$

This is a categorical distribution which is often 1, sometimes 0, and rarely 2 or 3.

Round 2

We want to determine

$$p(\dot{g} \mid y \dot{\leftarrow} 0, c = 0, b > 0, g = 2).$$

Can we determine any noise-variables directly? Can you determine any relationship between noise-variables?

Yes, first of all we can conclude that $x = 0$ since

$$C \leftarrow X \cdot (1 + N_c) \quad \text{and} \quad c = 0.$$

Furthermore since $b > 0$, then $N_b = 1$ and $N_y > 0$ because

$$B \leftarrow (Y + X) \cdot N_b.$$

Now let us write G only as a function of noise-variables

$$\begin{aligned} G &\leftarrow A + B + N_g \\ &= Y \cdot N_a + (Y + X) \cdot N_b + N_g \\ &= N_y \cdot N_a + N_y + x + N_g \\ &= N_y \cdot (N_a + 1) + N_g. \end{aligned}$$

Let us also write the counterfactual \dot{G} as a function of noise-variables and counterfactual variables

$$\begin{aligned} \dot{G} &\leftarrow A + B + N_g \\ &= \dot{y} \cdot N_a + (\dot{y} + X) \cdot N_b + N_g \\ &= \dot{y} \cdot N_a + \dot{y} + N_x + N_g \\ &= \dot{y} \cdot (N_a + 1) + N_g. \end{aligned}$$

2.6 Counterfactuals 3

Round 3

Determine

$$p(\dot{\omega} \mid g = 2, h = 2, c < 2, b \leftarrow 0).$$

First of all, if $h = 2$ then $N_h = 1$ and $B + C = 2$.

We can therefore assert C if we know the value of B .

At the level of the visible causal variables the equations becomes

$$G \leftarrow A + B + N_g \quad (11)$$

$$H \leftarrow B + C. \quad (12)$$

The combinations which produce $h = 2$, $g = 2$ and $c < 2$ are

| B | C | A | N_g |
|-----|-----|-----|-------|
| 2 | 0 | 0 | 0 |
| 1 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 |

If we further expand the equations for G and H we get

$$G \leftarrow \underbrace{N_y \cdot N_a}_A + \underbrace{(N_y + N_x) \cdot N_b}_B + N_g \quad (13)$$

$$H \leftarrow \underbrace{(N_y + N_x) \cdot N_b}_B + \underbrace{N_x \cdot (1 + N_c)}_C \quad (14)$$

Depending on the value of B we can determine combinations of the noise-variables

| $b = 2$ | | | $b = 1$ | | |
|---------|-------|-------|---------|-------|-------|
| N_b | N_x | N_y | N_b | N_x | N_y |
| 1 | 2 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 | 0 | 1 |
| 1 | 0 | 2 | | | |

Let's consider the first case where $b = 2$. In this case N_x must be 0, because otherwise h would be more than 2. This leaves $N_y = 2$, which further implies that $N_a = 0$, because otherwise g would be more than 2. Thus the only combination remaining is where $N_a = 0, N_x = 0, N_y = 2$ and N_c can be anything.

The other case is where $b = 1$. For $h = 2$ we must have that $c = 1$, which can only happen when $N_x = 1$ and $N_c = 0$. This only allows the case where $N_b = 1, N_x = 1$ and $N_y = 0$. The set of possible combinations are:

| A | B | C | N_a | N_c | N_g | N_x | N_y | $P(\cdot)$ | $P(\cdot \cdot)$ |
|-----|-----|-----|-------|-------|-------|-------|-------|---|------------------------------|
| 0 | 2 | 0 | 0 | 0 | 0 | 0 | 2 | $\frac{2}{3} \cdot \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{1}{4} \cdot \frac{2}{4} = \frac{8}{432}$ | $\frac{8}{24} = \frac{1}{3}$ |
| 0 | 2 | 0 | 0 | 1 | 0 | 0 | 2 | $\frac{2}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{4} \cdot \frac{2}{4} = \frac{4}{432}$ | $\frac{4}{24} = \frac{1}{6}$ |
| 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | $\frac{1}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{1}{4} \cdot \frac{1}{4} = \frac{4}{432}$ | $\frac{4}{24} = \frac{1}{6}$ |
| 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | $\frac{2}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{1}{4} \cdot \frac{1}{4} = \frac{8}{432}$ | $\frac{8}{24} = \frac{1}{3}$ |

Where we have computed the probability $P(\cdot)$ of each combination, as well as the conditional probability (on all known information) $P(\cdot | \cdot)$. If we enforce $B \dot{\leftarrow} 0$, then the final probabilities of G and H becomes

| A | \dot{B} | C | N_g | N_h | G | H | $P(\cdot \cdot)$ |
|-----|-----------|-----|-------|-------|-----|-----|---|
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | $\frac{1}{6} + \frac{1}{3} = \frac{1}{2}$ |
| 0 | 0 | 1 | 1 | 1 | 1 | 1 | $\frac{1}{6} + \frac{1}{3} = \frac{1}{2}$ |

So there is $\frac{1}{2}$ chance of $G = H = 0$ and $\frac{1}{2}$ chance of $G = H = 1$.

We can now answer the counterfactual question. The counterfactual outcomes and distribution of Ω are

| G | H | N_ω | Ω | $P(\cdot \cdot)$ |
|-----|-----|---------------|----------|---|
| 0 | 0 | $\{0, 1, 2\}$ | 0 | $\frac{1}{2} \cdot 1 = \frac{1}{2}$ |
| 1 | 1 | $\{0, 1\}$ | 0 | $\frac{1}{2} \cdot \frac{2}{4} = \frac{1}{4}$ |
| 1 | 1 | 2 | 2 | $\frac{1}{2} \cdot \frac{2}{4} = \frac{1}{4}$ |

$$\Omega = 2 \cdot K, \quad K \sim \mathcal{B}(1/4)$$

