

0001
0002
0003
0004
0005
0006
0007
0008
0009
0010
0011
0012
0013
0014
0015
0016
0017
0018
0019
0020
0021
0022
0023
0024
0025
0026
0027
0028
0029
0030
0031
0032
0033
0034
0035
0036
0037
0038
0039
0040
0041
0042
0043
0044
0045
0046
0047
0048
0049
0050
0051
0052
0053
0054
0055
0056
0057
0058
0059
0060
0061
0062
0063
0064

GT Causalité

Séance 2

This session is based on Brady Neal's lecture notes (Chapter 3, and part of Chapter 4). It also uses [KF09, Chapter 3], [PJS17, Chapters 2 and 6], and [Pea22, Chapters 1 and 2]. Sometimes I refer to [Lau96, Chapter 3] for some proofs.

1 Bayesian networks

Here I note that the explanations in Brady Neal's lecture notes are quite cryptic, so I am mostly following the story in [KF09]. Something that does not help in Neal's notes is that he amalgamates causality concepts with bayesian networks terminology which are unrelated with causality. Indeed, in this section we shall not think at all about causality and only think about independence relations.

To make things a bit more simple, here we consider random variables $(X_1, \dots, X_d) \in \mathcal{X}_1 \times \dots \times \mathcal{X}_d$ with distribution P and we do assume that P admits a density p with respect to the product measure $\mu_1 \times \dots \times \mu_d$.

We shall use the following vocabulary about DAGs: $\text{Pa}(X)$ denote the set of *parents* of X , *ie.* nodes Y such that there is a directed edges $Y \rightarrow X$; $\text{NonDesc}(X)$ the set of non descendants of X , *ie.* nodes in $V \setminus X$ that cannot be reached from X through a directed path.

Definition 1. *Let G be a DAG on vertices $V = \{X_1, \dots, X_d\}$. Then P factorizes according to G if there are conditional densities $p_i : \mathcal{X}_i \times \prod_{j \in \text{Pa}(X_i)} \mathcal{X}_j \rightarrow \mathbb{R}_+$ such that*

$$p(x_1, \dots, x_n) = \prod_{i=1}^d p_i(x_i \mid (x_j)_{j \in \text{Pa}(X_i)}) \quad P - \text{as.} \quad (1)$$

Definition 2. *A bayesian network is a pair (G, P) where G is a DAG and P factorizes according to G .*

Example 1. *Consider the DAG $X_1 \rightarrow X_2 \rightarrow X_3$. The distribution P with density $p(x_1, x_2, x_3) = p_1(x_1)p_{2|1}(x_2 \mid x_1)p_{3|2}(x_3 \mid x_2)$ factorizes according to the DAG.*

The graph G can be viewed in two very different ways:

1. as a data structure that provides the skeleton for representing a joint distribution compactly in a factorized way;

0065
0066
0067
0068
0069
0070
0071
0072
0073
0074
0075
0076
0077
0078
0079
0080
0081
0082
0083
0084
0085
0086
0087
0088
0089
0090
0091
0092
0093
0094
0095
0096
0097
0098
0099
0100
0101
0102
0103
0104
0105
0106
0107
0108
0109
0110
0111
0112
0113
0114
0115
0116
0117
0118
0119
0120
0121
0122
0123
0124
0125
0126
0127
0128

2. as a compact representation for a set of conditional independence assumptions about a distribution.

As we will see, these two views are, in a strong sense, equivalent.

Given three subsets of variables $\mathbf{X}, \mathbf{Y}, \mathbf{Z} \subset V$ we say that $(\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z})$ holds in P if \mathbf{X} and \mathbf{Y} are independent given \mathbf{Z} ; if \mathbf{X} and \mathbf{Y} are independent unconditionally, we say $(\mathbf{Y} \perp \mathbf{Z})$ or $(\mathbf{X} \perp \mathbf{Y} \mid \emptyset)$.

Definition 3. *The set of all independence relations of P , written $\mathcal{I}(P)$ is defined as the set of relations $(A \perp B \mid C)$ that hold in P .*

As said in Item 2, we will see that the statement “ P factorizes according to G ” is indeed equivalent to some statements about $\mathcal{I}(P)$. More precisely, we shall see (eventually under mild positivity conditions on p) that Definition 1 is equivalent to the following two definitions:

Definition 4 (Local Markov property). *P satisfies the local Markov property with respect to G iff*

$$\mathcal{I}_{\text{loc}}(G) := \{(X_i \perp \text{NonDesc}(X_i) \mid \text{Pa}(X_i)) : i = 1, \dots, d\} \subset \mathcal{I}(P).$$

Definition 5 (Global Markov property). *P satisfies the global Markov property with respect to G iff*

$$\mathcal{I}(G) := \{(\mathbf{X}, \mathbf{Y} \mid \mathbf{Z}) : \text{d-sep}_G(\mathbf{X}, \mathbf{Y} \mid \mathbf{Z})\} \subset \mathcal{I}(P).$$

To fully understand the global Markov property, we must first define d-separation, which will be addressed later; for now, we just point that the set $\mathcal{I}(G)$ is richer than the set $\mathcal{I}_{\text{loc}}(G)$ and typically contains more independence relations.

1.1 Local Markov property and factorization

We wish to show that $\mathcal{I}_{\text{loc}}(G) \subset \mathcal{I}(P)$ (aka. *local Markov property*) iff P factorizes according to G .

Theorem 1. *The following statements are equivalent:*

1. P factorizes according to G ;
2. P satisfies the local Markov property relative to G (ie. $\mathcal{I}_{\text{loc}}(G) \subset \mathcal{I}(P)$).

Proof. (1) \implies (2) See [KF09, Theorem 3.2] or [Lau96, Section 3.2].

(2) \implies (1) [KF09, Theorem 3.1] or [Lau96, Section 3.2]. □

Example 2. *Consider $A \rightarrow B \rightarrow C, A \rightarrow D \rightarrow E$. Then $(C \perp \{A, D, E\} \mid B)$ is a local independence relation. On the other hand, if P factorizes according to this graph, then $(A \perp \{C, E\} \mid \{B, D\})$ which is non-local (but we will see is captured by the global Markov property and d-separation).*

0129
0130
0131
0132
0133
0134
0135
0136
0137
0138
0139
0140
0141
0142
0143
0144
0145
0146
0147
0148
0149
0150
0151
0152
0153
0154
0155
0156
0157
0158
0159
0160
0161
0162
0163
0164
0165
0166
0167
0168
0169
0170
0171
0172
0173
0174
0175
0176
0177
0178
0179
0180
0181
0182
0183
0184
0185
0186
0187
0188
0189
0190
0191
0192

1.2 Global Markov property, d -separation

As we discussed, a graph structure G encodes a certain set of conditional independence assumptions $\mathcal{I}_{\text{loc}}(G)$. Knowing only that a distribution P factorizes over G , we can conclude that it satisfies the local Markov property. An immediate question is whether there are other independencies that we can “read-off” directly from G . That is, are there other independencies that hold for every distribution P that factorizes over G ?

1.2.1 Basic building blocks of G and intuitions

Here we take inspiration from Brady Neal’s sections 3.5 and 3.6; our goal is to understand the ideas motivating the definition of d -separation.

We use the example to define various motifs of interest, and give a hint on why they are interesting.

Example 3 (Chains). *A chain is a motif of the form $X_1 \rightarrow X_2 \rightarrow X_3$. Consider the graph on 3 vertices which is a chain and P factorizing according to it (equivalently P satisfies the local Markov property relative to G). Usually X_1 and X_3 are dependent. But if we look at $X_1, X_3 \mid X_2$, we can see that*

$$p_{123}(x_1, x_2, x_3) = p_1(x_1)p_{2|1}(x_2 \mid x_1)p_{3|2}(x_3 \mid x_2) \quad (2)$$

so that $X_1, X_3 \mid X_2 = x_2$ admits the density

$$p_{13|2}(x_1, x_3 \mid x_2) = \frac{p_{123}(x_1, x_2, x_3)}{p_2(x_2)} \quad (3)$$

$$= \frac{p_1(x_1)p_{2|1}(x_2 \mid x_1)}{p_2(x_2)} p_{3|2}(x_3 \mid x_2) \quad (4)$$

$$= p_{13|2}(x_1, x_3 \mid x_2)p_{3|2}(x_3 \mid x_2). \quad (5)$$

In other words, $(X_1 \perp X_3 \mid X_2)$ holds in P . We shall say that X_2 blocks the path between $X_1 \rightarrow X_2 \rightarrow X_3$.

Example 4 (Forks). *A fork is a motif of the form $X_1 \leftarrow X_2 \rightarrow X_3$. It is easily seen that forks and chains encodes the same type of conditional independencies.*

Example 5 (Colliders and immoralities). *A collider is a motif of the form $X_1 \rightarrow X_2 \leftarrow X_3$, if there is no edge between X_1 and X_3 , the motif is called an immorality¹. Consider the graph on 3 vertices which is an immorality and P factorizing according to it. Then,*

¹The term “immoral” is used humorously to indicate that these parents have not formed a relationship, even though they share a child.

0193 we can see that $X_1 \perp X_3$ holds in P because

0194
0195
$$p_{13}(x_1, x_3) = \int p_{123}(x_1, x_2, x_3) d\mu_2(x_2) \tag{6}$$

0196
0197
0198
$$= \int p_{123}(x_1, x_2, x_3) d\mu_2(x_2) \tag{7}$$

0199
0200
0201
$$= \int p_1(x_1) p_3(x_3) p_{2|13}(x_2 | x_1, x_3) d\mu_2(x_2) \tag{8}$$

0202
0203
0204
$$= p_1(x_1) p_3(x_3) \int p_{2|13}(x_2 | x_1, x_3) d\mu_2(x_2) \tag{9}$$

0205
0206
0207
$$= p_1(x_1) p_3(x_3). \tag{10}$$

0208
0209 But, if we look at $X_1, X_3 | X_2$, then

0210
0211
$$p_{13|2}(x_1, x_3 | x_2) = \frac{p_{123}(x_1, x_2, x_3)}{p_2(x_2)} = p_1(x_1) p_2(x_2) \frac{p_{2|13}(x_2 | x_1, x_3)}{p_2(x_2)}. \tag{11}$$

0212
0213
0214 So oddly-enough, conditional on X_2 the variables X_1 and X_3 may become dependent if
0215 X_2 is really depending on (X_1, X_3) . Brady Neal has the following “concrete” example:
0216 An example is the easiest way to see why this is the case. Imagine that you’re out dating
0217 men, and you notice that most of the nice men you meet are not very good-looking,
0218 and most of the good-looking men you meet are jerks. It seems that you have to choose
0219 between looks and kindness. In other words, it seems like kindness and looks are negatively
0220 associated. However, what if I also told you that there is an important third variable here:
0221 availability (whether men are already in a relationship or not)? And what if I told you that
0222 a man’s availability is largely determined by their looks and kindness; if they are both good-
0223 looking and kind, then they are in a relationship. The available men are the remaining
0224 ones, the ones who are either not good-looking or not kind. You see an association
0225 between looks and kindness because you’ve conditioned on a collider (availability). You’re
0226 only looking at men who are not in a relationship. The structure of this example is
0227 looks \rightarrow availability \leftarrow kindness.
0228
0229
0230
0231
0232
0233

0234 **Example 6** (Descendant of immoralities). Similarly, conditioning on a descendant of
0235 an immorality can induce association in between the parents of the collider. The intuition
0236 is that if we learn something about a collider’s descendant, we usually also learn some-
0237 thing about the collider itself because there is a direct causal path from the collider to its
0238 descendants, and we know that nodes in a chain are eventually dependent.
0239
0240
0241

0242 What is important to retain here is that conditioning on the middle vertex of a chain
0243 or a fork can “block” the flow of dependencies along it. In contrast, conditioning on the
0244 middle vertex of an immorality can “unblock” the flow.
0245
0246
0247

0248 **1.2.2 Paths, blocked paths, and d -separation**

0249
0250 A (undirected) path is a sequence of distinct vertices $(X_{i_1}, \dots, X_{i_m})$, $m \geq 2$, such that
0251 for each $k = 1, \dots, m - 1$ there is an edge $X_{i_k} \rightarrow X_{i_{k+1}}$ or $X_{i_k} \leftarrow X_{i_{k+1}}$.
0252

0253 **Definition 6** (Blocked path). A (undirected) path $(X_{i_1}, \dots, X_{i_m})$ is blocked by a set of
0254 vertices Z (not containing X_{i_1} nor X_{i_m}) if and only if:
0255
0256

0257
0258
0259
0260
0261
0262
0263
0264
0265
0266
0267
0268
0269
0270
0271
0272
0273
0274
0275
0276
0277
0278
0279
0280
0281
0282
0283
0284
0285
0286
0287
0288
0289
0290
0291
0292
0293
0294
0295
0296
0297
0298
0299
0300
0301
0302
0303
0304
0305
0306
0307
0308
0309
0310
0311
0312
0313
0314
0315
0316
0317
0318
0319
0320

- $(X_{i_1}, \dots, X_{i_m})$ contains a chain $X_{i_{k-1}} \rightarrow X_{i_k} \rightarrow X_{i_{k+1}}$ or a fork $X_{i_{k-1}} \leftarrow X_{i_k} \rightarrow X_{i_{k+1}}$ such that $X_{i_k} \in \mathbf{Z}$, or
- $(X_{i_1}, \dots, X_{i_m})$ contains an immorality $X_{i_{k-1}} \rightarrow X_{i_k} \leftarrow X_{i_{k+1}}$ such that $X_{i_k} \notin \mathbf{Z}$ and no descendant of X_{i_k} is in \mathbf{Z} .

[draw picture! in particular compare the difference between immorality and a collider $X \rightarrow Y \leftarrow Z$ with an edge between X and Z]

Definition 7 (d-separation). Let \mathbf{X} , \mathbf{Y} and \mathbf{Z} be three disjoint subsets of vertices. \mathbf{X} and \mathbf{Y} are d-separated by \mathbf{Z} if every path between vertices of \mathbf{X} and \mathbf{Y} is blocked by \mathbf{Z} . We then write

$$\text{d-sep}(\mathbf{X}, \mathbf{Y} \mid \mathbf{Z}).$$

1.2.3 The global Markov property

Recall $\mathcal{I}(G) := \{(\mathbf{X}, \mathbf{Y} \mid \mathbf{Z}) : \text{d-sep}_G(\mathbf{X}, \mathbf{Y} \mid \mathbf{Z})\} \subset \mathcal{I}(P)$.

Theorem 2. If $p > 0$ the following are equivalent:

1. P factorizes according to G ;
2. P satisfies the global Markov property relative to G (ie. $\mathcal{I}(G) \subset \mathcal{I}(P)$).

Proof. (1) \implies (2) [Lau96, Corollary 3.23].

(1) \implies (2) It is enough to observe that for every vertex X_i the sets $\{X_i\}$ and $\text{NonDesc}(X_i)$ are d-separated by $\text{Pa}(X_i)$. In other words the global Markov property implies the local Markov property, which in turn implies (1) by Theorem 1. \square

So indeed, we have equivalence between factorization, local Markov, and global Markov properties. This tells us that in a Bayesian network, the DAG gives indication about the independence relations in $\mathcal{I}(G)$ (which might not be all of $\mathcal{I}(P)$, see below).

To make it interesting, however, we shall ensure that the global Markov condition is indeed stronger than the local Markov conditions. This can be seen because the parents of a node always d-separate the node from its non-descendants.

Interestingly, there are efficient algorithms to find out if two sets are d-separated given a third set [KF09, Algorithm 3.1]; so given the DAG we can immediately test if some independence relation is in $\mathcal{I}(G)$.

1.3 Faithfulness, aka converse global Markov property

Definition 8 (Faithfulness). P is faithful to G if $\mathcal{I}(P) \subset \mathcal{I}(G)$.

It is interesting to wonder if factorization implies faithfulness. This is because since factorization implies global Markov, then we would have that factorization implies $\mathcal{I}(G) = \mathcal{I}(P)$ and be happy to have encoded in the DAG all the independence relations of P . Unfortunately, factorization does not imply faithfulness, as seen in the examples below:

Example 7. Suppose $X_1 \perp X_2$. Yet, the law of (X_1, X_2) factorizes according to the graph $X_1 \rightarrow X_2$, but $\{X_1\}$ and $\{X_2\}$ are not d-separated, so we cannot read the independence relation from the DAG $X_1 \rightarrow X_2$.

0321
0322
0323
0324
0325
0326
0327
0328
0329
0330
0331
0332
0333
0334
0335
0336
0337
0338
0339
0340
0341
0342
0343
0344
0345
0346
0347
0348
0349
0350
0351
0352
0353
0354
0355
0356
0357
0358
0359
0360
0361
0362
0363
0364
0365
0366
0367
0368
0369
0370
0371
0372
0373
0374
0375
0376
0377
0378
0379
0380
0381
0382
0383
0384

Example 8. See also [PJS17, Example 6.34].

This shows that in general we cannot read all the independence from the DAG in a Bayesian network. But, faithfulness is an important concept since, factorization + faithfulness implies that $\mathcal{I}(G) = \mathcal{I}(P)$; in other words that d -separation permits to characterize all independence relations.

1.4 Completeness of d -separation

The previous section shows that there can exist independence relations in $\mathcal{I}(P)$ that we cannot read from $\mathcal{I}(G)$; ie. some independence relations cannot be uncovered using d -separation. It is natural to ask whether or not d -separation is the best we can do. Indeed it is.

Theorem 3 (Completeness). *If \mathbf{X} and \mathbf{Y} are not d -separated given \mathbf{Z} in G , then there exists a distribution P that factorizes according to G and in which \mathbf{X} and \mathbf{Y} are dependent.*

Proof. [KF09, Theorem 3.4] □

We can view the completeness result as telling us that our definition of $\mathcal{I}(G)$ is the maximal one. For any independence relation n that is not a consequence of d -separation in G , we can always find a counterexample distribution P that factorizes over G .

1.5 Markov equivalence

Definition 9. *Two graphs G_1 and G_2 are Markov equivalent if $\mathcal{I}(G_1) = \mathcal{I}(G_2)$.*

Example 9 (Chains and forks). $X_1 \rightarrow X_2 \rightarrow X_3$, $X_1 \leftarrow X_2 \leftarrow X_3$, and $X_1 \leftarrow X_3 \rightarrow X_2$ are Markov equivalent (notice the importance for causality: in the causal interpretation, they are very different graphs since they don't suppose the same cause-effects relations, but probabilistically, they represent the same set of independence relations).

So the previous graphs are Markov equivalent, but they are not equivalent to the immorality $X_1 \rightarrow X_2 \leftarrow X_3$. This is because the immorality implies that $X_1 \perp X_3$ (see for instance Example 5), but there are distributions P that factorizes according to chains or forks in which X_1 and X_3 are dependent.

Theorem 4. *Two DAG G_1 and G_2 are Markov equivalent if and only if they have the same skeleton and the same set of immoralities.*

Proof. [KF09, Section 3.3.4]. □

1.6 Minimality

Definition 10 (Minimality). *If P factorizes according to G , but doesn't factorize according to any proper sgraph of G , then G is minimal.*

This also can be understood has: if we were to remove any edges from the DAG, P would factorizes according to the graph with the removed edges.

Theorem 5. *For every P there exists a minimal DAG.*

0385
0386
0387
0388
0389
0390
0391
0392
0393
0394
0395
0396
0397
0398
0399
0400
0401
0402
0403
0404
0405
0406
0407
0408
0409
0410
0411
0412
0413
0414
0415
0416
0417
0418
0419
0420
0421
0422
0423
0424
0425
0426
0427
0428
0429
0430
0431
0432
0433
0434
0435
0436
0437
0438
0439
0440
0441
0442
0443
0444
0445
0446
0447
0448

Proof. Let us exhibit such a DAG. Starting point is to notice that we can always write

$$p(x_1, \dots, x_d) = p_1(x_1)p_{2|1}(x_2 | x_1) \dots p_{d|1\dots(d-1)}(x_d | x_1, \dots, x_{d-1}).$$

Let $\text{Pa}_1 = \emptyset$. For each $j = 2, \dots, d$ find the minimal subset Pa_j of $\{1, \dots, j - 1\}$ for which

$$p_{j|1\dots(j-1)}(X_j | X_1, \dots, X_{j-1}) = p_{j|\text{Pa}_j}(X_j | (X_k)_{k \in \text{Pa}_j}) \quad P\text{-as.}$$

Now build the graph on $\{X_1, \dots, X_d\}$ recursively starting from isolated vertex X_1 , then adding X_2 and edges $X_2 \rightarrow X_1$ if $\text{Pa}_2 \neq \emptyset$, then adding X_3 and edges $X_3 \rightarrow X_k$ if $k \in \text{Pa}_3$, etc. Clearly this graph is a DAG, factors P , and is minimal. \square

Interestingly, the constructive proof of the previous theorem immediately shows that minimal DAGs exist but not guaranteed to be unique. In fact, we have exhibited a DAG using a specific ordering of the random variables, but choosing a different ordering can lead to another minimal DAG. We illustrate this in the following example:

Example 10. Taken from [KF09, Sections 3.2.1 and 3.4.1 for details]. We consider a set of random variables $\{I, D, G, S, L\}$ (Intelligence, Difficulty, Grade, SAT, Letter), and P whose density is given by

$$p(i, d, g, s, l) = p_I(i)p_D(d)p_{G|ID}(g | i, d)p_{S|I}(s | i)p_{L|G}(l | g) \quad (12)$$

The following are the minimal DAGs obtained from the algorithm used in the proof of Theorem 5, using the ordering (I, D, G, S, L) , (L, S, G, I, D) , and (L, D, S, I, G) .

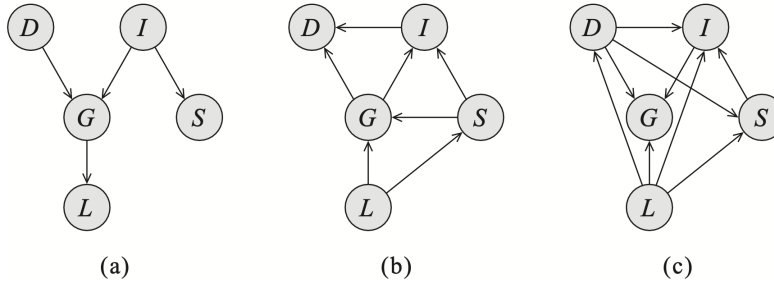


Figure 3.8 Three minimal I-maps for P_{student} , induced by different orderings: (a) D, I, S, G, L ; (b) L, S, G, I, D ; (c) L, D, S, I, G .

Note that, perhaps surprisingly, the three minimal DAGs G_1, G_2, G_3 in the previous example are not Markov equivalent (they don't have the same skeleton). There is nothing contradictory here: this is because they encode different subsets of $\mathcal{I}(G_1), \mathcal{I}(G_2), \mathcal{I}(G_3) \subset \mathcal{I}(P)$ but $\mathcal{I}(G_1) \neq \mathcal{I}(G_2) \neq \mathcal{I}(G_3)$.

This also shows that minimality and faithfulness are distinct notions, since for P faithful to G we must have $\mathcal{I}(G) = \mathcal{I}(P)$. Minimality is, however, a necessary condition for faithfulness.

Proposition 1. If P is faithful and factorizes according to G , then G is minimal.

Proof. A sketch of proof is given in [PJS17, Proposition 6.35], but it relies on the premise that two vertices with no edge between them can always be d-separated, which is still unclear to me how to prove formally. \square

2 Causal models, causes, effects, interventions

A “causal network” is a Bayesian network. The distinction between causal network and Bayesian network is purely semantic, they are mathematically the same object. The distinction stem from the fact that in the causal network, the edges are interpreted as representing a cause-effect relation. It is indeed the DAG that defines what is a *cause* and an *effect*.

Definition 11 (Cause). *In a causal network (G, P) , X_i is a cause of X_j if there is a directed path from X_i to X_j . It is a direct cause if there is an edge $X_i \rightarrow X_j$.*

2.1 “Classical” causal modeling principles

What distinguish a causal network from a Bayesian network is the meaning of the arrows. This implies that certain principles are assumed when modeling a causal networks. Those principles are not mathematical, but rather philosophical.

Principle 1 (Reichenbach’s common cause principle). *If two random variables X and Y are statistically dependent, then either X causes Y , or Y causes X , or there exists a third variable Z that causally influences both. Furthermore, this variable Z screens X and Y from each other in the sense that $X \perp Y \mid Z$.*

Principle 2 (Principle of independent mechanisms). *The Principle of independent mechanisms posits that the mechanisms governing the generation of a system’s variables are autonomous and do not influence each other. Specifically, the causal process that determines the effect of a variable X on another variable Y is independent of the processes governing X itself. This principle implies that changes to one causal mechanism (e.g., intervening on X) should not alter the mechanisms governing the remaining variables.*

The principle of independent mechanisms underpins many causal inference methods by ensuring that the causal structure can be disentangled into distinct, modular components, facilitating robust predictions and transferability across different contexts.

The following is an application example of the principle of independent mechanisms.

Example 11. *Borrowed from [PJS17, Section 2.1]. Suppose we have estimated the joint density $p(a, t)$ of the altitude A and the average annual temperature T of a sample of cities in some country. Consider the following ways of expressing $p(a, t)$:*

$$p(a, t) = p(a \mid t)p(t) = p(t \mid a)p(a). \quad (13)$$

The first decomposition describes T and the conditional $A \mid T$. It corresponds to a factorization of $p(a, t)$ according to the graph $T \rightarrow A$. The second decomposition corresponds to a factorization according to $A \rightarrow T$. Can we decide which of the two structures is the causal one (i.e., in which case would we be able to think of the arrow as causal)?

A first idea is to consider the effect of interventions. Imagine we could change the altitude A of a city by some hypothetical mechanism that raises the grounds on which the city is built. Suppose that we find that the average temperature decreases. Let us next imagine that we devise another intervention experiment. This time, we do not change the altitude, but instead we build a massive heating system around the city that raises

0513 the average temperature by a few degrees. Suppose we find that the altitude of the city is
0514 unaffected.

0515 Intervening on A has changed T , but intervening on T has not changed A . We would
0516 thus reasonably prefer $A \rightarrow T$ as a description of the causal structure.

0517 Why do we find this description of the effect of interventions plausible, even though
0518 the hypothetical intervention is hard or impossible to carry out in practice? If we change
0519 the altitude A , then we assume that the physical mechanism $p(t | a)$ responsible for pro-
0520 ducing an average temperature is still in place and leads to a changed T . This would
0521 hold true independent of the distribution from which we have sampled the cities, and thus
0522 independent of $p(a)$. Austrians may have founded their cities in locations subtly different
0523 from those of the Swiss, but the mechanism $p(t | a)$ would apply in both cases. If, on the
0524 other hand, we change T , then we have a hard time thinking of $p(a | t)$ as a mechanism
0525 that is still in place — we probably do not believe that such a mechanism exists in the
0526 first place. Given a set of different city distributions $p(a, t)$, while we could write them all
0527 as $p(a | t)p(t)$, we would find that it is impossible to explain them all using an invariant
0528 $p(a | t)$.

0529 Our intuition can be rephrased and postulated in two ways: If $A \rightarrow T$ is the correct
0530 causal structure, then

0531 (i) it is in principle possible to perform a localized intervention on A , in other words,
0532 to change $p(a)$ without changing $p(t | a)$, and

0533 (ii) $p(a)$ and $p(t | a)$ are autonomous, modular, or invariant mechanisms or objects in
0534 the world.

0535 2.2 Alternative causal modeling principles

0536 Occam razor's: minimum description length!

0537 2.3 Interventions

0538 3 Causal effects

0539 3.1 Nonparametric identifiability of causal effects

0540 3.2 Parametric identifiability

0541 4 SCM

0542 References

0543 [KF09] Daphne Koller and Nir Friedman. *Probabilistic Graphical Models: Principles*
0544 *and Techniques*. Adaptive Computation and Machine Learning. MIT Press,
0545 Cambridge, MA, 2009.

0546 [Lau96] Steffen L. Lauritzen. *Graphical Models*. Oxford Statistical Science Series. Oxford
0547 University Press, Oxford, New York, May 1996.

0577 [Pea22] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University
0578 Press, Cambridge New York, NY Port Melbourne New Delhi Singapore, second
0579 edition, reprinted with corrections edition, 2022.
0580
0581
0582 [PJS17] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of Causal*
0583 *Inference: Foundations and Learning Algorithms*. Adaptive Computation and
0584 Machine Learning. The MIT press, Cambridge, Mass, 2017.
0585
0586
0587
0588
0589
0590
0591
0592
0593
0594
0595
0596
0597
0598
0599
0600
0601
0602
0603
0604
0605
0606
0607
0608
0609
0610
0611
0612
0613
0614
0615
0616
0617
0618
0619
0620
0621
0622
0623
0624
0625
0626
0627
0628
0629
0630
0631
0632
0633
0634
0635
0636
0637
0638
0639
0640