

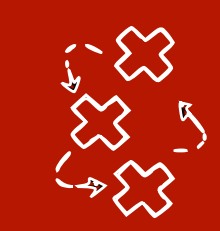


Deep learning 2: Causality & DL

1.3: Causal graphs

Lecturer: Sara Magliacane

UvA - Spring 2022

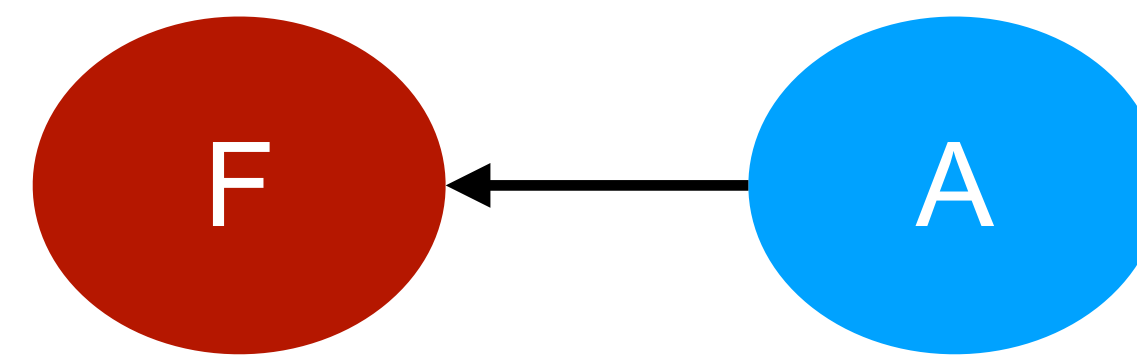
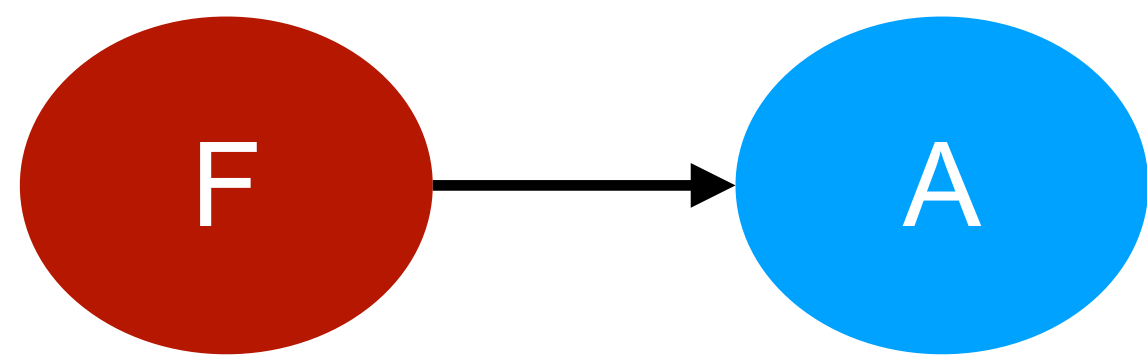


BNs vs causal BNs - example

- Fire (F) and Alarm (A) with $p(F, A)$ and $A \not\perp\!\!\!\perp F$ can be factorized as:

$$p(F, A) = p(F) p(A|F)$$

$$p(F, A) = p(A) p(F|A)$$

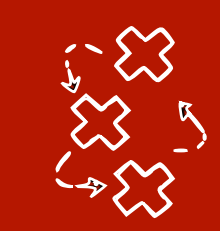


CAUSAL

NOT-CAUSAL

(lighting a fire triggers alarm)

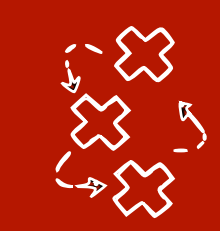
(triggering alarm does not light a fire)



The $do(X = x)$ operator [Pearl 2009]

- We introduce a new operator that can represent a **hypothetical intervention** on the whole population, i.e. a perturbation of the system:

$$do(X = x)$$



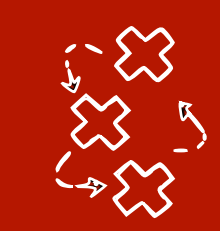
The $do(X = x)$ operator [Pearl 2009]

- We introduce a new operator that can represent a **hypothetical intervention** on the whole population, i.e. a perturbation of the system:

$$do(X = x)$$

- For a BN $((\mathbf{V}, \mathbf{E}), p)$ and intervention on $i \in \mathbf{V}$, we define the **interventional distribution**:

$$p(X_{\mathbf{V}} | do(X_i = x_i)), \text{ which in general } \neq p(X_{\mathbf{V}})$$



The $do(X = x)$ operator [Pearl 2009]

- We introduce a new operator that can represent a **hypothetical intervention** on the whole population, i.e. a perturbation of the system:

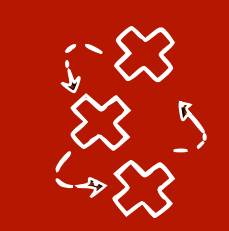
$$do(X = x)$$

- For a BN $((\mathbf{V}, \mathbf{E}), p)$ and intervention on $i \in \mathbf{V}$, we define the **interventional distribution**:

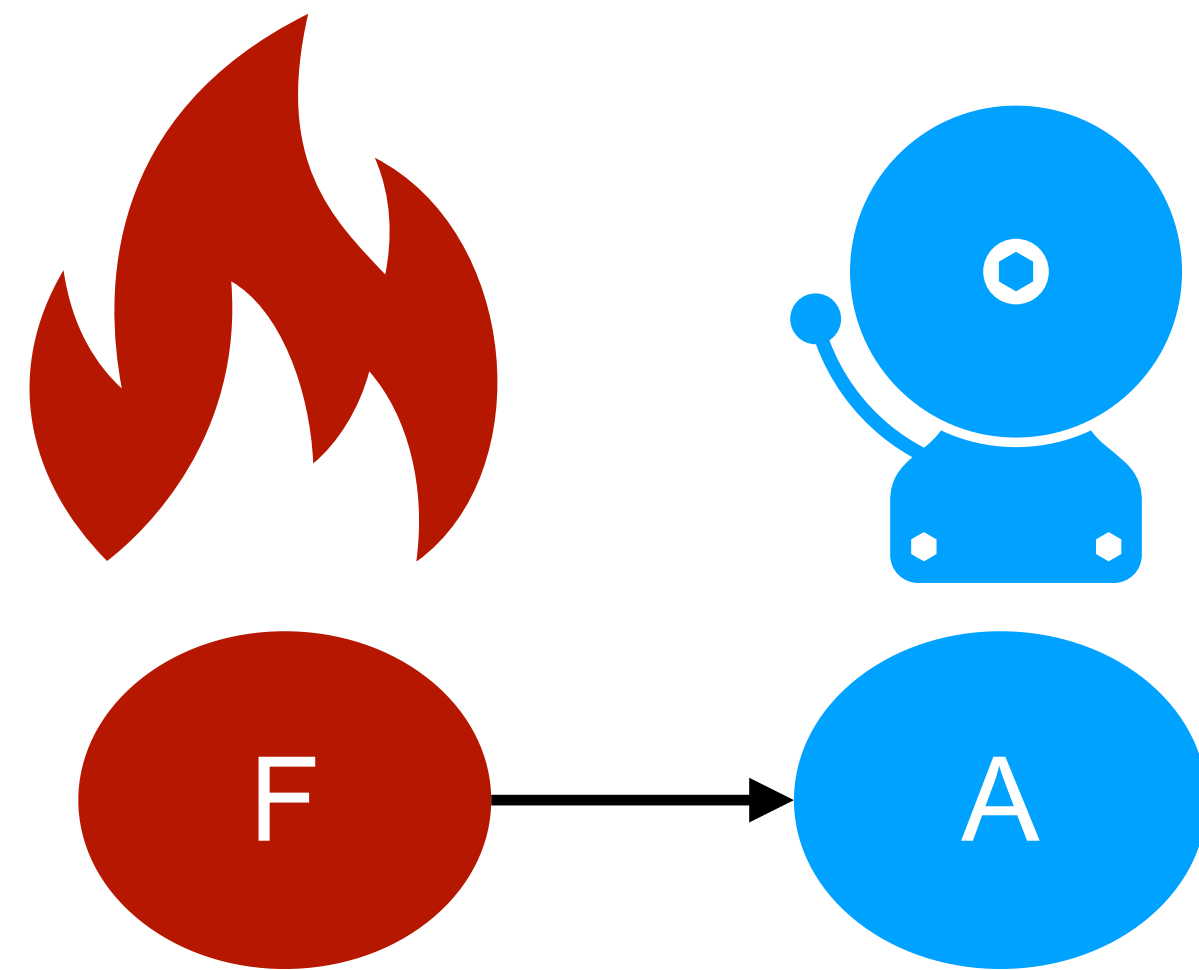
$$p(X_{\mathbf{V}} | do(X_i = x_i)), \text{ which in general } \neq p(X_{\mathbf{V}})$$

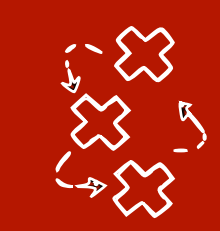
- We can also define conditional and marginal versions:

$$p(X_j | do(X_i = x_i)), \text{ which in general } \neq p(X_j)$$



Seeing is not doing





Causal Bayesian networks

- Given DAG $G = (\mathbf{V}, \mathbf{E})$ and distribution p , (G, p) is a Bayesian network if

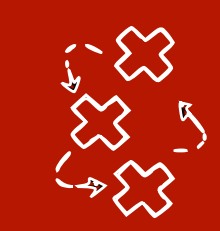
$$p(X_1, \dots, X_p) = \prod_{i \in \mathbf{V}} p(X_i | \mathbf{X}_{\text{Pa}_G(i)})$$

- If for any $\mathbf{W} \subset \mathbf{V}$:

$$p(X_{\mathbf{V}} | \text{do}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})) = \begin{cases} 0 & \text{if } X_{\mathbf{W}} \neq \tilde{x}_{\mathbf{W}} \\ \prod_{i \in \mathbf{V} \setminus \mathbf{W}} p(X_i | X_{\text{Pa}_G(i)}) & \text{if } X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}} \end{cases}$$

doesn't change *consistent with intervention*

then (G, p) is a **causal Bayesian network**



Causal Bayesian networks

- Given DAG $G = (\mathbf{V}, \mathbf{E})$ and distribution p , (G, p) is a Bayesian network if

$$p(X_1, \dots, X_p) = \prod_{i \in \mathbf{V}} p(X_i | \mathbf{X}_{\text{Pa}_G(i)})$$

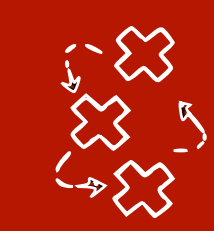
- If for any $\mathbf{W} \subset \mathbf{V}$:

$$p(X_{\mathbf{V}} | \text{do}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})) = \prod_{i \in \mathbf{V} \setminus \mathbf{W}} p(X_i | X_{\text{Pa}_G(i)}) \cdot \mathbf{1}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})$$

indicator function

(G, p) is a **causal Bayesian network**

Parents in G are now direct causes



Truncated factorisation formula [Pearl 2009]

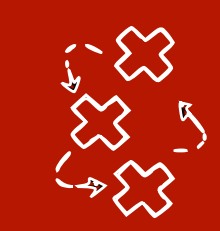
- If for any $\mathbf{W} \subset \mathbf{V}$:

$$p(X_{\mathbf{V}} | \text{do}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})) = \prod_{i \in \mathbf{V} \setminus \mathbf{W}} p(X_i | X_{\text{Pa}_G(i)}) \cdot \mathbf{1}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})$$

doesn't change

consistent with
intervention

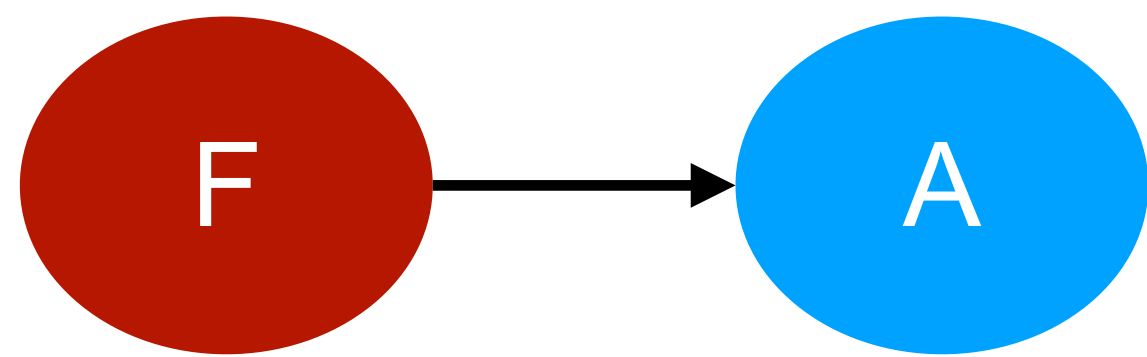
- Includes also observational data $\mathbf{W} = \emptyset$,
- Includes also multiple **intervention targets** $|\mathbf{W}| \geq 1$



BNs vs causal BNs - example

- Fire (F) and Alarm (A) with $p(F, A)$ and $A \not\perp\!\!\!\perp F$ can be factorized as:

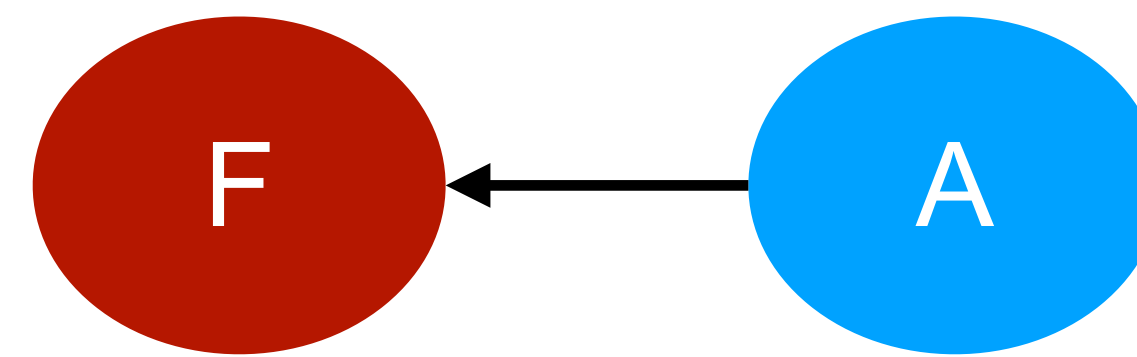
$$p(F, A) = p(F) p(A|F)$$



$do(A=1)$:

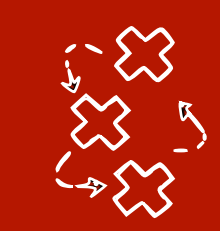
$$p(F, A | do(A=1)) = p(F) \cdot 1(A=1)$$

$$p(F, A) = p(A) p(F|A)$$

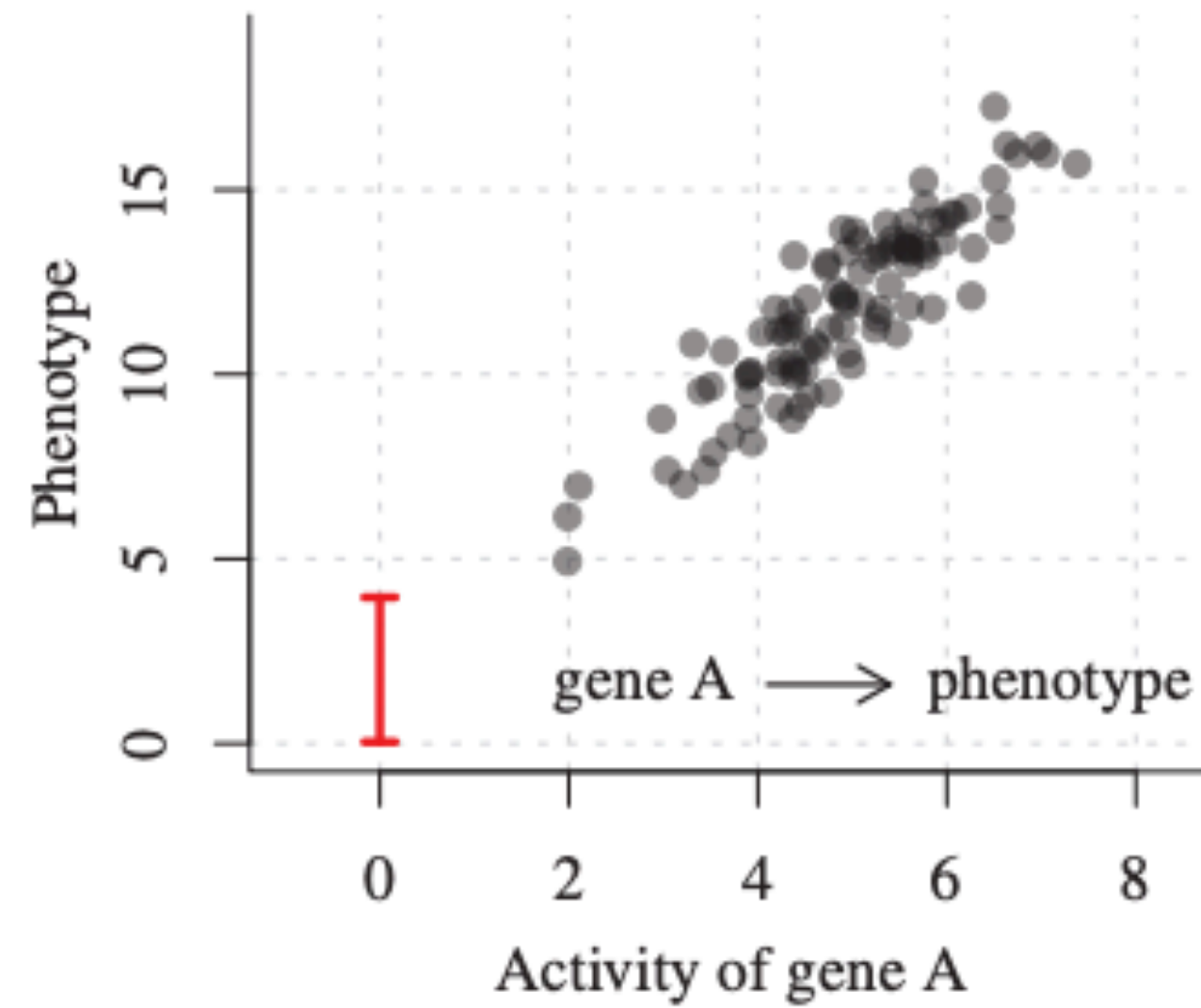
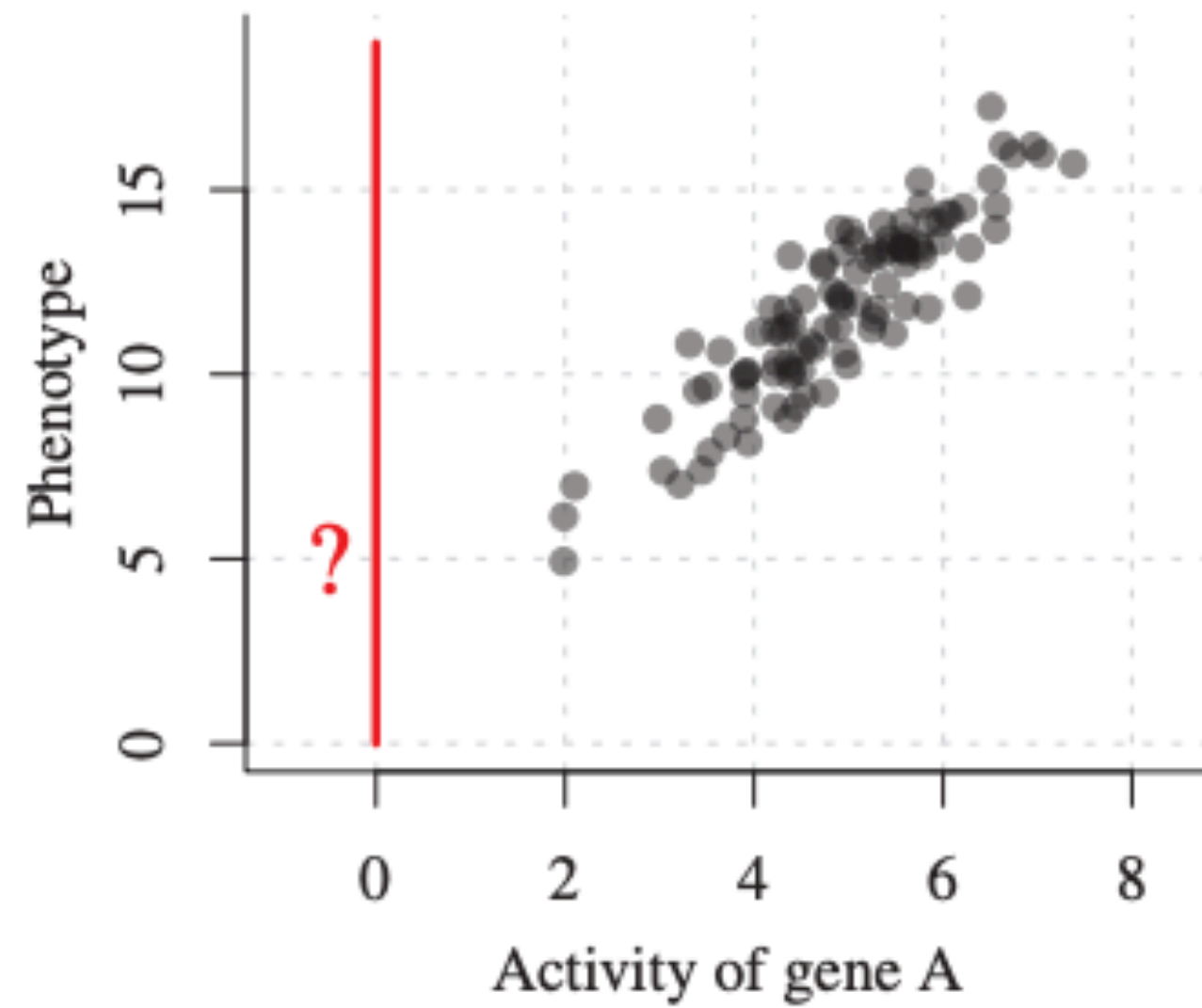


$do(A=1)$:

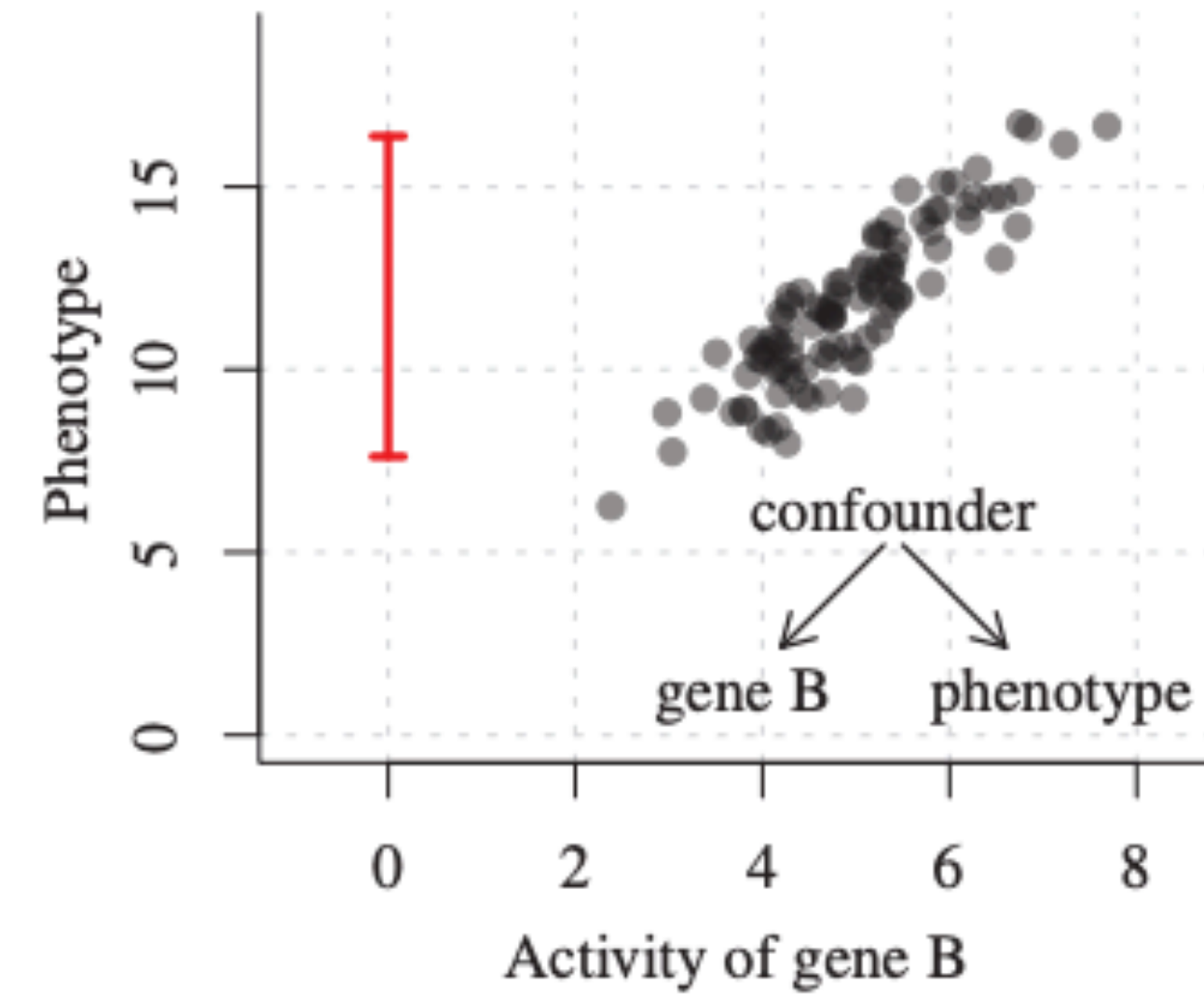
$$p(F, A | do(A=1)) = 1(A=1) \cdot p(F|A)$$



Another example of causal effect vs no effect for $do(X = 0)$

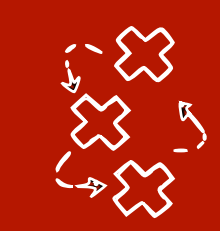


$do(A = 0)$



$do(B = 0)$

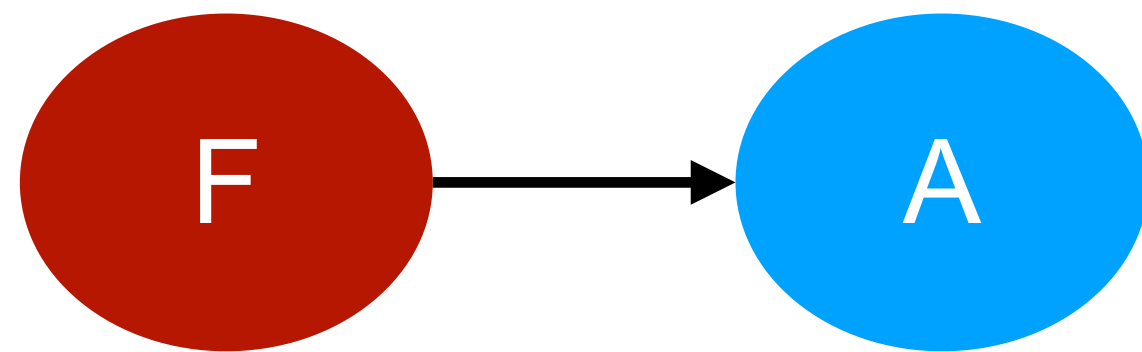
Fig 1.4 in Elements of Causal Inference (http://web.math.ku.dk/~peters/jonas_files/ElementsOfCausalInference.pdf)

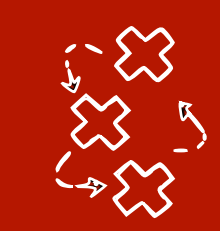


Mutilated/manipulated graphs

Graphically $p(X_{\mathbf{V}} | \text{do}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})) = \prod_{i \in \mathbf{V} \setminus \mathbf{W}} p(X_i | X_{\text{Pa}_G(i)}) \cdot \mathbf{1}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})$ can

be represented as cutting the incoming edges to $X_{\mathbf{W}}$

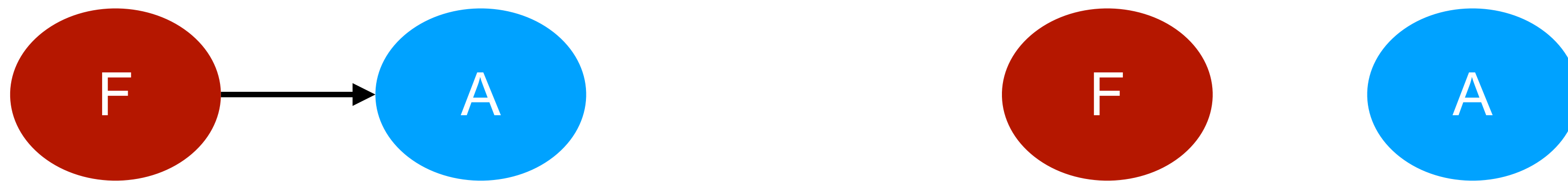


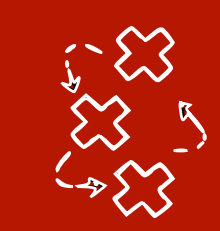


Mutilated/manipulated graphs

Graphically $p(X_{\mathbf{V}} | \text{do}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})) = \prod_{i \in \mathbf{V} \setminus \mathbf{W}} p(X_i | X_{\text{Pa}_G(i)}) \cdot \mathbf{1}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})$ can

be represented as cutting the incoming edges to $X_{\mathbf{W}}$

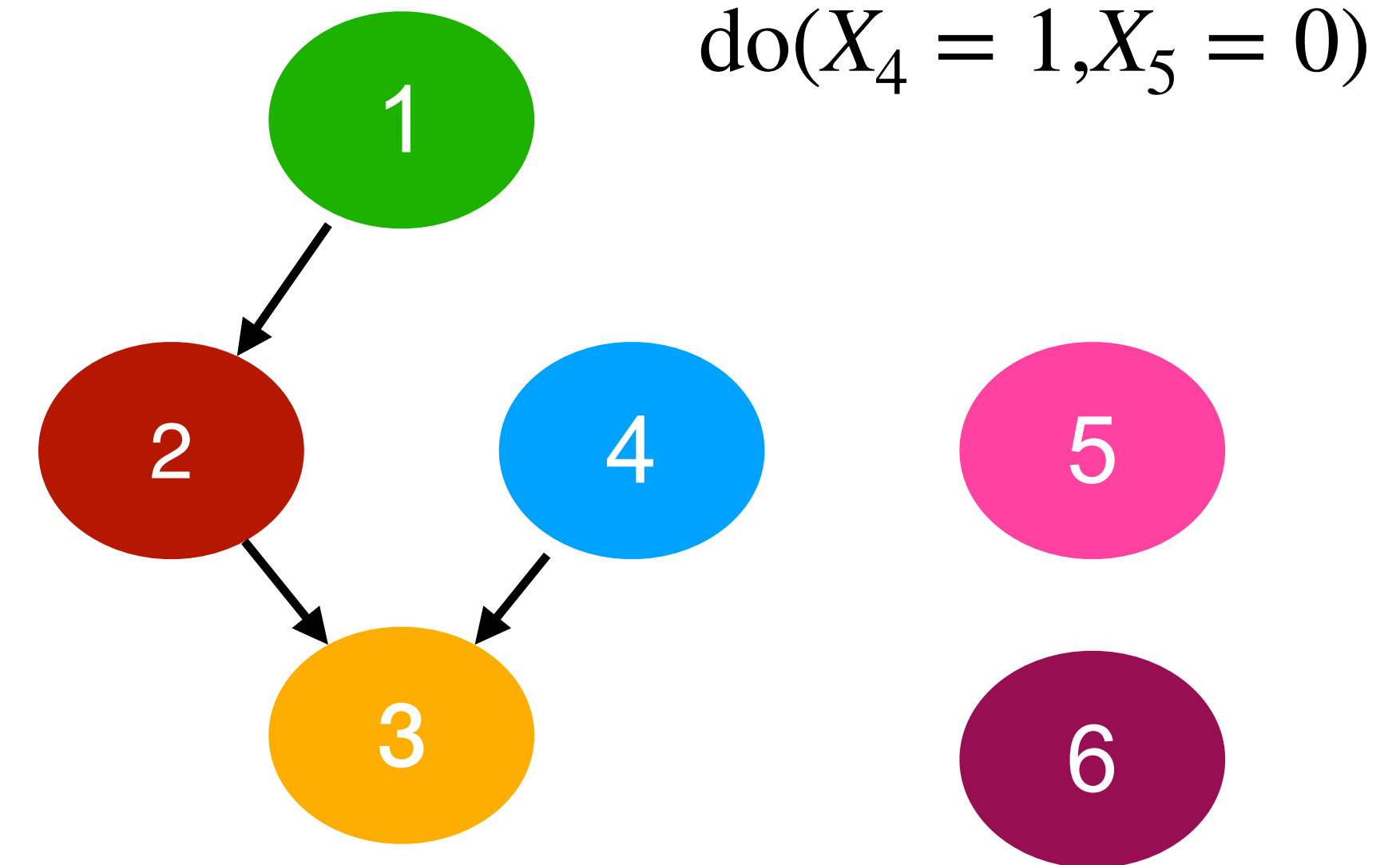
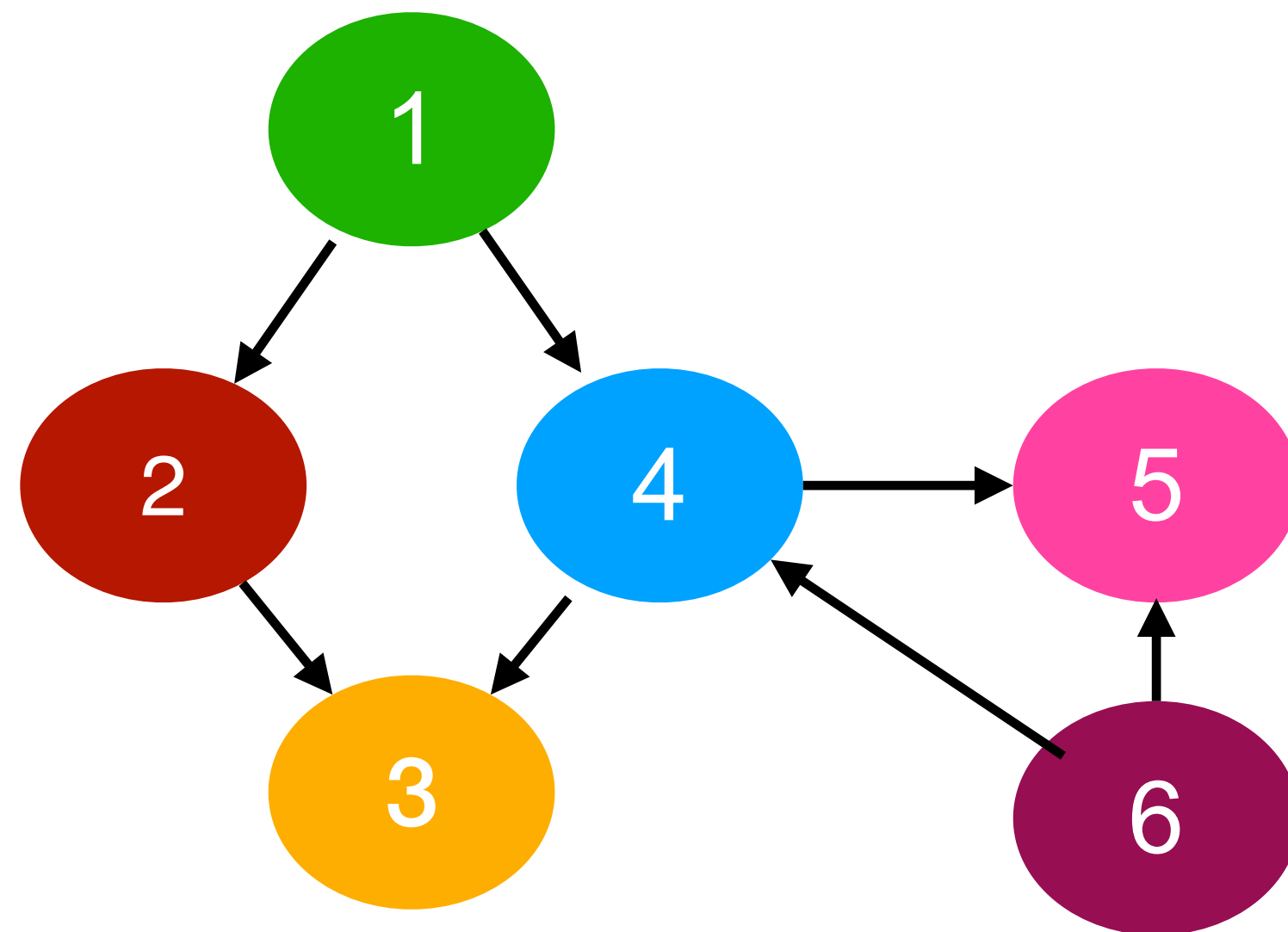


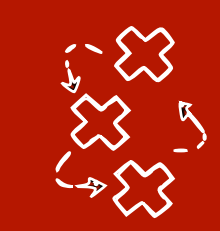


Mutilated/manipulated graphs

Graphically $p(X_{\mathbf{V}} | \text{do}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})) = \prod_{i \in \mathbf{V} \setminus \mathbf{W}} p(X_i | X_{\text{Pa}_G(i)}) \cdot \mathbf{1}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})$ can

be represented as cutting the incoming edges to $X_{\mathbf{W}}$



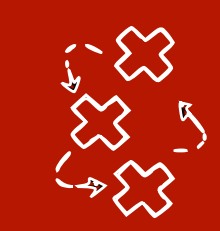


Perfect vs soft interventions

- We introduce a new operator that can represent a **hypothetical intervention** on the whole population, i.e. a perturbation of the system:

$$\text{do}(X_i = x_i) \text{ which changes } p(X_i | X_{\text{Pa}(i)}) \rightarrow \mathbf{1}(X_i = x_i)$$

- This is called a **perfect** (or surgical) **intervention**
- There are also other types of intervention, e.g. **soft interventions which change** $p(X_i | X_{\text{Pa}(i)}) \rightarrow \tilde{p}(X_i | X_{\text{Pa}(i)})$



Truncated factorisation formula [Pearl 2009]

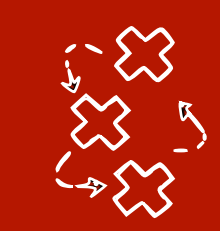
- If for any $\mathbf{W} \subset \mathbf{V}$:

$$p(X_{\mathbf{V}} | \text{do}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})) = \prod_{i \in \mathbf{V} \setminus \mathbf{W}} p(X_i | X_{\text{Pa}_G(i)}) \cdot \mathbf{1}(X_{\mathbf{W}} = \tilde{x}_{\mathbf{W}})$$



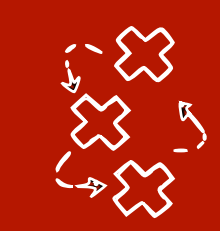
doesn't change from
the observational distr.

MODULARITY
ASSUMPTION



Causal mechanisms and Modularity

- In a causal BN (G, p) , each $p(X_i | \mathbf{X}_{\text{Pa}(i)})$ is the **causal mechanism of X_i**
- **Modularity assumption:** intervening on X_j will not change any causal mechanism $p(X_i | \mathbf{X}_{\text{Pa}(i)})$ for any $i \neq j$



Causal mechanisms and Modularity

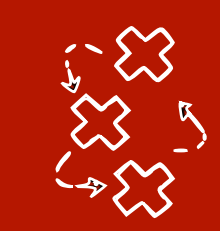
- In a causal BN (G, p) , each $p(X_i | \mathbf{X}_{\text{Pa}(i)})$ is the **causal mechanism of X_i**
- **Modularity assumption:** intervening on X_j will not change any causal mechanism $p(X_i | \mathbf{X}_{\text{Pa}(i)})$ for any $i \neq j$
- **Independent Causal Mechanism Principle:** the generative process is composed of **autonomous models** that do **not inform or influence** each other

Knowing
Changing

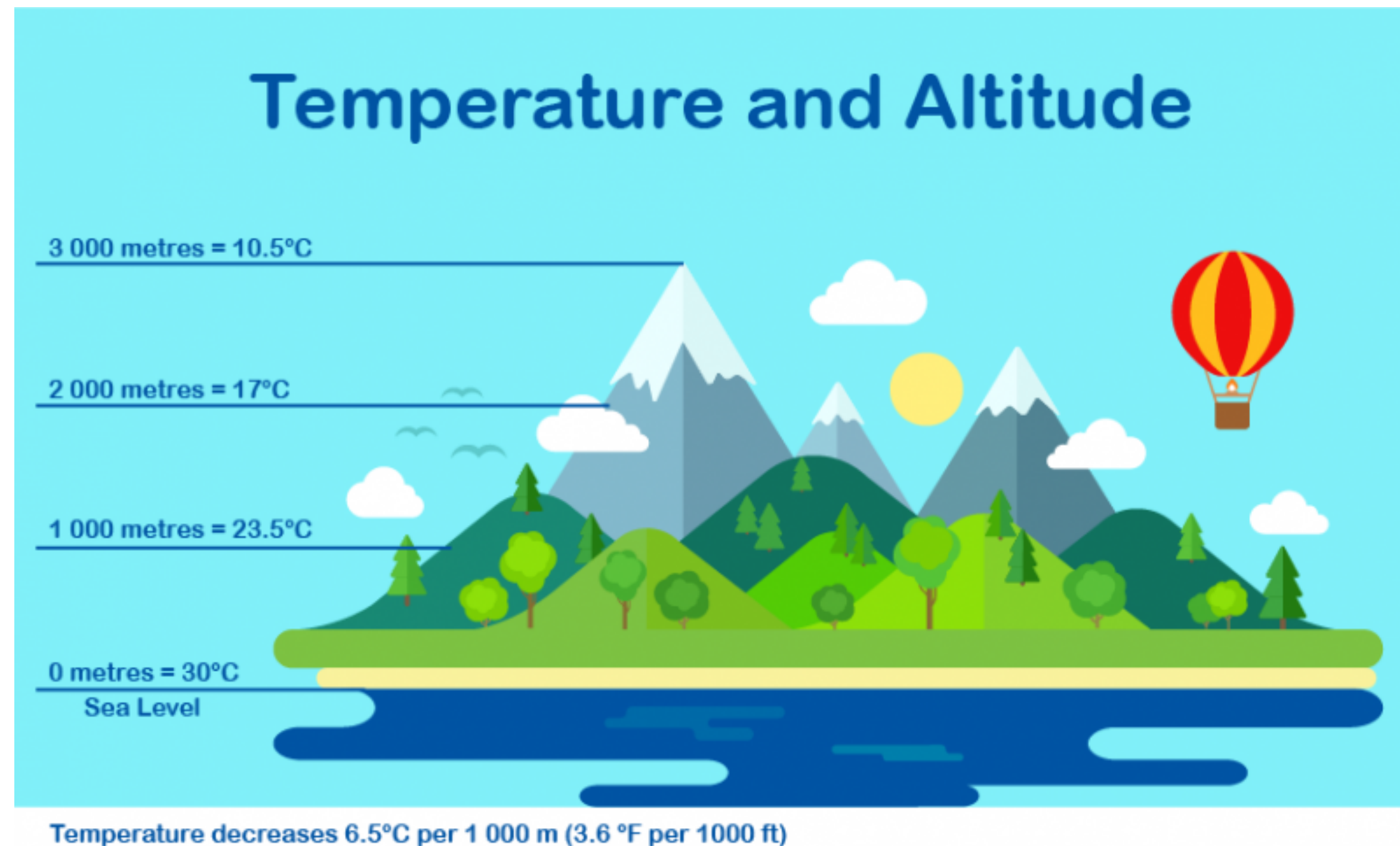
$$p(X_j | \mathbf{X}_{\text{Pa}(j)})$$

Does not give info
Does not change

$$p(X_i | \mathbf{X}_{\text{Pa}(i)}) \quad i \neq j$$

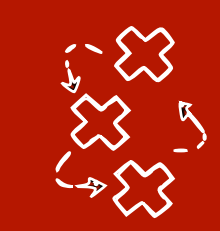


Modularity example

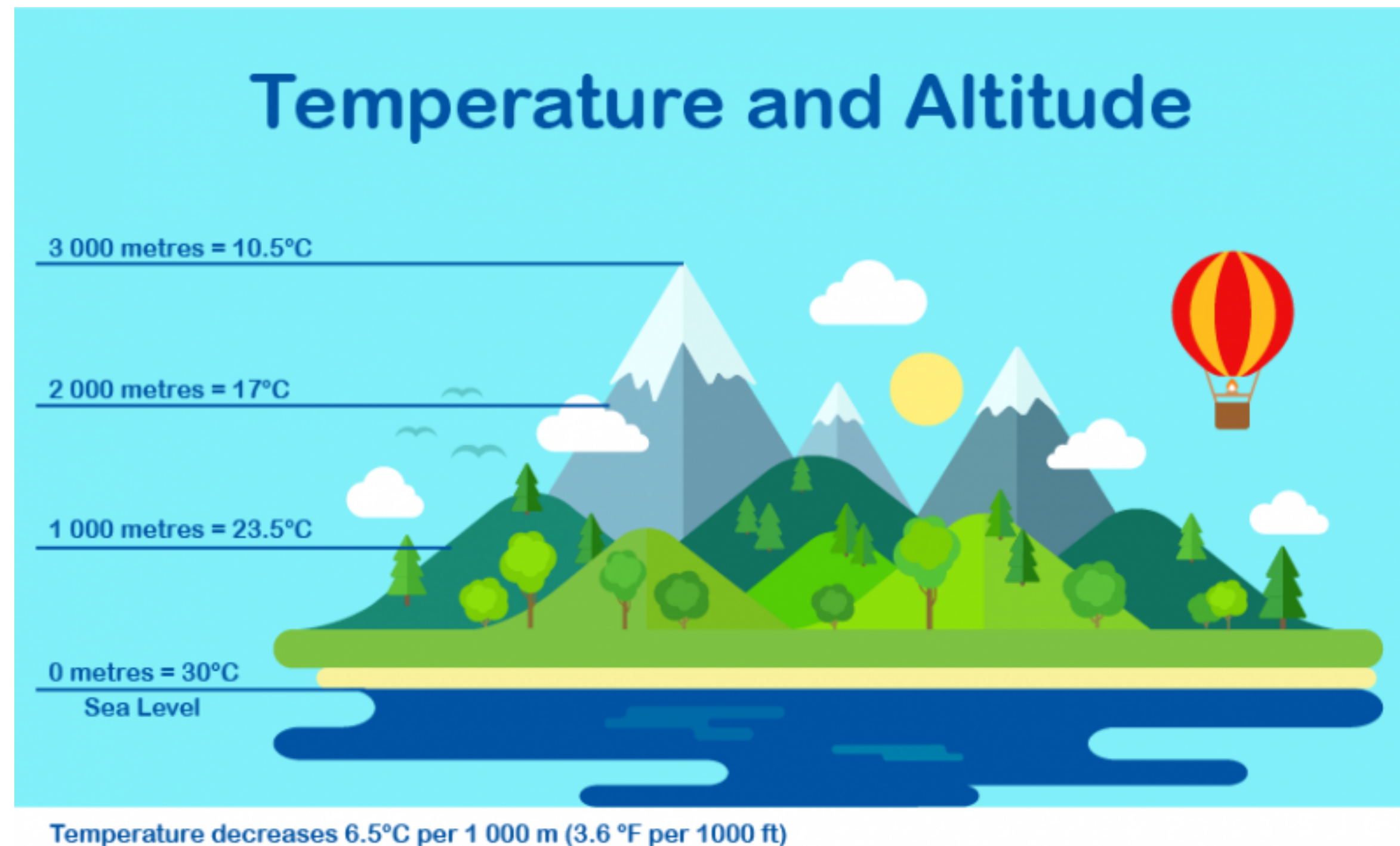


$$P(A, T) = P(T|A)P(A)$$

$$P(A, T) = P(A|T)P(T)$$



Modularity example

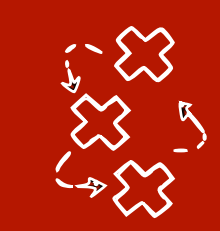


$$P(A, T) = P(T|A)P(A)$$

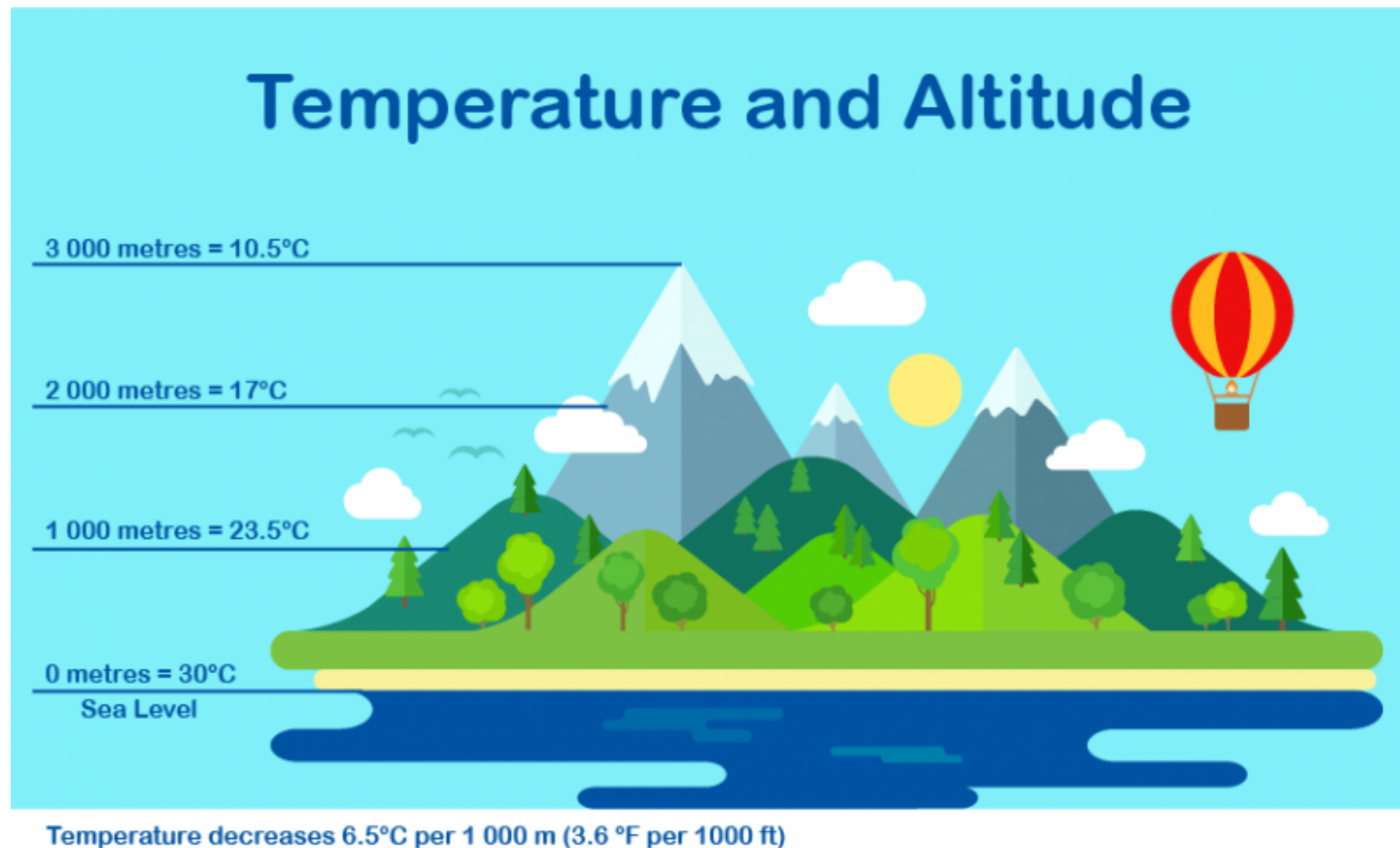
changing $P(A)$ does not
change $P(T|A)$

$$P(A, T) = P(A|T)P(T)$$

Changing $P(T)$ might
change $P(A|T)$



Modularity example



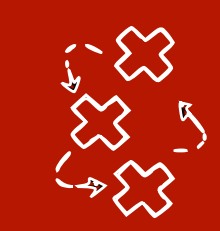
$P(T|A)$ is an invariant physical mechanism

$$P(A, T) = P(T|A)P(A)$$

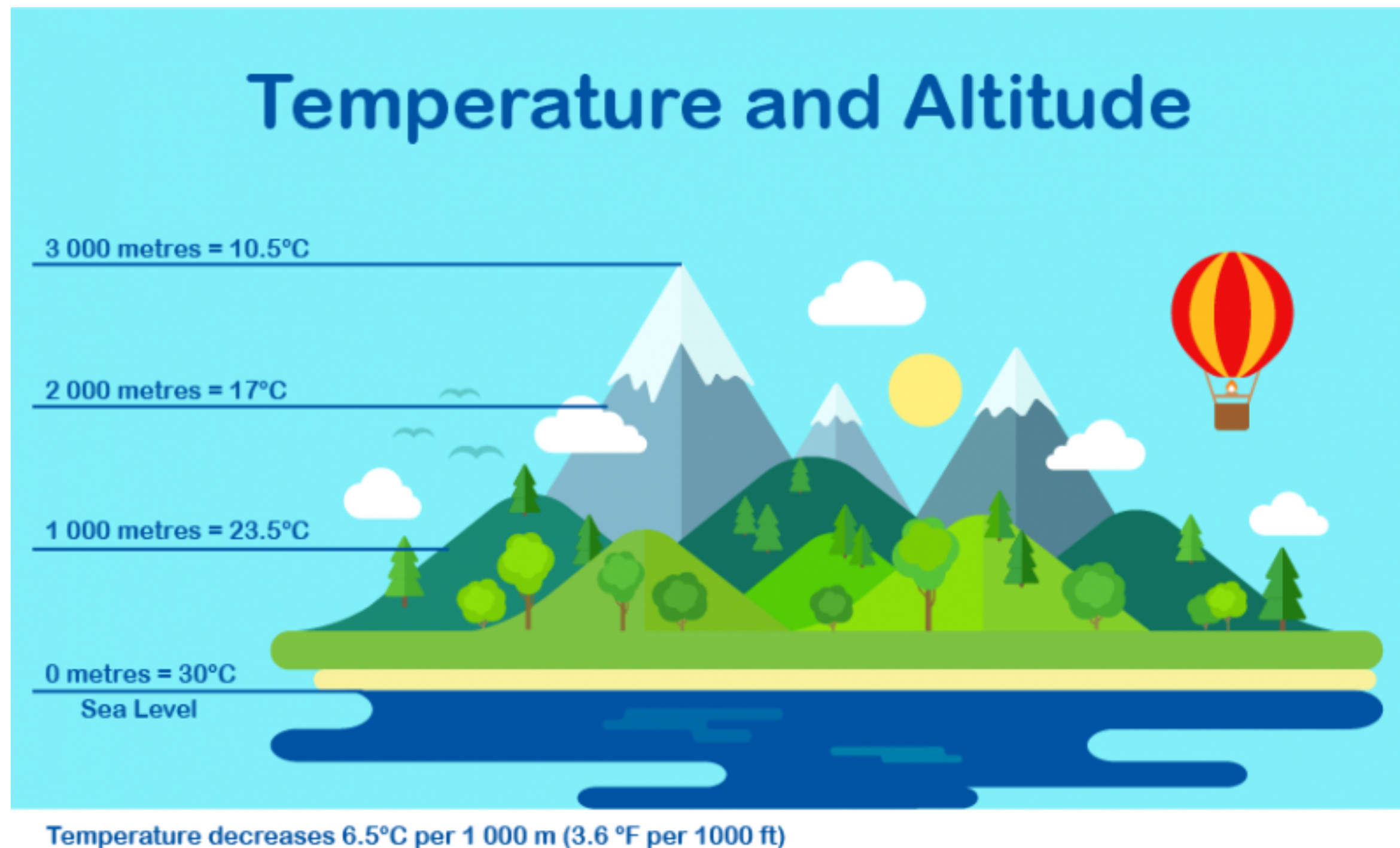
changing $P(A)$ does not change $P(T|A)$

$$P(A, T) = P(A|T)P(T)$$

Changing $P(T)$ might change $P(A|T)$



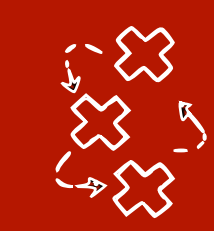
Modularity example



$$P(A, T) = P(T|A)P(A)$$

changing $P(A)$ does not
change $P(T|A)$

The **causal factorisation**
allows for **localised/sparse**
interventions



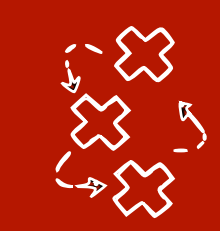
Structural **causal** models (SCMs)

- Let (G, p) be a **causal** Bayesian network
- We can write each variable X_i for $i \in \mathbf{V}$ as a **function of its parents** in G and an additional **noise term** ϵ_i in a **structural equation**:

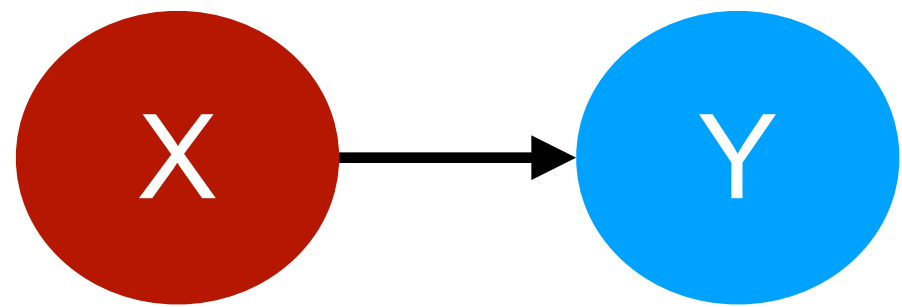
$$X_i \leftarrow h_i(X_{\text{Pa}(i)}, \epsilon_i)$$

often linear often Gaussian

- We often assume noises are **independent of each other** $\forall i \neq j : \epsilon_i \perp \epsilon_j$

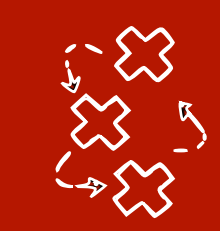


SCM: Example 3.2 in Elements of Causal Inference

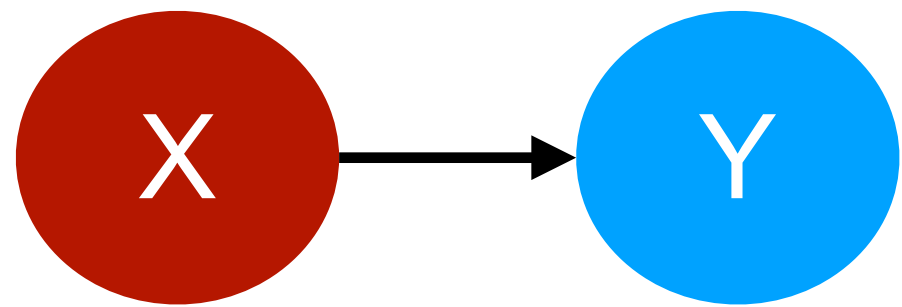


$$P(X) = \mathcal{N}(0,1)$$

$$P(Y|X = x) = \mathcal{N}(4 \cdot x, 1)$$



SCM: Example 3.2 in Elements of Causal Inference

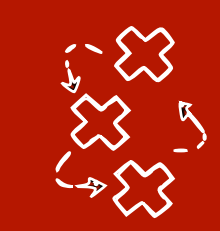


$$P(X) = \mathcal{N}(0,1)$$

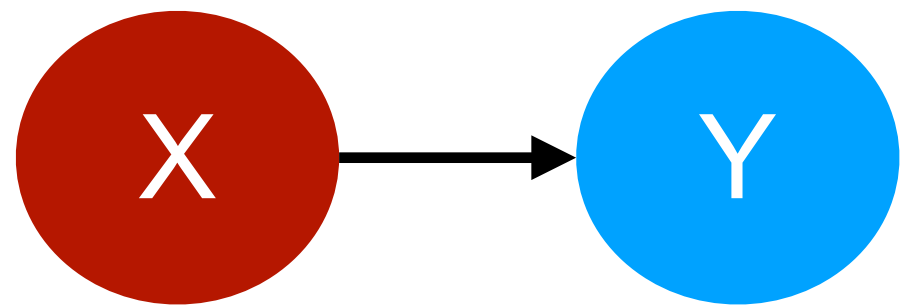
$$P(Y|X = x) = \mathcal{N}(4 \cdot x, 1)$$

$$\begin{cases} X \leftarrow \epsilon_X \\ Y \leftarrow 4 \cdot X + \epsilon_Y \end{cases}$$

$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$



SCM: Example 3.2 in Elements of Causal Inference



$$P(X) = \mathcal{N}(0,1)$$

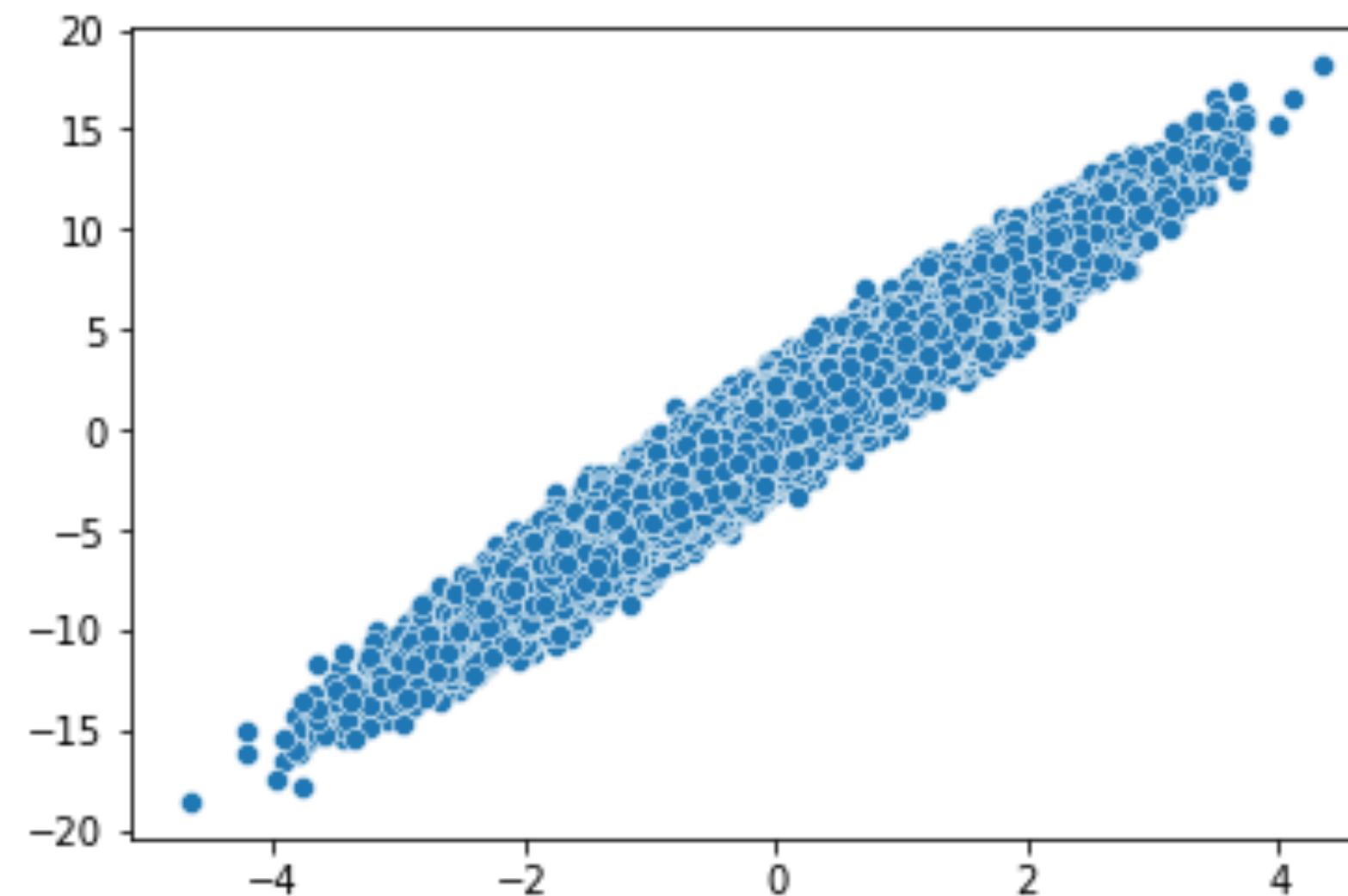
$$P(Y|X = x) = \mathcal{N}(4 \cdot x, 1)$$

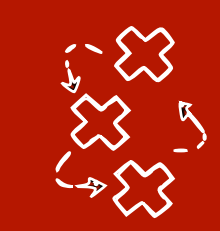
$$\begin{cases} X \leftarrow \epsilon_X \\ Y \leftarrow 4 \cdot X + \epsilon_Y \end{cases}$$

$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$

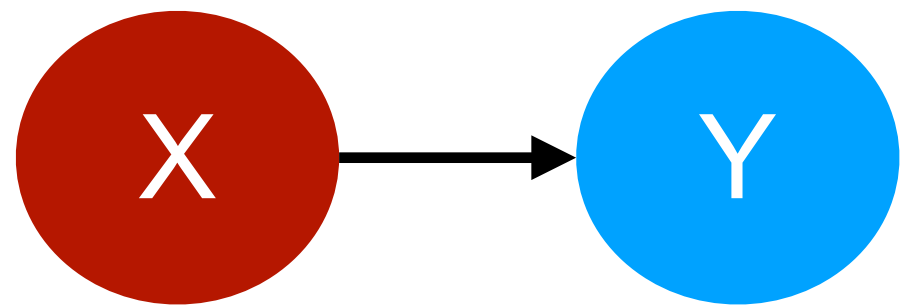
```
x = randn(n_samples)
y = 4 * x + randn(n_samples)
# plot P(X,Y)
sns.scatterplot(x=x,y=y)
```

<AxesSubplot:>





SCM: Example 3.2 in Elements of Causal Inference

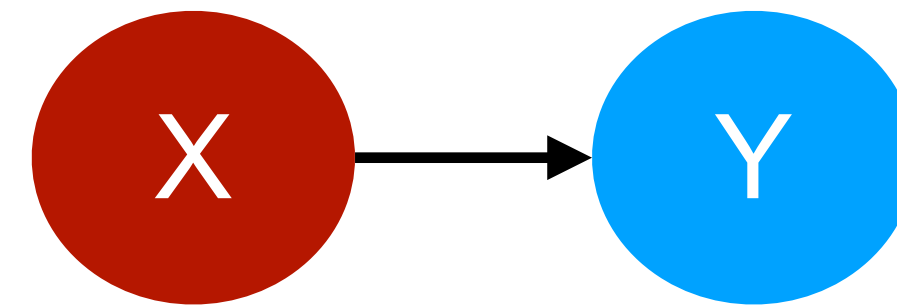


$$P(X) = \mathcal{N}(0,1)$$

$$P(Y|X = x) = \mathcal{N}(4 \cdot x, 1)$$

$$\begin{cases} X \leftarrow \epsilon_X \\ Y \leftarrow 4 \cdot X + \epsilon_Y \end{cases}$$

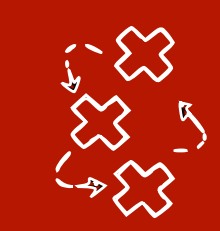
$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$



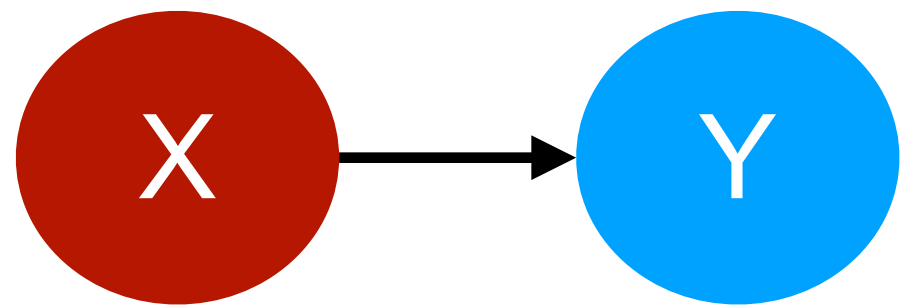
do($X = 2$) :

$$\begin{cases} X \leftarrow 2 \\ Y \leftarrow 4 \cdot X + \epsilon_Y \end{cases}$$

$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$



SCM: Example 3.2 in Elements of Causal Inference

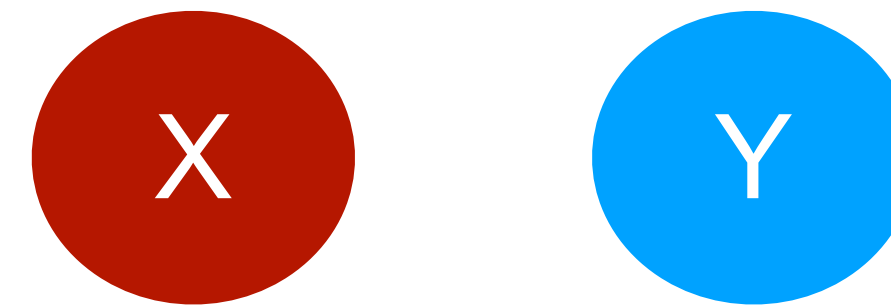


$$P(X) = \mathcal{N}(0,1)$$

$$P(Y|X = x) = \mathcal{N}(4 \cdot x, 1)$$

$$\begin{cases} X \leftarrow \epsilon_X \\ Y \leftarrow 4 \cdot X + \epsilon_Y \end{cases}$$

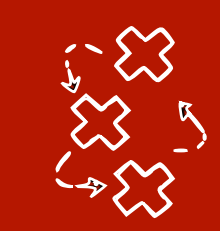
$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$



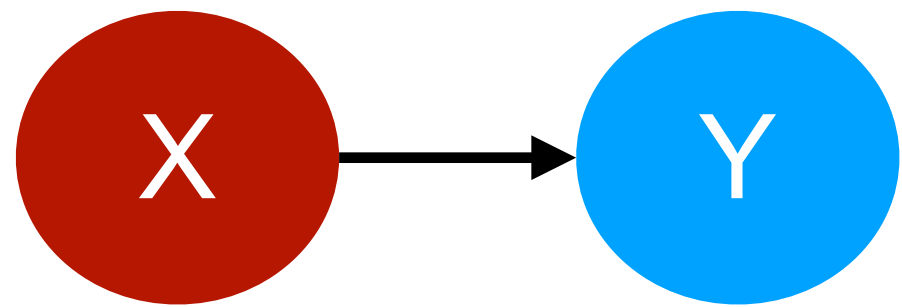
do($Y = 4$) :

$$\begin{cases} X \leftarrow \epsilon_x \\ Y \leftarrow 4 \end{cases}$$

$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$



Soft interventions example

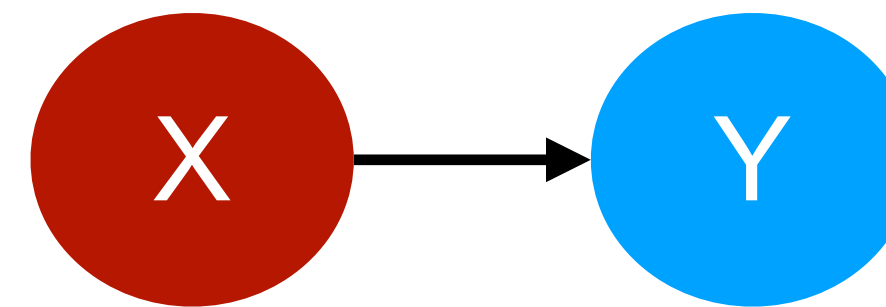


$$P(X) = \mathcal{N}(0,1)$$

$$P(Y|X = x) = \mathcal{N}(4 \cdot x, 1)$$

$$\begin{cases} X \leftarrow \epsilon_X \\ Y \leftarrow 4 \cdot X + \epsilon_Y \end{cases}$$

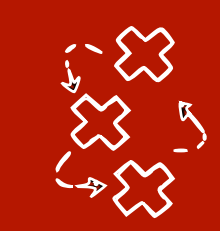
$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$



soft intervention on Y :

$$\begin{cases} X \leftarrow \epsilon_X \\ Y \leftarrow 3 \cdot X + 5 \cdot \epsilon_Y \end{cases}$$

$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$



Not in this module: Identification of causal effects

- Given a causal graph G , an **identification strategy** is a formula to estimate an interventional distribution from a combination of observational ones

- Backdoor criterion, Adjustment criterion**

$$p(x_j | \text{do}(x_i)) = \int_{x_Z} p(x_j | x_i, x_Z) p(x_Z) dx_Z$$

- Z does not contain any descendant of nodes $r \neq i$ on a directed path from i to j , AND vs Backdoor: $Z \cap \text{Des}(i) = \emptyset$
- Z blocks all paths from i to j that are not directed paths from i to j vs Backdoor: all backdoor paths

- Frontdoor criterion**

$$p(x_j | \text{do}(x_i')) = \int_{x_M} p(x_M | x_i') \int_{x_i} p(x_j | x_M, x_i') p(x_i) dx_i$$

- M blocks all directed paths from i to j , AND
- There are no unblocked backdoor paths from $i \leftarrow \dots$ to M , AND
- $Z = \{i\}$ blocks all backdoor paths from $M \leftarrow \dots$ to j

- Instrumental variables**

We can exploit the instrumental variable (IV) I $\beta = \frac{\text{Cov}(I, Y)}{\text{Cov}(I, X)}$

- Do-calculus (complete)**

- $I \rightarrow X$, but $I \not\rightarrow Y$ directly, $I \perp\!\!\!\perp W$