

GT. Lecture of
Introduction to Causal inference
from a Machine Learning Perspective.

Book by Brady Neal.

- ① Motivations.
- ② Potential Outcomes
- ③ The flow of Association & Causation in graphs
- ④ Causal Models
 - . do operator
 - . Structural Causal models
- ⑤ Randomized experiments.
- ⑥
- ⋮
- ⑪ Causal discovery from observational data
- ⑫ Causal ————— interventional data.

(non fini).

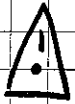
Chap 1 Motivation.

①

④ Simpson's paradox.
 Judge a population of 2 treatments

		Condition (C)		
		Mild	Severe	Total
Treatment (T)	A	$\boxed{15\%}$ 210/1400	$\boxed{30\%}$ 30/100	$\boxed{16\%}^*$ 240/1500
	B	$\boxed{10\%}^*$ 5/50	$\boxed{20\%}^*$ 100/500	$\boxed{19\%}$ 105/550

 % of deaths
 Y outcome (death)



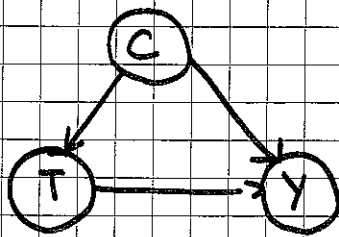
Treatment A more efficient for Mild Condition
 Severe Condition

Treatment B _____ for total population

What treatment should we use?

The answer comes from the causal structure we set on the data.
 (could come?)

Scenario 1



ex doctors keep treatment B (more expensive) for patients with severe conditions.

In this case, treatment B is associated to ↑ mortality because it is given to people more likely to die.

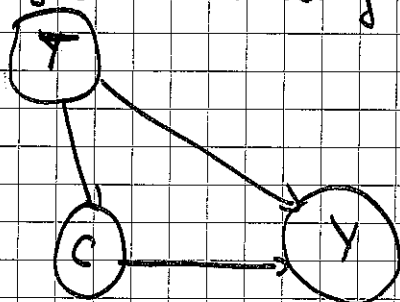
→ choose B because in each condition it works better

Scenario 2

Treatment is a cause of condition



For instance treatment B is rare (scarce) so people have to wait for it and it worsens their condition => higher mortality.



=> Prefer A.

Conclusion Getting a causal assumption will help to decide.

[=> We need to be able to get the causal graph to solve the problem.

② How to capture the causal relation?

⚠ correlation is not causality
association

2 curves. purely coincidence. spurious correlation

see site

(although AI will find an explanation).

Why correlation replaced by association.

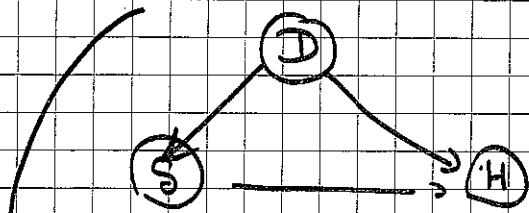
(5)

because correlation = linear.

Here association = statistical dependence

Ex people going to bed with shoes are more likely to have headaches.

⚠ Common cause: drinking the night before.



↳ confounding association.

If association = causality \Rightarrow standard statistics would be sufficient.

descript^o of main themes \rightarrow sublimage.

Consider the situation / experiment where you get a treatment (T) and observe an outcome (Y) after treatment.

ex

T	Y
medic / no medic	dead or not
dog / no dog	being more happy or not

• If getting a dog makes you happy how can you be sure that you are happy BECAUSE of the DOG.

• If you get the medicine you are better. But what would have happened without the treatment? Maybe you would be ok too?

We can not go back in time!!!

We define 4 variables.

	T	<u>treatment</u>
	Y	<u>outcome of the experiment</u>
$Y(x)$	$Y(1)$	<u>potential outcome with the treatment 1</u>
	$Y(0)$	<u>potential outcome with the treatment 0</u>

Assumption 1 consistency

$$Y = Y(T)$$

if treatment T is given the the observed outcome Y is the potential outcome T.

Individual treatment effect.

$$\tau = Y(1) - Y(0)$$

△ if population $i = 1 \dots n$ Y_i, X_i, T_i
↑
covariates.

$$\tau_i = Y_i(1) - Y_i(0)$$

Pb observations

fundamental problem of causal inference

we only observe $Y_i(1)$ or $Y_i(0)$ never both.

i	T	Y		Y(1) - Y(0)
		Y(1)	Y(0)	
1	0	?	0	?
2	1	1	?	?
3	1	0	?	?
4	0	?	1	?

The potential outcome that we do not observe are called counterfactuals

The counterfactuals only exists once the outcome has been observed.

2.3 Getting around the fundamental problem.

(5)

We can't access the individual T.E.

Can we get the average treatment effect over a population?

$$ATE = \tau = E[Y(1) - Y(0)] \quad \text{where the average is over the individuals.}$$

$$= E[Y(1)] - E[Y(0)] \quad (\text{linearity of expectation}).$$

Is it equal to $E[Y|T=1] - E[Y|T=0]$?

• No in general.

• Yes if

$$\boxed{(Y(1), Y(0)) \perp\!\!\!\perp T} \quad (A.1)$$

Ignorability
Exchangeability

Indeed If this is true. \Rightarrow

$$\begin{aligned} ATE &= E[Y(1)] - E[Y(0)] = E[Y(1)|T=1] - E[Y(0)|T=0] \\ &= E[Y|T=1] - E[Y|T=0] \end{aligned}$$

\Rightarrow on peut estimer ATE en faisant des moyennes par colonnes ds le tableau.
Stat classiques.

Ignorability \Leftrightarrow T independent des potential outcome. ⑤
or T lié au fait que $Y(1)$ ou $Y(0)$ est non observé.

Le fait que i est manquant ne dépend pas des outcome \Rightarrow on peut faire des moyennes en "ignorant" le pb des données manquantes.

Exchangeability

$$E[Y(1) | T=0] = E[Y(1) | T=1]$$

ie. les groupes sont échangeables.
Traitement.

Group treatments are comparable.

the same in all relevant aspects other than treatment.

linked to the concept of "controlling for"

We used the assumption A1 to identify causal effect

def 1 A causal quantity (ex $E[Y(t)]$) is identifiable if we can compute it from a purely statistical quantity of $E[Y | X=t]$.

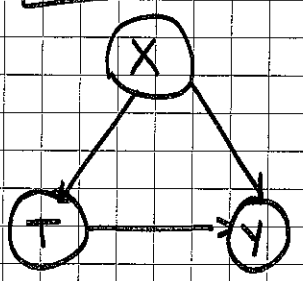
But A1 is unrealistic : surely there is confounding in most data.
ignorability

solution : randomized experiences. (chap 5).

- propose a less unrealistic assumption.
- conditional exchangeability.

Assumption 2

Covariates X which describes the individuals.
 $(Y(1), Y(0)) \perp\!\!\!\perp T \mid X$ Conditional Exchangeability or unconfoundedness.



⇒ we are able to identify the causal effects within levels of X

$$\begin{aligned}
 \bullet \ E [Y(1) - Y(0) \mid X] &= E [Y(1) \mid X] - E [Y(0) \mid X] \\
 \text{(CA2)} &= E [Y(1) \mid X, T=1] - E [Y(0) \mid X, T=0] \\
 &= E [Y \mid X, T=1] = E [Y \mid X, T=0]
 \end{aligned}$$

Marginal effect.

$$\begin{aligned}
 E [Y(1) - Y(0)] &= E_x [Y(1) - Y(0) \mid X] \\
 &= E_x [E [Y \mid X, T=1]] - E_x [E [Y \mid X, T=0]]
 \end{aligned}$$

$$E [Y(1) - Y(0)] = E_x [E [Y \mid X, T=1] - E [Y \mid X, T=0]] \quad \checkmark \Delta \text{ positivity}$$

demo more detailed after (*)

⚠ exchangeability \Rightarrow conditional exchangeability to be more realistic. ⑦ bis

But it may exist unobserved confounders that are not in X.

\Rightarrow observational data (Chap 3 & 4)

\Rightarrow randomized trials - not a pb.

We are tempted to add as many covariates as possible.

BUT in fact we need an assumption of positivity.

$$\|A3\| \quad \forall x \quad 0 < \mathbb{P}(T=1 | X=x) < 1$$

$\mathbb{P}(X=x) > 0.$

Positivity

Overlap.

Common support.

Why is it important

$$\mathbb{E}_x \left[\mathbb{E}[Y | T=1, X] - \mathbb{E}[Y | T=0, X] \right]$$

$$= \sum_x \mathbb{P}(X=x) \left[\sum_y y \mathbb{P}(Y=y | T=1, X=x) - \sum_y y \mathbb{P}(Y=y | T=0, X=x) \right]$$

$$= \sum_x \mathbb{P}(X=x) \left[\sum_y y \frac{\mathbb{P}(Y=y, T=1, X=x)}{\mathbb{P}(T=1 | X=x) \mathbb{P}(X=x)} - \sum_y y \frac{\mathbb{P}(Y=y | T=0, X=x)}{\mathbb{P}(T=0 | X=x) \mathbb{P}(X=x)} \right]$$

$\Delta = 0$ ⚠

intuition means you do observe some ~~observed~~ covariates with treatment or not.

why "same support" $\text{Supp}(\mathbb{P}(X | T=0)) = \text{Supp}(\mathbb{P}(X | T=1))$

Positivity-unconfoundedness tradeoff of LATE

Other assumptions

Ass 4 No interference
 $Y_i(t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_m) = Y_i(t_i)$

independence !!

Note 2.4 => terms of statistics

- estimand = quantity that we want to estimate
- Statistical estimand: does not imply the potential outcome.

$$\underbrace{E[Y(1) - Y(0)]}_{\text{causal estimand}} = \underbrace{E_x [E(Y | T=1, X) - E(Y | T=0, X)]}_{\text{statistical estimand}}$$

causal estimand \rightsquigarrow statistical estimand
identification:

3:30 2:30

9

6:00 x 5

How to estimate.

$E_x [E(Y|T=1, X)] \rightsquigarrow$ model - assisted estimators

Empirical over data. $\frac{1}{n} \sum_{i=1}^n \underbrace{E(Y|T=1, X=x_i)}_{\text{linear model}}$.

Continuous treatment

We care about $E[Y(x)]$: how it changes with x
 $= E[Y | T=x]$

if linear
 $\alpha x + \beta$ $\frac{\partial}{\partial x} E[Y(x)] = \alpha$

\Rightarrow pb if not linear ~~of~~

Chap 3 The flow of Associations & Causation in Graphs.

• Graphs have seemed to be convenient to represent and intuition

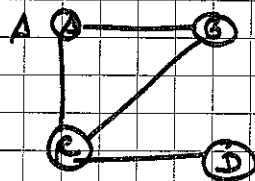
⇒ utility of graphical models.

3.1 Graph terminology

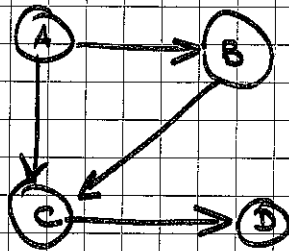
Definitions

• Graph = collection of edges & nodes.

• undirected



• directed



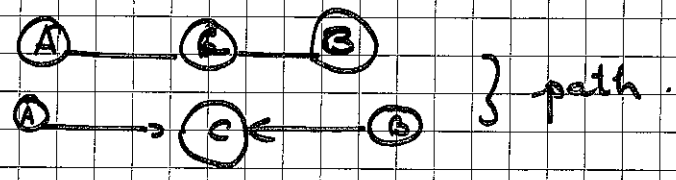
• No directed edge goes from parent node to child node.

• $\begin{cases} pa(X) = \text{ensemble of the parents of } X. \\ pa_i = pa(X_i) \end{cases}$

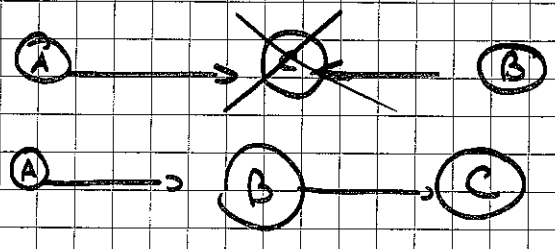
• 2 Nodes are adjacent if there \exists an edge between them.

B & D non adjacent
A & B adjacent

Path = any sequence of adjacent nodes regardless the direction of edges.



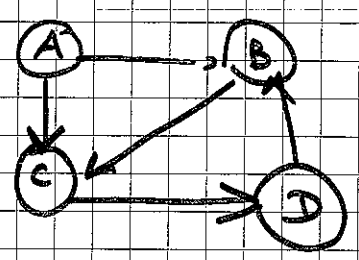
Directed path is a path with directed edges all in the same direction



If a directed path starts in X and ends in Y then

$$\begin{cases} X = \text{ancestor} \\ Y = \text{descendant} \end{cases}$$

Cycle if a path starts and ends at the same node.

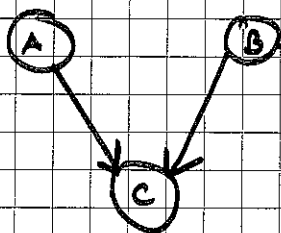


B - C - D - B is a cycle

If there is no cycles in a directed graph,
graph = DAG
directed acyclic graph.

13

Immortality: two parents and no edges between the parents



2. Bayesian Networks.

prob graphical models: broad class of models that uses graphs to represent the relationships between dependence random variables.

→ directed graphs: ex Bayesian networks
non-directed graph: Markov field

Bayesian network the graph is a DAG

(is a prob graphical model with particular graph).

$$\text{Factorization: } p(X_1, \dots, X_m) = \prod_{i=1}^m p(X_i | pa(X_i))$$

$$15 + 26 + 1 = 42 + 45 \quad \boxed{14}$$

- Let $P(X_1, \dots, X_n)$ be a joint distribution.
- and G a DAG with n nodes.
- with node i representing X_i .

Def P is Markov with respect to G if
(locally)

any random variable X_i is independent of
all its non-descendants in G given its parents

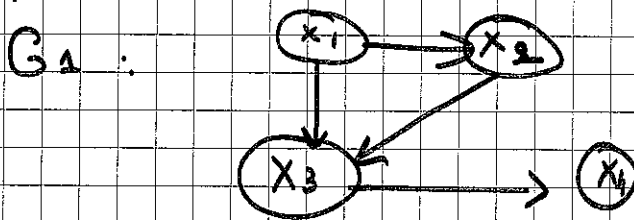
Consequence

Take 4 nodes

$$P(x_1, x_2, x_3, x_4) = P(x_1) P(x_2 | x_1) P(x_3 | x_2, x_1) P(x_4 | x_3, x_2, x_1)$$

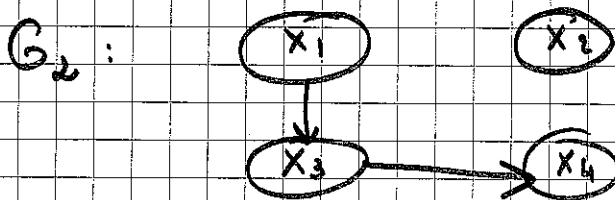
Always true
(Chain rule
of probability)

- If we know that P is Markov with respect to



$$\Rightarrow P(x_4 | x_3, \overbrace{x_2, x_1}^{\text{non descendants}}) = P(x_4 | x_3)$$

parent



$$P(x_2 | x_1) = P(x_2)$$

$$P(x_3 | x_1, x_2) = P(x_3 | x_1)$$

$$P(x_4 | x_3, x_2, x_1) = P(x_4 | x_3)$$

Consequence

If P is locally Markov with respect to G .

Then

$$P(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | PA(x_i)).$$

Bayesian Network factorisation.

Chain rule for Bayesian network.

Markov compatibility.

⚠ j'en pense que c'est pas direct comme DENO a que la DAG y joue un rôle. A voir avec un vrai livre de graphes.

in fact it is an equivalence.

Remarque The local Markov assumption only gives us information on the independences

It does not tell us nothing about dependence between adjacent nodes in the DAG.

Pour assurer cela il faut en dire un peu plus.

It's Minimality Assumption

B is locally Markov with respect to G & G is minimal if

- Given its parents a node is indep of its non descendant.
- Adjacent nodes in DAG are dependant.

if P is minimal there is no additional independencies

\Rightarrow if a ~~DAG~~ $(x) \longrightarrow (y)$.

if P is Markov with respect to P

$$\Rightarrow P(x, y) = P(y|x) P(x).$$

But maybe also $P(x, y) = P(y|p(x))$

if B is minimal

$$\Rightarrow P(x, y) \neq P(x) P(y).$$

We can't remove any other edges. if G is minimal

\Rightarrow Any edge is active

3. Causal Graphs.

Def - What is a cause

A variable X is said to be a cause of Y if Y can change in response to changes in X .

Causal Edge assumption: