

Diagramas Causais

Definições, Construção e Critério Backdoor

José Luiz Padilha – DEST/UFPR

25 de abril, 2025

Conteúdo

- 1 Preliminares
- 2 DAGs: algumas definições
- 3 Construindo um DAG
- 4 O critério backdoor
- 5 Um pouco de teoria
- 6 Bibliografia

Preliminares

Introdução: associação, causalidade e grafos

Em seu livro de **Causality: Models, Reasoning, and Inference (2009)**, Judea Pearl apresenta uma poderosa e extensa teoria gráfica de causalidade.

O trabalho de Pearl fornece um *framework* para causalidade que difere da abordagem de respostas potenciais. Contudo, o autor demonstra que os conceitos fundamentais subjacentes à perspectiva de resposta potencial e à perspectiva de grafos causais são equivalentes.

Os **grafos acíclicos dirigidos (DAGs)** são usados para desenvolver justificativas para métodos de estimação de efeitos causais.

Discutiremos como os DAGs permitem uma representação gráfica de relações causais e seu uso para identificação dos efeitos causais.

Como duas variáveis podem estar associadas na população?



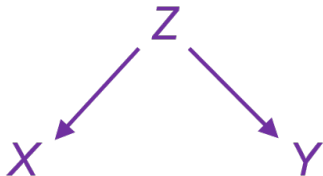
- Duas variáveis X e Y serão **associadas** na população se X causa Y .

Como duas variáveis podem estar associadas na população?



- X e Y serão associadas na população se Y causa X .

Como duas variáveis podem estar associadas na população?



- Por fim, X e Y serão associadas na população se existir alguma variável Z que causa **ambas** X e Y .

Como duas variáveis podem estar associadas na população?



- X e Y não podem ser associadas na população por qualquer outra razão.
- Se X e Y são associadas na população, então **pelo menos uma** das situações acima deve ser verdade.

O que queremos dizer com associação “na população”?

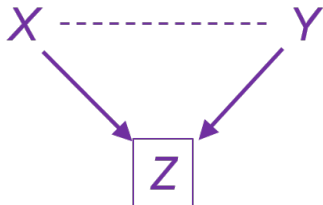
- Na terminologia estatística, X e Y são associadas “na população” significa que estas variáveis são **marginalmente associadas**.
- Se X e Y são marginalmente associadas, então, para um indivíduo em particular, saber a respeito de X nos dá alguma informação sobre o valor provável de Y , e vice-versa.
- Suponha, por simplicidade, X e Y dicotômicas. Se X e Y são marginalmente associadas, então

$$\Pr(Y = 1|X = 1) \neq \Pr(Y = 1|X = 0)$$

e

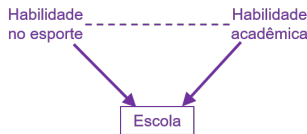
$$\Pr(X = 1|Y = 1) \neq \Pr(X = 1|Y = 0).$$

Como duas variáveis podem estar associadas em uma subpopulação?



- Suponha que Z é um **efeito** tanto de X como de Y.
- Então, X e Y serão **associadas dentro do estrato de Z**, mesmo se na população estas variáveis forem independentes.
- X e Y serão condicionalmente associadas (dado Z), mesmo que sejam marginalmente independentes (não associadas).
- A caixa ao redor de Z denota que estamos estratificando (condicionando) em Z.
- A reta tracejada denota a associação condicional induzida.

Como duas variáveis podem estar associadas em uma subpopulação?



- Suponha que uma escola aceita alunos ou porque são “bons” no **esporte**, ou porque são “bons” **academicamente**; ou ainda, alunos que são “bons” nos dois.
- Suponha que a habilidade acadêmica e a habilidade no esporte sejam **independentes** na população.
- **Dentro da escola**, existirá uma associação (negativa) entre habilidade acadêmica e habilidade no esporte.
- Por quê? Suponha que escolhemos um aluno ao acaso e percebemos que ele não tem habilidade nos esportes. Então, este deve ser “bom” academicamente.

Exemplo: Dados Simulados

```
library(tidyverse)
n <- 10000
X <- rnorm(n) # Habilidade no esporte
Y <- rnorm(n) # Habilidade acadêmica
cor(X, Y) # X e Y não são associadas na população
```

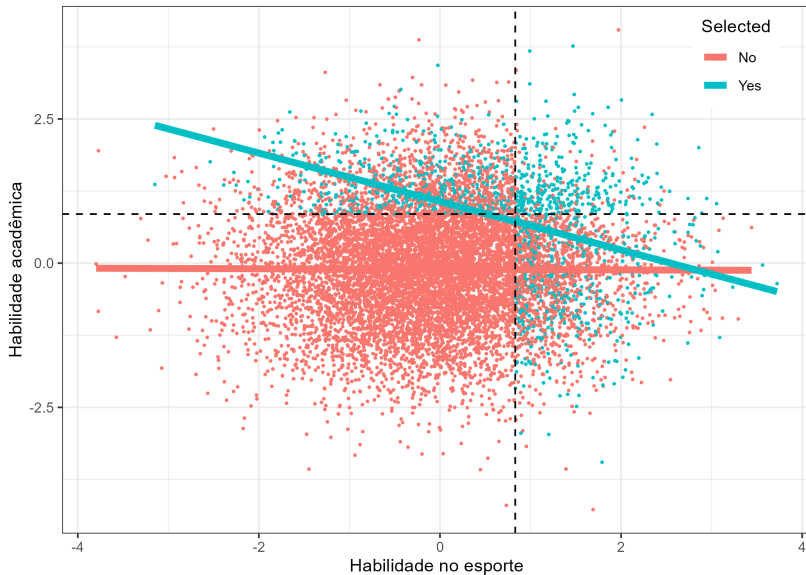
```
[1] 0.02158104
```

```
cor.test(X, Y)$p.value
```

```
[1] 0.03092122
```

```
probs <- 0 + 0.3*I(X>quantile(X, probs = 0.8)) +
  0.3*I(Y>quantile(Y, probs = 0.8)) +
  0.2*I(X>quantile(X, probs = 0.8) & Y>quantile(Y, probs = 0.8))
Z <- rbinom(n, 1, probs)
dat <- data.frame(X, Y, Z = factor(Z, levels = c(0, 1),
                                   labels = c("No", "Yes")))
```

Exemplo: Dados Simulados



Exemplo: Dados Simulados

```
cor.test(X[Z==1],Y[Z==1])
```

Pearson's product-moment correlation

data: X[Z == 1] and Y[Z == 1]

t = -15.915, df = 1233, p-value < 2.2e-16

alternative hypothesis: true correlation is not equal to 0

95 percent confidence interval:

-0.4580448 -0.3654441

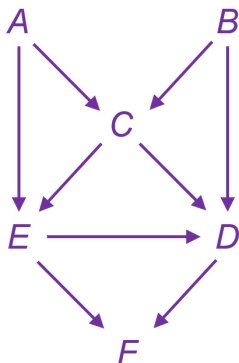
sample estimates:

cor

-0.4128106

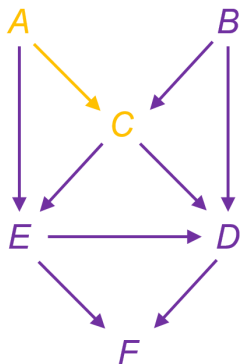
DAGs: algumas definições

Um exemplo



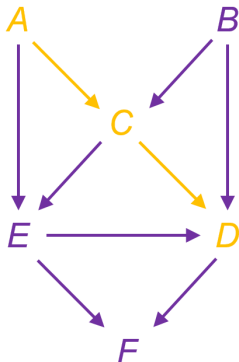
Grafo acíclico dirigido

- Este é um exemplo de um **grafo acíclico dirigido** (DAG) causal (diagrama causal).
- É **dirigido**, pois cada aresta é uma seta de ponta única.
- É **causal**, pois as setas representam nossas suposições a respeito da direção da influência causal.
- É **acíclico**, pois não contém ciclos: nenhuma variável causa a si mesma.



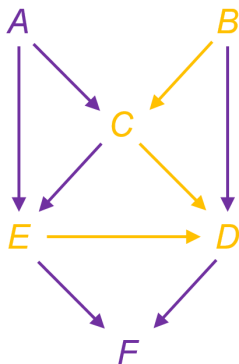
Pais e filhos

- A é pai (ou mãe) de C.
- C é filho (ou filha) de A.



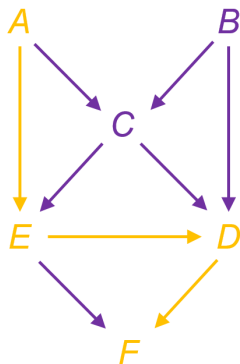
Ancestrais e descendentes

- A é um **ancestral** de D .
- D é **descendente** de A .
- A também é um **ancestral** de C .
- C também é um **descendente** de A .
 - Ou seja, pais são ancestrais, e filhos são descendentes.



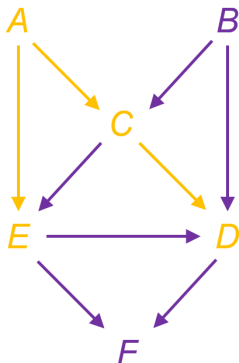
Caminho

- Este é um **caminho** de E para B .



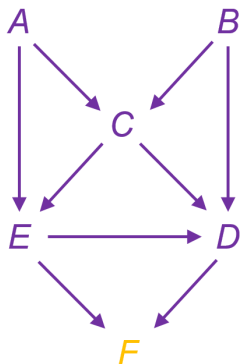
Caminho dirigido

- Este é um **caminho dirigido** de A para F (todas as setas no caminho apontam “para frente”).



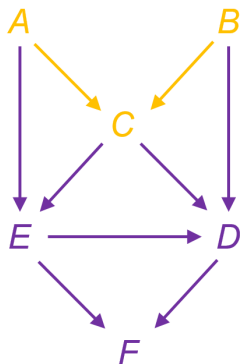
Caminho backdoor

- Este é um **caminho porta dos fundos** de E para D (o caminho começa com uma seta chegando em E).



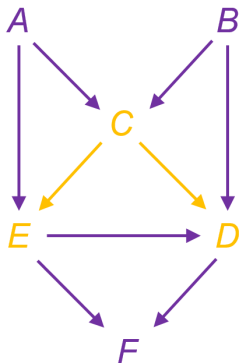
Collider

- F é um **colisor** desde que duas pontas de setas se encontram em F .



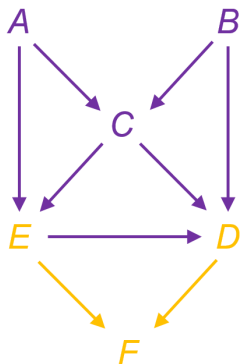
Nota

- Note que C é um colisor no caminho $A \rightarrow C \leftarrow B$.



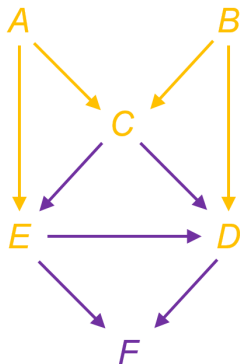
Nota

- No entanto, C NÃO É um colisor no caminho $E \leftarrow C \rightarrow D$.
- Assim, a definição de um colisor é em relação ao caminho que está sendo considerado.



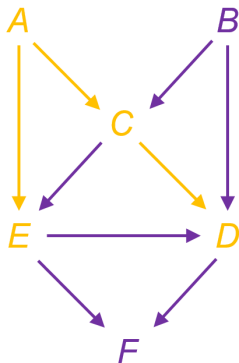
Caminho bloqueado

- O caminho $E \rightarrow F \leftarrow D$ é **bloqueado** desde que este contenha um colisor (F).



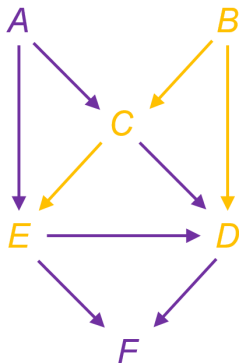
Caminho bloqueado

- Este caminho também é bloqueado (em C).



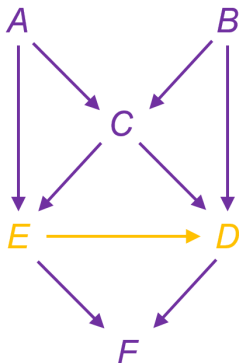
Caminho aberto

- Um caminho que não contém um colisor está **aberto**. Aqui temos um exemplo.



Caminho aberto

- E outro.



Caminho aberto

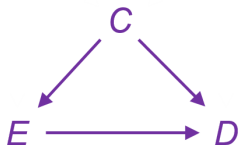
- E outro.

Construindo um DAG



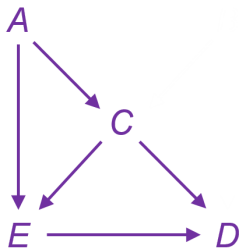
Passo 1

- O primeiro passo na construção de um DAG para um problema particular é escrever a **exposição** e o **desfecho** de interesse, com uma seta da exposição para o desfecho.
- Esta seta representa o **efeito causal** que queremos estimar.



Passo 2

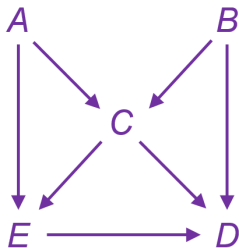
- Se existir qualquer **causa comum** C de E e D , devemos colocá-lo no grafo, com setas de C para E e de C para D .
- Devemos incluir C no grafo, independentemente deste ter sido ou não mensurado em nosso estudo.



Passo 3

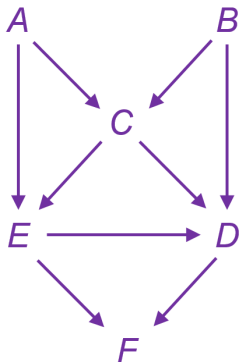
- Continuamos assim, adicionando ao diagrama qualquer variável (observada ou não observada) que é uma **causa comum** de duas ou mais variáveis já existentes no diagrama.

Construindo um DAG



Passo 3

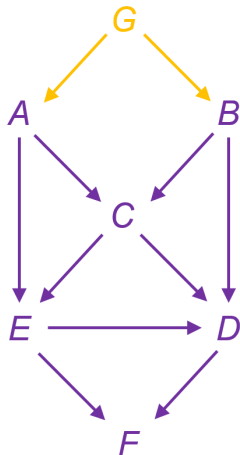
- Continuamos assim, adicionando ao diagrama qualquer variável (observada ou não observada) que é uma **causa comum** de duas ou mais variáveis já existentes no diagrama.



Passo 3

- Se quisermos, podemos também incluir **outras variáveis**, mesmo que eles não sejam causas comuns de outras variáveis no diagrama.
- Por exemplo, F .
- Vamos supor que finalizamos nesse ponto. As variáveis e setas que **NÃO** estão em nosso grafo representam nossas **suposições causais**.

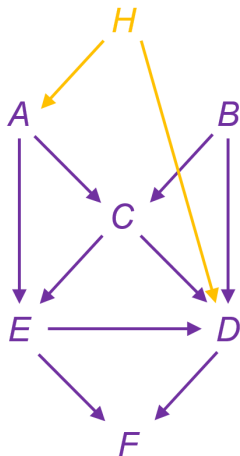
Construindo um DAG



Quais são as nossas suposições?

- Por exemplo, estamos fazendo a suposição que não há uma causa comum G de A e B .

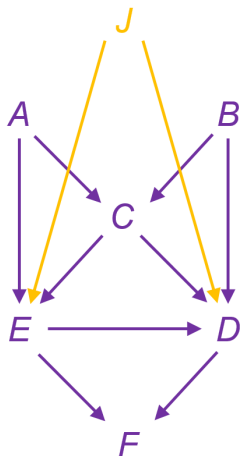
Construindo um DAG



Quais são as nossas suposições?

- E que não há uma causa comum H de A e D .

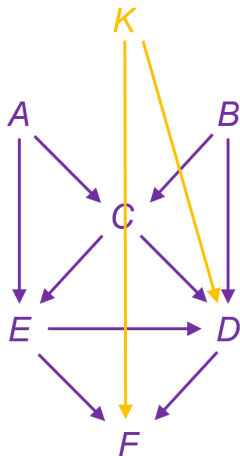
Construindo um DAG



Quais são as nossas suposições?

- E que A , B e C representam TODAS as causas comuns de E e D ; não há uma causa comum adicional J .

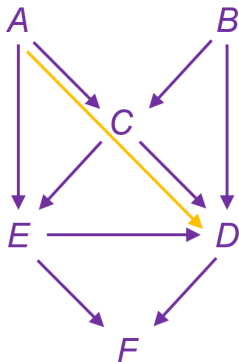
Construindo um DAG



Quais são as nossas suposições?

- E que não há uma causa comum adicional K de D e F .

Construindo um DAG

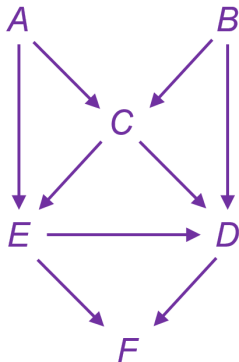


Quais são as nossas suposições?

- Portanto, cada seta omitida também representa uma suposição.
- Por exemplo, estamos assumindo que todo o efeito de *A* em *D* atua por meio de *C* e *E*.

O critério backdoor

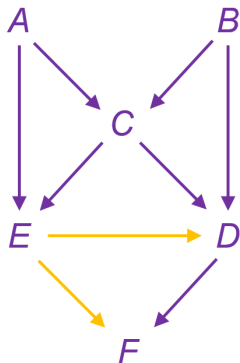
O critério backdoor: existe confundimento?



Qual o próximo passo?

- SE acreditarmos em nosso diagrama causal, podemos proceder para determinar se a relação $E \rightarrow D$ está **confundida** ou não.
- Isto é feito utilizando o **critério porta dos fundos**.
- O critério porta dos fundos é aplicado em duas partes:
 - 1 a primeira parte define se existe ou não confundimento.
 - 2 se existir, a segunda parte determina se é possível controlar o confundimento.

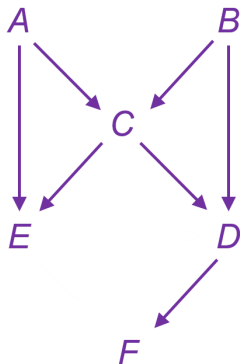
O critério backdoor: existe confundimento?



Passo 1

- Primeiro removemos todas as setas **saindo da exposição**.

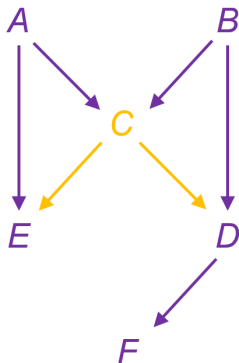
O critério backdoor: existe confundimento?



Passo 2

- Em seguida, procuramos por caminhos abertos a partir da exposição até o desfecho.
- Relembrando: um caminho aberto não contém um colisor.

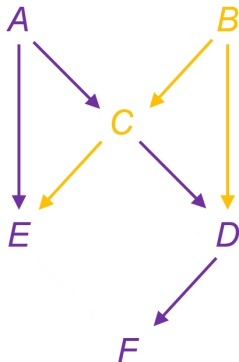
O critério backdoor: existe confundimento?



Passo 2

- Este é um caminho aberto?

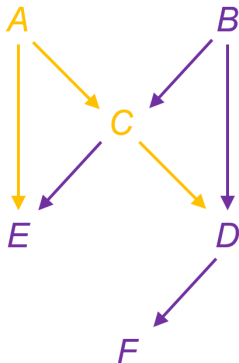
O critério backdoor: existe confundimento?



Passo 2

- Este é um caminho aberto?

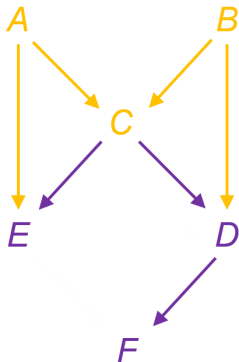
O critério backdoor: existe confundimento?



Passo 2

- Este é um caminho aberto?

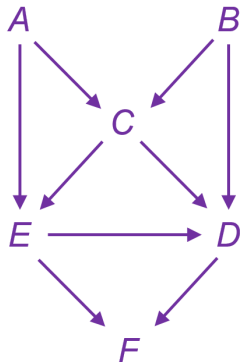
O critério backdoor: existe confundimento?



Passo 2

- Este é um caminho aberto?

O critério backdoor: existe confundimento?



Existe confundimento?

- Identificamos três caminhos porta dos fundos abertos de E para D . Assim, há confundimento.
- Próxima pergunta: podemos usar alguns ou todos de A , B , C , F para controlar esse confundimento?
- Existe um conjunto \mathcal{S} de variáveis tal que se estratificarmos (ajustarmos) por elas, podemos concluir que o efeito causal existe no estrato?

O critério backdoor

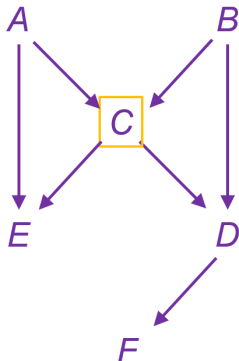
A segunda parte do critério da porta dos fundos nos permite determinar, com base em nosso diagrama causal, se um conjunto de covariáveis candidato é ou não suficiente para controlar o confundimento:

O critério backdoor

- (i) Primeiro, o conjunto candidato \mathcal{S} não deve conter descendentes da exposição.
- (ii) Em seguida, removemos todas as setas que saem da exposição.
- (iii) Então, nós juntamos com uma linha tracejada quaisquer duas variáveis que compartilham um filho que esteja ela mesma em \mathcal{S} ou que tenha um descendente em \mathcal{S} .
- (iv) Existe um caminho aberto de E para D que não passa por um membro de \mathcal{S} ?

Se NÃO, então \mathcal{S} é suficiente para controlar para confundimento.

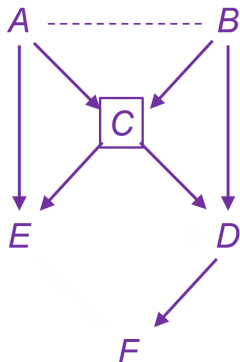
O critério backdoor: podemos controlar o confundimento?



O critério backdoor: passos (i) e (ii)

- C é suficiente?
- C não é um descendente de E , então o passo (i) é satisfeito.
- Todas as setas saindo da exposição já foram removidas (passo (ii)).

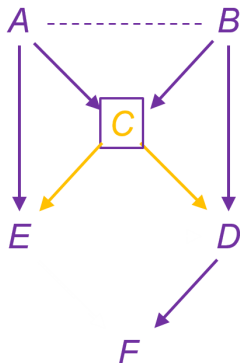
O critério backdoor: podemos controlar o confundimento?



passo (iii)

- Conectamos A e B com uma linha tracejada, pois eles compartilham um filho (C) que está em nosso conjunto candidato (C).
- Nenhuma outra variável precisa ser conectada desta forma.

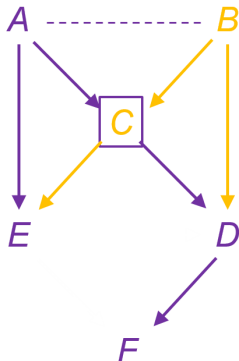
O critério backdoor: podemos controlar o confundimento?



passo (iv)

- Agora procuramos por caminhos abertos de E para D e vemos se estes todos passam por C .
- Este está OK!

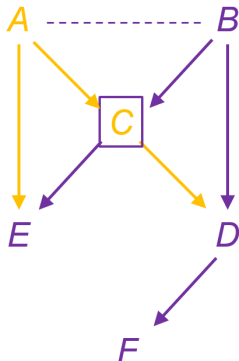
O critério backdoor: podemos controlar o confundimento?



passo (iv)

- Este também!

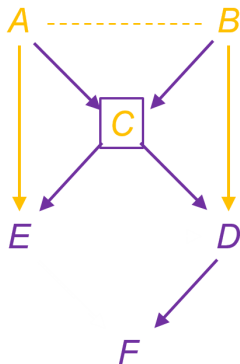
O critério backdoor: podemos controlar o confundimento?



passo (iv)

- Este também!

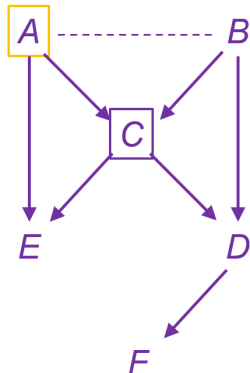
O critério backdoor: podemos controlar o confundimento?



passo (iv)

- PORÉM, aqui está um caminho aberto de E para D que NÃO passa por C
- Assim, controlar apenas por C NÃO é suficiente.

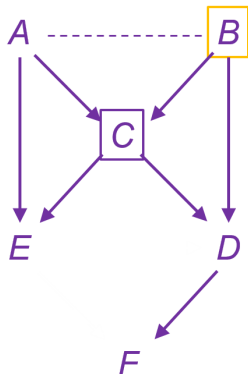
O critério backdoor: podemos controlar o confundimento?



Qual é a solução?

- Devemos controlar adicionalmente para A.

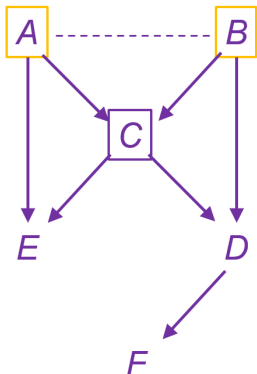
O critério backdoor: podemos controlar o confundimento?



Qual é a solução?

- Ou B .

O critério backdoor: podemos controlar o confundimento?



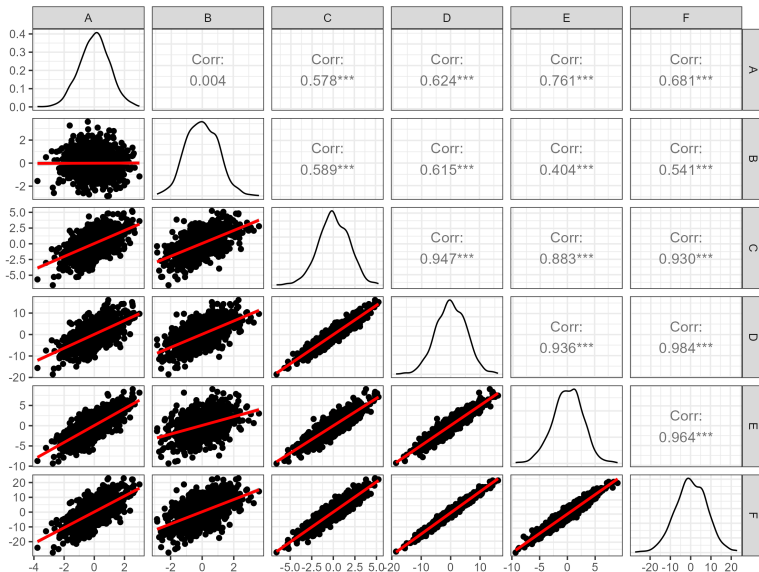
Qual é a solução?

- Ou ambos A e B para controlar para o confundimento.

Exemplo: Dados Simulados

```
set.seed(123456)
n <- 1000
A <- rnorm(n, mean = 0)
B <- rnorm(n, mean = 0)
C <- rnorm(n, mean = A + B)
E <- rnorm(n, mean = A + C)
D <- rnorm(n, mean = E + C + B)
F <- rnorm(n, mean = E + D)
```

Exemplo: Dados Simulados



Exemplo: Dados Simulados

```
# Sem controlar por confundidoras (modelo errado)  
# Beta real para a exposição E é um  
printCoefmat(coef(summary(f1 <- lm(D ~ E))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.8945e-05	5.7078e-02	0.0009	0.9993
E	1.7584e+00	2.0979e-02	83.8190	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exemplo: Dados Simulados

```
# Controlando apenas por C (modelo errado)  
printCoefmat(coef(summary(f2 <- lm(D ~ E + C))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.024917	0.039264	-0.6346	0.5258
E	0.852752	0.030747	27.7349	<2e-16 ***
C	1.572563	0.047142	33.3581	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exemplo: Dados Simulados

```
# Controlando por todas as variáveis (modelo errado)  
printCoefmat(coef(summary(f3 <- lm(D ~ E + A + B + C + F))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.012905	0.022092	-0.5841	0.5593
E	0.056234	0.037664	1.4930	0.1357
A	-0.038847	0.039956	-0.9723	0.3312
B	0.487645	0.034200	14.2587	<2e-16 ***
C	0.491745	0.035293	13.9331	<2e-16 ***
F	0.491357	0.015108	32.5220	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exemplo: Dados Simulados

```
# Controlando por todas as variáveis pré-tratamento  
# (desnecessário)  
printCoefmat(coef(summary(f4 <- lm(D ~ E + A + B + C))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.0012754	0.0317176	0.0402	0.9679
E	1.0295130	0.0328393	31.3500	<2e-16 ***
A	0.0236395	0.0573091	0.4125	0.6801
B	0.9439513	0.0447865	21.0767	<2e-16 ***
C	1.0065989	0.0452955	22.2230	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exemplo: Dados Simulados

```
# Controlando pelas variáveis do critério backdoor  
printCoefmat(coef(summary(f5 <- lm(D ~ E + A + C))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.026824	0.038094	-0.7042	0.4815
E	1.057980	0.039442	26.8235	< 2.2e-16 ***
A	-0.494846	0.062221	-7.9531	4.918e-15 ***
C	1.454513	0.048084	30.2494	< 2.2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exemplo: Dados Simulados

```
# Controlando pelas variáveis do critério backdoor  
printCoefmat(coef(summary(f6 <- lm(D ~ E + B + C))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.00096507	0.03169550	0.0304	0.9757
E	1.03774975	0.02605995	39.8216	<2e-16 ***
B	0.93602139	0.04043367	23.1496	<2e-16 ***
C	1.00576143	0.04523108	22.2361	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exemplo: Dados Simulados

```
tab %>% mutate(across(where(is.numeric), ~ round(., 3))) %>%  
  print()
```

	formula	coef	se	t_stat	p_val
1	D ~ E	1.758	0.021	83.819	0.000
2	D ~ E + C	0.853	0.031	27.735	0.000
3	D ~ E + A + B + C + F	0.056	0.038	1.493	0.136
4	D ~ E + A + B + C	1.030	0.033	31.350	0.000
5	D ~ E + A + C	1.058	0.039	26.824	0.000
6	D ~ E + B + C	1.038	0.026	39.822	0.000

Um pouco de teoria

Um pouco de teoria

Como vimos, as DAGs fornecem uma maneira conveniente de representar dependências estatísticas e causais entre uma coleção de variáveis.

Dadas as variáveis (X_1, \dots, X_p) , podemos escrever sua distribuição conjunta como

$$P(X_1, \dots, X_p) = \prod_{j=1}^p P(X_j | X_1, \dots, X_{j-1}), \quad (1)$$

em que X_0 é o conjunto vazio. As variáveis (X_1, \dots, X_{j-1}) são predecessoras (ancestrais) de X_j .

Suponha $P(X_j | X_1, \dots, X_{j-1}) = P(X_j | pa_j)$, em que pa_j é um subconjunto das predecessoras de X_j , os pais de X_j . Suponha que pa_j é o conjunto mínimo. Para X_1 , pa_j é o conjunto vazio.

No grafo simulado do exemplo

Os ancestrais de E são A , B , e C , enquanto os pais são A e C .

```
printCoefmat(coef(summary(lm(E ~ A + B + C))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.0017523	0.0306039	0.0573	0.9544
A	1.0611501	0.0438996	24.1722	<2e-16 ***
B	0.0560922	0.0431773	1.2991	0.1942
C	0.9922634	0.0303576	32.6858	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Note o efeito nulo de B sobre E , condicional em A e C .

No grafo simulado do exemplo

Os ancestrais de D são A , B , C e E , enquanto os pais são B , C e E .

```
printCoefmat(coef(summary(lm(D ~ A + B + C + E))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.0012754	0.0317176	0.0402	0.9679
A	0.0236395	0.0573091	0.4125	0.6801
B	0.9439513	0.0447865	21.0767	<2e-16 ***
C	1.0065989	0.0452955	22.2230	<2e-16 ***
E	1.0295130	0.0328393	31.3500	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Note o efeito nulo de A sobre D , condicional em B , C e E .

Podemos escrever a distribuição conjunta (1) como

$$P(X_1, \dots, X_p) = \prod_{j=1}^p P(X_j | pa_j). \quad (2)$$

Observe que (2) estabelece independências condicionais não especificadas por (1).

Essas relações condicionais são representadas de forma direta em um DAG.

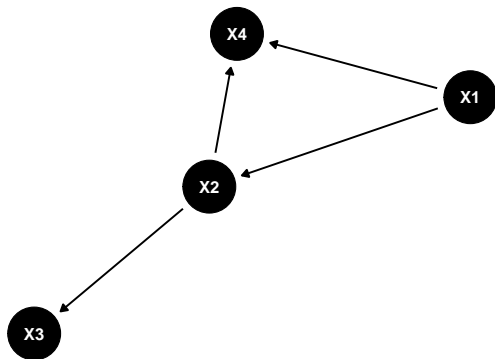
Por exemplo,

$$P(X_1, X_2, X_3, X_4) = P(X_4|X_2, X_1)P(X_3|X_2)P(X_2|X_1)P(X_1) \quad (3)$$

corresponde ao seguinte DAG.

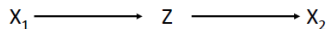
```
library(ggdag)
dag <- dagify(X2 ~ X1,
              X3 ~ X2,
              X4 ~ X2 + X1)
ggdag(dag) + theme_dag()
```

Um pouco de teoria

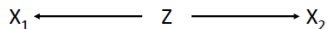


- O grafo compreende independências condicionais que são obtidas diretamente de (3): $X_4 \perp X_3 | X_2, X_1$, e $X_3 \perp X_1 | X_2$.
- É possível mostrar também que $X_4 \perp X_3 | X_2$ e $X_3 \perp X_1 | X_2, X_4$, que seriam difíceis de determinar diretamente da fatorização (3).

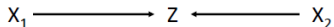
Considere as quatro estruturas básicas de DAGs:



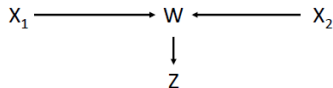
(a) Chain: Z is an Intermediate Variable



(b) Fork: Z is a Common Cause



(c) Collider: Z is a Common Effect



(d) Descendant of a Collider: Z is an Effect of a Common Effect

Considere 3 conjuntos de variáveis A , B , e C . Um caminho (em qualquer direção), é dito *d-separado*, ou bloqueado, por um conjunto de variáveis C se, e somente se,

- ❶ contém uma “chain” (como no painel (a)), tal que a variável Z está em C , ou um “fork” (como no painel (b)), tal que a variável Z está em C , ou
- ❷ contém um “fork” invertido, ou colisor (como no painel (c)), tal que a Z não está em C e tal que nenhum descendente de um colisor está em C (i.e. no painel (d), Z não pode estar em C , e nem ser W .)

O conjunto C é dito *d-separado* A de B se e somente se C bloqueia todo caminho de uma variável em A para uma variável em B .

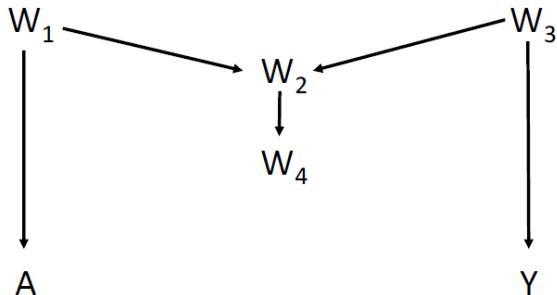
Theorem 1

Se A e B são d-separados por C em um DAG, então $A \perp B | C$. Por outro lado, se A e B não são d-separados por C , então A e B são dependentes, condicional em C , a menos que as dependências representadas pelas setas se cancelem exatamente.

De volta à figura, podemos ver que:

- No painel (a), temos que $X_1 \perp X_2 | Z$, mas não vale $X_1 \perp X_2$.
- No painel (b), temos que $X_1 \perp X_2 | Z$, mas não vale $X_1 \perp X_2$.
- No painel (c), temos que $X_1 \perp X_2$, mas não vale $X_1 \perp X_2 | Z$.
- No painel (d), temos que $X_1 \perp X_2$, $Z \perp X_1 | W$ e $Z \perp X_2 | W$, mas não vale $X_1 \perp X_2 | W$ ou $X_1 \perp X_2 | Z$.

Para ilustrar o teorema, considere o DAG da seguinte figura, que inclui um colisor.



Algumas das independências implicadas pelo DAG são:

$$A \perp Y$$

$$A \perp W_3$$

$$A \perp Y|W_3$$

$$A \perp W_3|Y$$

$$A \perp Y|W_1, W_2$$

$$A \perp Y|W_4, W_3$$

$$W_1 \perp W_3$$

$$W_1 \perp Y$$

$$W_1 \perp Y|A$$

$$W_4 \perp Y|W_2.$$

Algumas das dependências implicadas pelo DAG são:

$$A \not\perp W_1$$

$$A \not\perp W_2$$

$$A \not\perp Y|W_2$$

$$A \not\perp Y|W_4$$

$$A \not\perp W_1|W_2, W_4$$

$$W_1 \not\perp W_2$$

$$W_1 \not\perp W_4$$

$$W_1 \not\perp W_3|W_2$$

$$W_1 \not\perp W_3|W_2, W_4$$

$$W_3 \not\perp Y.$$

Theorem 2

Dados um DAG contendo as variáveis A e Y assim como um conjunto de variáveis C (excluindo A e Y) que não contém descendentes de A ou Y , o conjunto C é suficiente para ajuste de confundimento do efeito de A em Y se e somente se não houver nenhum caminho backdoor não bloqueado de A para Y . Isto é,

$$\{Y(a)\}_{a \in \mathcal{A}} \perp A | C,$$

em que $\{Y(a)\}$ é o conjunto de respostas potenciais de todos os valores $a \in \mathcal{A}$ de A .

No exemplo, o conjunto vazio (nenhuma variável em C) é suficiente para remover o confundimento do efeito de A em Y . São suficientes também os conjuntos W_1 , ou o conjunto (W_1, W_2) , ou (W_3, W_4) .

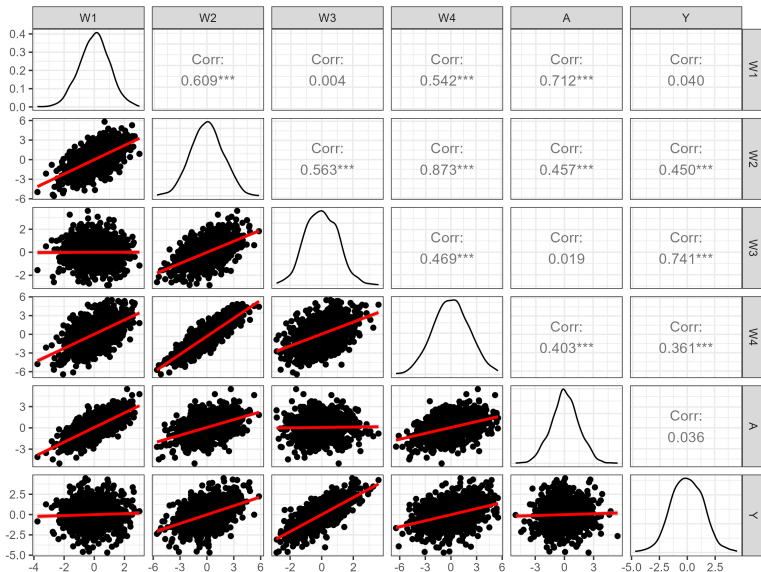
Como vimos, W_2 sozinho é insuficiente. Isso prova que uma definição tradicional de confundidor como uma variável que está associada com A e Y é inadequada (W_2 é um exemplo!). O DAG implica que ajustar por W_2 é pior do que não ajustar para nenhuma variável.

Ajustar para W_2 ou W_4 sem também ajustar para W_1 ou W_3 pode causar um viés chamado *viés de colisor*.

Exemplo: Dados Simulados

```
set.seed(123456)
n <- 1000
W1 <- rnorm(n, mean = 0)
W3 <- rnorm(n, mean = 0)
A <- rnorm(n, mean = W1)
Y <- rnorm(n, mean = W3)
W2 <- rnorm(n, mean = W1 + W3)
W4 <- rnorm(n, mean = W2)
```

Exemplo: Dados Simulados



Exemplo: Dados Simulados

```
# Não existe associação entre A e Y  
# Beta real para a exposição A é zero  
printCoefmat(coef(summary(f1 <- lm(Y ~ A))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.0086703	0.0455781	-0.1902	0.8492
A	0.0362925	0.0317562	1.1428	0.2534

Exemplo: Dados Simulados

Seria suficiente ajustar por W1

```
printCoefmat(coef(summary(f2 <- lm(Y ~ A + W1))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.0078544	0.0456103	-0.1722	0.8633
A	0.0159940	0.0452338	0.3536	0.7237
W1	0.0412380	0.0654226	0.6303	0.5286

Exemplo: Dados Simulados

```
# Ou ajustar por W1 e W2
```

```
printCoefmat(coef(summary(f3 <- lm(Y ~ A + W1 + W2))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.0036572	0.0384708	-0.0951	0.9243
A	-0.0158132	0.0381854	-0.4141	0.6789
W1	-0.5260916	0.0619584	-8.4910	<2e-16 ***
W2	0.5474490	0.0271887	20.1352	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exemplo: Dados Simulados

```
# Ou ajustar por W3 e W4
```

```
printCoefmat(coef(summary(f4 <- lm(Y ~ A + W3 + W4))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.0005002	0.0306201	0.0163	0.9870
A	0.0192256	0.0238385	0.8065	0.4202
W3	1.0426234	0.0347992	29.9611	<2e-16 ***
W4	0.0054817	0.0185292	0.2958	0.7674

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exemplo: Dados Simulados

```
# Ajustar apenas por W2 induz o viés de colisor  
printCoefmat(coef(summary(f5 <- lm(Y ~ A + W2))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.0037939	0.0398086	0.0953	0.9241
A	-0.2151174	0.0311748	-6.9004	9.21e-12 ***
W2	0.4424636	0.0250632	17.6539	< 2.2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exemplo: Dados Simulados

```
# Ajustar apenas por W4 também induz o viés de colisor  
printCoefmat(coef(summary(f6 <- lm(Y ~ A + W4))))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.0082494	0.0421980	-0.1955	0.845
A	-0.1309562	0.0321191	-4.0772	4.921e-05 ***
W4	0.2852735	0.0220562	12.9339	< 2.2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Exemplo: Dados Simulados

```
tab %>% mutate(across(where(is.numeric), ~ round(., 3))) %>%  
  print()
```

	formula	coef	se	t_stat	p_val
1	Y ~ A	0.036	0.032	1.143	0.253
2	Y ~ A + W1	0.016	0.045	0.354	0.724
3	Y ~ A + W1 + W2	-0.016	0.038	-0.414	0.679
4	Y ~ A + W1 + W2	-0.016	0.038	-0.414	0.679
5	Y ~ A + W2	-0.215	0.031	-6.900	0.000
6	Y ~ A + W4	-0.131	0.032	-4.077	0.000

W_2 não é um *confundidor verdadeiro* (variável que influencia tanto a exposição como o desfecho via um caminho direcionado que não inclui a exposição), porque não influencia (não causa) nem A nem Y .

Neste exemplo, não há *confundidores verdadeiros*, e, portanto, não há confundimento do efeito de A em Y .

Em alguns casos, nem todos os confundidores verdadeiros são necessários para formar um *conjunto suficiente de confundidores verdadeiros*.

Bibliografia

Principais Referências Usadas

Pearl, J., “Causality. Models, Reasoning and Inference”, *Cambridge University Press*, 2009.

Reis, R. C. P. “MAT02010 - Tópicos Avançados em Estatística II. Introdução a Inferência Causal”. *UFRGS*, 2019.

Brumback, B. A., “Fundamentals of Causal Inference With R”, *CRC Press*, 2022.

Morgan, S. L. and Winship, C., “Counterfactuals and Causal Inference”, *Cambridge University Press*, 2015.