

CONSIDERAÇÕES SOBRE INFERÊNCIA CAUSAL & LINCA



Marcelo Magalhães Taddeo

marcelo.magalhaes@ufba.br

28 de setembro de 2023



Laboratório de Inferência Causal

- ▶ Marcelo M. Taddeo (IME/UFBa)
- ▶ Leila Amorim (IME/UFBa)
- ▶ Rosemeire L. Fiaccone (IME/UFBa)
- ▶ Lilia Costa (IME/UFBa)
- ▶ Raydonal Ospina (IME/UFBa)
- ▶ Rosana Aquino (ISC/UFBa)
- ▶ Vinicius Mendes (FE/UFBa)
- ▶ Elzo Júnior (Cidacs/Fiocruz)

+ colaboradores e estudantes de graduação e pós-graduação.

Grupo de pesquisa

LInCa - Laboratório de Inferência Causal e Aplicações.

Endereço para acessar este espelho: dgp.cnpq.br/dgp/espelhogrupo/9296811361381732

Identificação

Endereço / Contato

Repercussões

Linhas de pesquisa

Recursos humanos

Instituições parceiras

Indicadores de RH

Equipamentos e Softwares

Identificação

Situação do grupo: Certificado

Ano de formação: 2023

Data da Situação: 15/09/2023 09:04

Data do último envio: 23/09/2023 16:04

Líder(es) do grupo: Leila Denise Alves Ferreira Amorim

Marcelo Magalhães Taddeo

Área predominante: Ciências Exatas e da Terra; Probabilidade e Estatística



Endereço para acessar este espelho: dgp.cnpq.br/dgp/espelhogrupo/9296811361381732

Repercussões dos trabalhos do grupo

O LInCa é um grupo interdisciplinar com enfoque intersetorial composto de pesquisadores com vasta experiência no desenvolvimento de metodologias estatísticas (inferência estatística, análise de dados longitudinais e de sobrevivência, aprendizado de máquina, inferência bayesiana, séries temporais etc.) e com forte componente aplicado, especialmente, em bioestatística, epidemiologia e econometria. Os pesquisadores deste grupo têm avançado na produção científica sobre modelagem estatística e suas aplicações em diversas áreas e estão inseridos em ampla rede de colaborações nacionais e internacionais. O LInCa se propõe a estudar e desenvolver metodologias estatísticas para a inferência causal, especialmente em de estudos observacionais (não aleatorizados), e suas aplicações nas análises de dados reais, e também na formação de recursos humanos para atuar na avaliação do impacto de políticas públicas, facilitando a tomada de decisões.

Linhas de pesquisa

| Nome da linha de pesquisa | Quantidade de Estudantes | Quantidade de Pesquisadores |
|---|--------------------------|-----------------------------|
| Análise de mediação causal | 0 | 2 |
| Avaliação de impacto em epidemiologia e econometria | 4 | 10 |
| Métodos estatísticos em inferência causal | 3 | 7 |

ÁREAS DE INTERESSE:

1. Metodologias de identificação causal em diferentes contextos
 - (a) Mediação causal
 - (b) Dados de sobrevivência
 - (c) Modelos com variáveis latentes
 - (d) Dados agrupados
 - (e) Dados funcionais
2. Metodologias de estimação
 - (a) Variações e extensões dos métodos clássicos
 - (b) Métodos bayesianos
 - (c) Técnicas de machine learning no aprendizado causal
3. Aplicações
 - (a) Epidemiologia
 - (b) Econometria
 - (c) Educação



Capítulo 8

Diagramas causais e equações estruturais na avaliação de políticas públicas

SAVE THE DATE

O Departamento de Ciência e Tecnologia convidou para o Lançamento do livro

AVALIAÇÃO DE IMPACTO DAS POLÍTICAS DE SAÚDE: UM GUIA PARA O SUS

Marcelo M. Taddeo¹
Leila Denise Amorim¹
Rosana Aquino²

frontiers
in Pharmacology

REVIEW
published: 18 September 2019
doi: 10.3389/fphr.2019.00073

Propensity Score Methods in Health Technology Assessment: Principles, Extended Applications, and Recent Advances

M Sanni Ali^{1,2*}, Daniel Prieto-Alhambra^{1,4}, Luciane Cruz Lopes¹, Dandara Ramos⁵, Nivea Bispo⁶, Maria Y. Ichihara^{1,5}, Julia M. Pescarini³, Elizabeth Williamson¹, Rosemeire L. Fiaccone^{1,7}, Mauricio L. Barreto^{1,8} and Liam Smeeth^{1,2}

STATISTICS AND ITS INTERFACE Volume 15 (2022) 399–413

Causal measures using generalized difference-in-difference approach with nonlinear models

MARCELO M. TADDEO*, LEILA D. AMORIM^{†‡§}, AND ROSANA AQUINO^{*†}

Revista Colombiana de Estadística - Applied Statistics
January 2022, volume 45, issue 1, pp. 161-191
<http://doi.org/10.15446/rce.v45n1.94553>

Causal Mediation for Survival Data: A Unifying Approach via GLM

Mediación causal para datos de supervivencia: un enfoque unificador a través de GLM

MARCELO M. TADDEO*, LEILA D. AMORIM^b

DEPARTAMENTO DE ESTATÍSTICA, INSTITUTO DE MATEMÁTICA E ESTATÍSTICA, UNIVERSIDADE FEDERAL DA BAHIA, SALVADOR, BRAZIL

Por que entender, caracterizar e determinar relações de causalidade?

Propósito central da ciência: **explicar** e **prever** fenômenos de interesse.

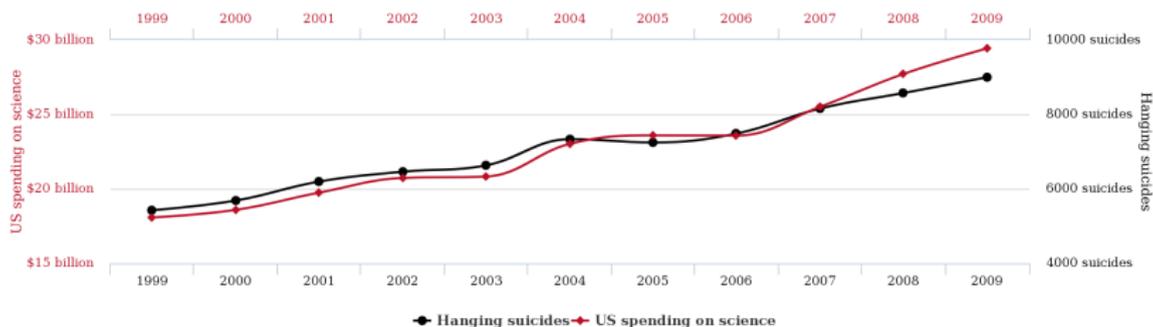
Popper: o cientista é movido pela busca de explicações causais com o objetivo “de encontrar *teorias explicativas* (...) que descrevam certas propriedades estruturais do mundo e que nos permitam deduzir, com o auxílio de condições iniciais, os efeitos que se pretende explicar.”

⇒ encontrar leis (universais) tais que, combinadas com certas condições iniciais (**causas**), permitam explicar determinados acontecimentos (**efeitos**).

⇒ identificar propriedades causais de um conjunto **especificado** de estruturas, variáveis ou circunstâncias.

CORRELAÇÃO E CAUSAÇÃO

US spending on science, space, and technology correlates with Suicides by hanging, strangulation and suffocation



tylervigen.com

Fonte: U.S. Office of Management and Budget and Centers for Disease Control & Prevention

► <http://www.tylervigen.com/spurious-correlations>

INFERÊNCIA CAUSAL E ESTATÍSTICA

“The aim of standard **statistical analysis**, typified by regression, estimation, and hypothesis testing techniques, is to assess parameters of a distribution from samples drawn of that distribution. (...)” (Pearl)

TRADUZINDO:

- ▶ Inferência estatística é sobre os parâmetros da distribuição populacional...
- ▶ ... e não sobre relações causais.

INFERÊNCIA CAUSAL E ESTATÍSTICA

“The aim of standard **statistical analysis**, typified by regression, estimation, and hypothesis testing techniques, is to assess parameters of a distribution from samples drawn of that distribution. (...)” (Pearl)

TRADUZINDO:

- ▶ Inferência estatística é sobre os parâmetros da distribuição populacional...
- ▶ ... e não sobre relações causais.

“**Causal analysis** (...) [aims] to infer not only beliefs or probabilities under static conditions, but also the dynamics of beliefs under **changing conditions**, for example, changes induced by **treatments** or **external interventions**” (Pearl)

INFERÊNCIA CAUSAL E ESTATÍSTICA

“The aim of standard **statistical analysis**, typified by regression, estimation, and hypothesis testing techniques, is to assess parameters of a distribution from samples drawn of that distribution. (...)” (Pearl)

TRADUZINDO:

- ▶ Inferência estatística é sobre os parâmetros da distribuição populacional...
- ▶ ... e não sobre relações causais.

“**Causal analysis** (...) [aims] to infer not only beliefs or probabilities under static conditions, but also the dynamics of beliefs under **changing conditions**, for example, changes induced by **treatments** or **external interventions**” (Pearl)

- ▶ Isto não é totalmente verdadeiro... *(por quê?)*

On the Application of Probability Theory to Agricultural Experiments. Essay on Principles. Section 9.

Jerzy Splawa-Neyman

Translated and edited by D. M. Dabrowska and T. P. Speed from the Polish original, which appeared in *Roczniki Nauk Rolniczych Tom X (1923) 1–51 (Annals of Agricultural Sciences)*

m : # áreas / campos cultiváveis (plots)

U_k : rendimentos **observados** nos campos

$$k = 1, \dots, m$$

n : # variedades de culturas

↪ cada campo contém só uma variedade!

Objetivo: comparar os rendimentos das culturas.

▶ Para cada cultura $i = 1, \dots, n$, defina

$$U_{i1}, \dots, U_{im}.$$

▶ **Retorno médio da variedade i :**

$$a_i = \frac{1}{m} \sum_{k=1}^m U_{ik}.$$

“The goal of a field experiment which consists of the comparison of v varieties will be regarded as equivalent to the problem of comparing the numbers

$$a_1, \dots, a_n.$$

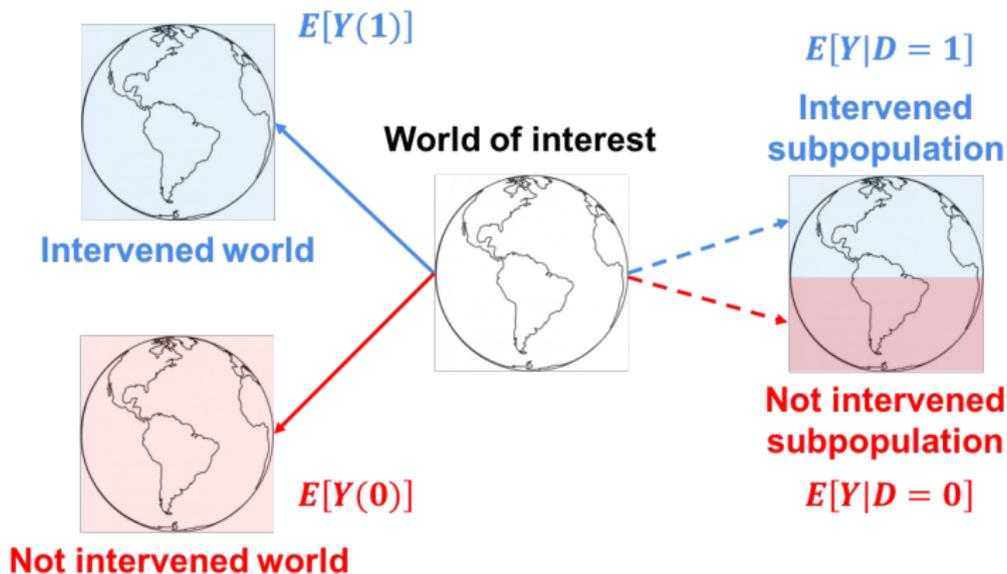
or their estimates (...).”

Obs.: para cada i , somente um U_{ik} é observado...

▶ os outros são **contrafactuais!**

CONTRAFACTUALS

[A cause is] "an object followed by another, and where all the objects, similar to the first, are followed by objects similar to the second. Or, in other words, where, if the first object had not been, the second never had existed." (Hume, 1748)



RUBIN (1974)¹

A : atribuição de tratamento

▶ Para Rubin, **tratamento** \equiv **causa**

$A = 1 \Rightarrow$ 'Tratamento'

$A = 0 \Rightarrow$ 'Controle'

Y : resposta de interesse

$Y(1)$: resposta quando a unidade é exposta ao tratamento;

$Y(0)$: resposta quando a unidade é exposta ao controle.

Objetivo: avaliar as quantidades $Y(1)$ e $Y(0)$ ou contrastes entre elas em estudos observacionais.

¹Rubin, D. B. (1974), "Estimating causal effects of treatments in randomized and nonrandomized studies," J. Educ. Psychology, 66, 688–701.

RUBIN (1974)¹

A: atribuição de tratamento

▶ Para Rubin, **tratamento** \equiv **causa**

$A = 1 \Rightarrow$ 'Tratamento'

$A = 0 \Rightarrow$ 'Controle'

Y: resposta de interesse

$Y(1)$: resposta quando a unidade é exposta ao tratamento;

$Y(0)$: resposta quando a unidade é exposta ao controle.

Objetivo: avaliar as quantidades $Y(1)$ e $Y(0)$ ou contrastes entre elas em estudos observacionais.

PROBLEMA FUNDAMENTAL DA INFERÊNCIA CAUSAL

É IMPOSSÍVEL **OBSERVAR** SIMULTANEAMENTE $Y(1)$ E $Y(0)$! UMA DESSAS QUANTIDADES É OBRIGATORIAMENTE **CONTRAFACTUAL**.

¹Rubin, D. B. (1974), "Estimating causal effects of treatments in randomized and nonrandomized studies," J. Educ. Psychology, 66, 688–701.

EFEITOS CAUSAIS

▶ **Def.:** qualquer avaliação estatística de $Y(a)$.

▶ **Exemplo:**

$$P(Y = y | \text{do}(A = a)).$$

▶ **Tipicamente:** qualquer contraste $\tau(Y(0), Y(1); X, A)$ entre $Y(0)$ e $Y(1)$.

▶ **Ainda mais tipicamente:** contrastes em termos de médias.

▶ **Average Treatment Effect (ATE):**

$$E[Y(1) - Y(0)],$$

ou sua contraparte para populações finitas

$$\frac{1}{N} \sum_{i=1}^N (Y_i(1) - Y_i(0)).$$

▶ **Inclusão de covariáveis:**

$$E[Y(1) - Y(0) | X = x] \quad \text{ou} \quad \frac{1}{N} \sum_{i: X_i = x} (Y_i(1) - Y_i(0)).$$

▶ **Caso especial:** **Average Treatment (Effect) on the Treated (ATT)**

$$E[Y(1) - Y(0) | A = 1] \quad \text{ou} \quad \frac{1}{N_1} \sum_{i: A_i = 1} (Y_i(1) - Y_i(0)).$$

EFEITOS CAUSAIS — OUTRAS MEDIDAS

- ▶ Se Y é binária:

$$E[Y(1) - Y(0)] = P(Y(1) = 1) - P(Y(0) = 1),$$

- ▶ ou o **risco relativo**

$$\frac{P(Y(1) = 1)}{P(Y(0) = 1)},$$

- ▶ ou a **razão de chances**

$$\frac{P(Y(1) = 1)/P(Y(1) = 0)}{P(Y(0) = 1)/P(Y(0) = 0)}.$$

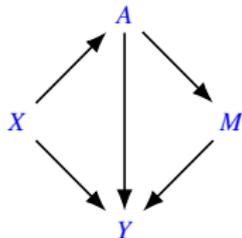
- ▶ As mesmas quantidades podem ser avaliadas entre os tratados, e.g.,

$$\frac{P(Y(1) = 1|A = 1)}{P(Y(0) = 1|A = 1)}.$$

DIAGRAMAS CAUSAIS — UMA OUTRA PERSPECTIVA

MÉTODO GRÁFICO PARA REPRESENTAR RELAÇÕES CAUSAIS ENTRE VARIÁVEIS.

- ▶ Um **grafo** \mathcal{G} é um par $(\mathcal{V}, \mathcal{A})$, onde \mathcal{V} é um conjunto finito de **vértices** e $\mathcal{A} \subset \mathcal{V} \times \mathcal{V}$ um conjunto de **arestas** em \mathcal{G} .
- ▶ Este grafo representa as relações causais entre as variáveis envolvidas no sentido de que $A \longrightarrow B$ significa A **causa** B .



- ▶ Uma aresta do tipo

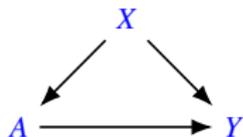


é chamada **direcionada**.

- ▶ Um grafo somente com arestas direcionadas e sem “loops” é chamado **grafo acíclico direcionado (DAG)**.
- ▶ Cada estrutura $U \longrightarrow V$ é um mecanismo.
- ▶ A totalidade de mecanismos representa o fenômeno de interesse.

- ▶ Distribuição conjunta: $P(a, m, x, y)$
- ▶ Decomposição induzida: $P(a, m, x, y) = P(y|a, m, x) P(m|a) P(a|x) P(x)$

IDENTIFICAÇÃO DO EFEITO CAUSAL



A: Exposição / Tratamento

X: Variáveis confundidoras (covariáveis)

Y: Desfecho

PERGUNTA:

- ▶ Como estimar o efeito causal médio (ATE):

$$E[Y(1) - Y(0)|X]?$$

HIPÓTESES:

Consistência: $Y = AY(1) + (1 - A)Y(0)$

Ignorabilidade: $\{Y(0), Y(1)\} \perp\!\!\!\perp A|X$

E o ATE?

- ▶ Dadas as condições acima:

$$E[Y(a)|X = x] \stackrel{\text{Ign}}{=} E[Y(a)|A = a, X = x]$$

$$\stackrel{\text{Con}}{=} E[Y|A = a, X = x].$$

$$\begin{aligned} \therefore E Y(a) &= E E[Y(a)|X] \\ &= E E[Y|A = a, X]. \end{aligned}$$

- ▶ Ou seja,

$$\begin{aligned} E Y(1) - E Y(0) &= E E[Y|A = 1, X] \\ &\quad - E E[Y|A = 0, X]. \end{aligned}$$

ATE \neq ATT!

Sob as hipóteses de identificação (slide anterior):

$$\begin{aligned} \text{ATT} &\stackrel{\text{Def}}{=} E[Y(1)|A = 1] - E[Y(0)|A = 1] \\ &= E[Y|A = 1] - E E[Y(0)|A = 1, X] \\ &\stackrel{\text{Ign}}{=} E[Y|A = 1] - E E[Y(0)|A = 0, X] \\ &\stackrel{\text{Con}}{=} E[Y|A = 1] - E E[Y|A = 0, X] \end{aligned}$$

Ilustração:

► Desfecho: $Y = \beta_0 + \tau A + \beta_1 X + \varepsilon$

⇒ $ATE = \tau$.

► $E[Y|A = 0, X] = \beta_0 + \beta_1 X$

$$\Rightarrow E E[Y|A = 0, X] = \beta_0 + \beta_1 E X.$$

► $E[Y|A = 1] = \beta_0 + \tau + \beta_1 E[X|A = 1]$

⇒ $ATT = (\beta_0 + \tau + \beta_1 E[X|A = 1]) - (\beta_0 + \beta_1 E X) = \tau + \beta_1 \{E[X|A = 1] - E X\}$

► Conclusão:

$$ATT - ATE = \beta_1 \{E[X|A = 1] - E X\}$$

e

$$ATT = ATE \Leftrightarrow E[X|A = 1] = E X$$

Viés de Seleção

Sob a condição de consistência:

$$\begin{aligned} E[Y|A = 1, x] - E[Y|A = 0, x] &= E[Y(1)|A = 1, x] - E[Y(0)|A = 0, x] \\ &= E[Y(1)|A = 1, x] - E[Y(0)|A = 1, x] \\ &\quad + \{E[Y(0)|A = 1, x] - E[Y(0)|A = 0, x]\} \\ &= \text{ATT}(x) + \underbrace{\{E[Y(0)|A = 1, x] - E[Y(0)|A = 0, x]\}}_{\text{Viés de Seleção}} \end{aligned}$$

- ▶ Portanto, a fórmula de identificação vale somente se **NÃO** houver viés de seleção!
- ▶ No caso do ATE, temos algo parecido...

$$E[Y|A = 1, x] - E[Y|A = 0, x] = \underbrace{E[Y(1)|A = 1, x]}_{\neq E[Y(1)|x]} - \underbrace{E[Y(0)|A = 0, x]}_{\neq E[Y(0)|x]}$$

Escore de Propensão (ATE)

- ▶ Se, além das condições usuais

(A1) (Consistência) $Y = (1 - A)Y(0) + AY(1)$

(A2) (Ignorabilidade) $Y(0), Y(1) \perp\!\!\!\perp A|X$

valer

(A3) (Positividade) Para todo $x \in \mathcal{X}$, vale $0 < P(A = 1|X = x) < 1$.

Sob as condições (A1) – (A2),

$$ATE = E \left[\frac{AY}{p(X)} - \frac{(1 - A)Y}{1 - p(X)} \right] = E \left[\frac{A - p(X)}{p(X)(1 - p(X))} Y \right].$$

- ▶ Estimadores:

(A) Horvitz-Thompson:

$$\widehat{ATE} = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{A_i Y_i}{\pi(X_i|\hat{\theta})} - \frac{(1 - A_i) Y_i}{1 - \pi(X_i|\hat{\theta})} \right\}$$

(B) IPTW: pesos padronizados

$$\omega_i = \frac{A_i}{\pi(X_i|\hat{\theta})} / \sum_{j=1}^n \frac{A_j}{\pi(X_j|\hat{\theta})}$$

(no caso dos tratados)

Escore de Propensão (ATT)

Sob as mesmas condições,

$$\begin{aligned} \text{ATT} &= \frac{1}{EA} E \left[AY - (1 - A) \frac{p(X)}{1 - p(X)} Y \right] \\ &= \frac{1}{EA} E \left[\frac{A - p(X)}{1 - p(X)} Y \right]. \end{aligned}$$

► Estimador: tomando

$$\widehat{EA} = \widehat{P}(A = 1) = \frac{n_1}{n},$$

temos

$$\widehat{\text{ATT}} = \frac{1}{n_1} \left\{ \sum_{i:A_i=1} Y_i - \sum_{i:A_i=0} \frac{\pi(X_i|\widehat{\theta})}{1 - \pi(X_i|\widehat{\theta})} Y_i \right\}$$

Obs.: o ATT demanda uma **condição de positividade mais fraca** do que o ATE,

$$P(A = 0|X = x) > 0,$$

para todo x observável entre os não-tratados.

MÉTODOS DE IDENTIFICAÇÃO DO EFEITO CAUSAL (ATE)

1. Regressão (paramétrica, não paramétrica)
2. Métodos baseados em escores de propensão:
 - (a) Ponderação pelo inverso da probabilidade
 - (b) Matching / Pareamento
3. Balanceamento de covariáveis (Calibração)
4. Métodos duplo robustos
5. Diferenças em diferenças (diff-in-diff)
6. RDD
7. Variáveis instrumentais
8. Modelos marginais estruturais (MSM) para dados longitudinais
9. Métodos bayesianos