

# INFERÊNCIA CAUSAL COM MACHINE LEARNING

uma aplicação para evasão fiscal

---

Rafael Felipe Bressan

2021-01-09

Receita Federal do Brasil

# Motivação

---

- Limite de velocidade reduz as mortes no trânsito?
- Permissão para cobrança de bagagem aérea reduziu o preço das tarifas?
- O recebimento de uma carta-cobrança da Receita Federal faz com que o contribuinte recolha seus impostos devidos?
- Essas questões são **causais** em sua natureza. Requerem conhecimento do processo de geração dos dados. Suas respostas não podem ser calculadas apenas com os dados observados.

- Análise causal requer manipulação/intervenção no processo gerador
- Uma quebra estrutural é induzida
- Correlações anteriores não são mais válidas
- Dados puramente observacionais não carregam toda a informação necessária

$$Y_i = f(X_i, \epsilon_i; \theta)$$

- Causalidade requer inferência sobre parâmetros da distribuição,  $\theta$ 
  - *Machine Learning* tradicional oferece correlações a partir de dados observacionais
  - Inferência  $\neq$  previsão
    - ML: minimiza  $\hat{e} = \hat{y} - Y$
    - Análise causal: estima  $\hat{\theta}$  com intervalo de confiança
  - Boa previsão **não garante** correta estimação de parâmetros
  - **Viés de regularização**:  $\hat{f}_1(\cdot; \hat{\theta}_1) \approx \hat{f}_2(\cdot; \hat{\theta}_2)$  mesmo se  $\hat{\theta}_1 \neq \hat{\theta}_2$

- Como fazer com que algoritmos de ML façam estimação causal não-viesada?
- Fronteira do conhecimento em inferência causal
  - Chernozhukov et al. (2018) - *Double Machine Learning*
  - Wager and Athey (2018) - *Causal Forests*
  - Syrgkanis et al. (2019) - *Doubly Robust Instrumental Variables*

# **Experimento Randomizado**

---

# Experimento Randomizado

- Experimentos randomizados são o padrão-ouro para inferência causal
- Re-analisaremos o trabalho de Fellner, Sausgruber, and Traxler (2013)
- Correspondências fiscais para mais de 50.000 contribuintes
- Analisar efeitos de variação no conteúdo
  - Valores médios por tipo de carta
  - Heterogeneidade nos efeitos



# Modelo ForestDML

- Modelo parcialmente linear. Tratamento  $T$  é exógeno, não é necessária instrumentalização

$$Y = \theta(X) \cdot T + g(X, W) + \epsilon$$

$$\mathbb{E}[\epsilon \mid X, W] = 0$$

$$T = f(X, W) + \eta$$

$$\mathbb{E}[\eta \mid X, W] = 0$$

$$\mathbb{E}[\eta \cdot \epsilon \mid X, W] = 0$$

- Através de **DML** (ortogonalização de Neyman e *cross-fitting*)

$$\hat{\theta}(x) = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n K_x(X_i) \cdot \left( Y_i - \hat{q}(X_i, W_i) - \theta \cdot \left( T_i - \hat{f}(X_i, W_i) \right) \right)^2$$

- Kernel  $K_x$  é uma **floresta causal**

- Tratamento é endógeno. Necessita de variável instrumental

$$Y = \theta(X) \cdot T + g(X) + \epsilon, \quad \mathbb{E}[\epsilon \mid X, Z] = 0$$

$$Z = m(X) + \eta, \quad \mathbb{E}[\eta \mid X] = 0$$

$$\mathbb{E}[\eta \cdot \epsilon \mid X, Z] = 0$$

$$\mathbb{E}[T \cdot \epsilon \mid X] \neq 0$$

- Estimativa preliminar de  $\theta(x)$  e algoritmo *Doubly Robust*

$$\hat{\theta}_{DR}(x) = \operatorname{argmin}_{\theta} \sum_{i \in \mathcal{I}} \left( \theta_{\text{pre}}(x) + \frac{\left( \hat{Y}_i - \theta_{\text{pre}}(x) \hat{T}_i \right) \hat{Z}_i}{\hat{\beta}(X_i)} - \theta(X_i) \right)^2$$

# Resultados

- Receber uma correspondência **tem efeito positivo** sobre o registro para pagamento do tributo
- Uma **ameaça** na carta aumenta este efeito
- Informações e apelo moral não possui efeito estatisticamente significativo

	OLS ATE	ForestDML ATE	ATT	IV2SLS LATE	DRIV LATE
Correio	0,0650	0,0766	0,0766	0,0767	0,0588
<b>Ameaça</b>	<b>0,0750</b>	<b>0,0850</b>	<b>0,0848</b>	<b>0,0872</b>	<b>0,0650</b>
Info	0,0646	0,0762	0,0760	0,0728	0,0547
Moral	0,0648	0,0695	0,0695	0,0724	0,0513

Chernozhukov, Victor, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. 2018. “Double/debiased machine learning for treatment and structural parameters.” *The Econometrics Journal* 21 (1): C1–C68. <https://doi.org/10.1111/ectj.12097>.

Fellner, Gerlinde, Rupert Sausgruber, and Christian Traxler. 2013. “Testing Enforcement Strategies in the Field: Threat, Moral Appeal and Social Information.” *Journal of the European Economic Association* 11 (3): 634–60.

Syrgkanis, Vasilis, Victor Lei, Miruna Oprescu, Maggie Hei, Keith Battocchi, and Greg Lewis. 2019. “Machine Learning Estimation of Heterogeneous Treatment Effects with Instruments.” <http://arxiv.org/abs/1905.10176>.

Wager, Stefan, and Susan Athey. 2018. “Estimation and Inference of Heterogeneous Treatment Effects Using Random Forests.” *Journal of the American Statistical Association* 113 (523): 1228–42.