# Practical do-Shapley Explanations with Estimand-Agnostic Causal Inference

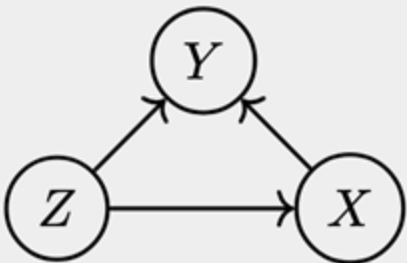Álvaro Parafita, Tomas Garriga,
Axel Brando, Francisco J. Cazorla

# Introduction

- **Goal**: make causal explanations practical.

- **do-SHAP** (Jung et al. 2022): Shapley Values with *causal interventional* effects.

$$\phi_{\mathbf{x}}(X_k) = \sum_{\mathbf{S} \subseteq \mathbf{X} \setminus \{X_k\}} w(|\mathbf{S}|) \cdot (\nu_{\mathbf{x}}(\mathbf{S} \cup \{X_k\}) - \nu_{\mathbf{x}}(\mathbf{S})), \quad \nu_{\mathbf{x}}(\mathbf{S}) = \mathbb{E}[Y \mid do(\mathbf{S} = \mathbf{x_S})]$$

- Two main **contributions**:

  1) A method to estimate $\nu_{\mathbf{x}}(\mathbf{S})$ in an automatable, practical way.

  2) A method to reduce the number of coalitions that need to be evaluated.

Jung, Y., Kasiviswanathan, S., Tian, J., Janzing, D., Blöbaum, P., & Bareinboim, E. (2022, June). On measuring causal contributions via do-interventions. In *International Conference on Machine Learning* (pp. 10476-10501). PMLR.
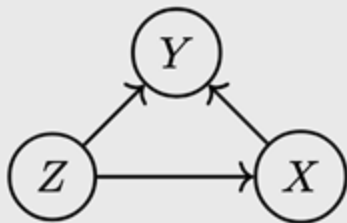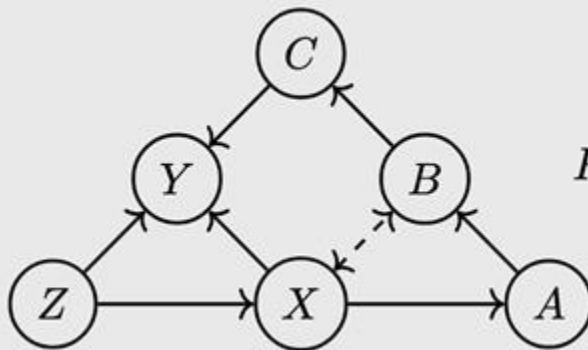
## 1) Why is do-SHAP impractical?

- By default, do-SHAP employs the **Estimand-Based approach**.

- Given a query (e.g., the effect of X on Y), compute an **estimand** of the query.
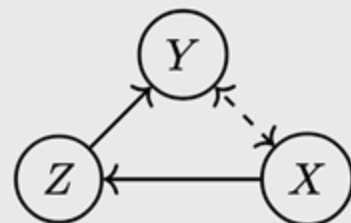
$$P_x(Y) = \mathbb{E}_Z\left[P(Y \mid x, Z)\right]$$

# Estimand-Based approach



$$P_x(Y) = \mathbb{E}_Z \left[ P(Y \mid x, Z) \right]$$

$$P_x(Y) = \mathbb{E}_{Z \mid x} \left[ \mathbb{E}_X \left[ P(Y \mid X, Z) \right] \right]$$

$$P_x(Y) = \mathbb{E}_{Z,C} \left[ \frac{P(Y \mid x, Z, C)}{P(Z, C)} \cdot \mathbb{E}_{A \mid x} \left[ \mathbb{E}_X \left[ P(Z, C \mid X, A) \right] \right] \right]$$
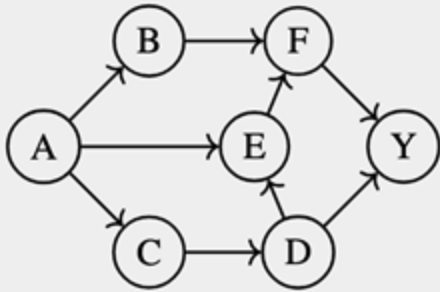
Parafita, Á. (2022). *Causality without Estimands: from Causal Estimation to Black-Box Introspection* (Doctoral dissertation, Universitat de Barcelona (Spain)).

# Estimands are impractical for do-SHAP

- **Problem**: do-SHAP needs to evaluate multiple queries, one per coalition:

$$\nu_{\mathbf{x}}(\mathbf{S}) = \mathbb{E}[Y \mid do(\mathbf{S} = \mathbf{x_S})]$$

- **Alternative:** Estimand-Agnostic Approach.

  - Train **a single Structural Causal Model**, following the known graph.

  - Use **general estimation procedures** to estimate any $\nu_{\mathbf{x}}(\mathbf{S})$.

- Guaranteed results as long as the queries are **identifiable.**

## 2) Reduce coalition evaluations

- Remove all **superfluous** nodes from any coalition.

$$\left.\begin{array}{l}\{A, B, C, E, F\}, \{A, B, C, F\}, \\ \{A, E, C, F\}, \{B, E, C, F\}, \\ \{A, C, F\}, \{B, C, F\}, \{E, C, F\}\end{array}\right\} = \{C, F\}$$

- **Cache**: only compute and store values for irreducible coalitions.

# How to find the irreducible coalition?

*Theorem.* Given a topological order $<_\mathcal{G}$ in $\mathcal{G}$ and $\mathbf{S} \subseteq \mathbf{X}$, let $\mathbf{Z} := \{X \in \mathbf{S} \mid \mathbf{S}_{>_\mathcal{G} X} \in Fr_\mathcal{G}(X, Y)\}$, with $\mathbf{S}_{>_\mathcal{G} X} := \{Z \in \mathbf{S} \mid Z >_\mathcal{G} X\}$. Then $\nu(\mathbf{S}) = \nu(\mathbf{S} \setminus \mathbf{Z})$, and $\mathbf{S} \setminus \mathbf{Z}$ is irreducible.

---

**Algorithm 1** Frontier-Reducibility Algorithm (FRA) – set version

---

**Require:** $\mathbf{S} \subseteq \mathbf{X}$, coalition.
**Require:** Fr, a map: tuple[int] $\rightarrow$ bool.
**Require:** $\mathcal{G}$, causal graph.

1: **procedure** FRA($\mathbf{S}, \text{Fr}; \mathcal{G}$)
2:     SORT($\mathbf{S}, <_\mathcal{G}$)
3:     $\mathbf{P} \leftarrow \varnothing$
4:     $\mathbf{Z} \leftarrow \varnothing$
5:     $k \leftarrow |\mathbf{S}|$
6:     **while** $k > 0$ **do**
7:         $X \leftarrow \mathbf{S}[k]$
8:         **if** $X \notin Pa_\mathcal{G}(Y)$ **then**
9:             $\mathbf{P'} \leftarrow \mathbf{P} \cap De_\mathcal{G}(X)$
10:             $\mathbf{T} \leftarrow (\mathbf{P'} \setminus \mathbf{Z}) \cup \{X\}$
11:             **if** $\mathbf{T} \notin \text{Fr}$ **then**
12:                 $\mathbf{C} \leftarrow \{X\}$
13:                 **while** $\mathbf{C} \neq \varnothing$ and $Y \notin \mathbf{C}$ **do**
14:                     $\mathbf{P'} \leftarrow \mathbf{P'} \cup \mathbf{C}$
15:                     $\mathbf{C} \leftarrow \bigcup_{C \in \mathbf{C}} Ch_\mathcal{G}(C) \setminus \mathbf{P'}$
16:                 **end while**
17:                 $\text{Fr}[\mathbf{T}] \leftarrow (\mathbf{C} = \varnothing)$
18:             **end if**
19:             **if** $\text{Fr}[\mathbf{T}]$ **then**
20:                 $\mathbf{Z} \leftarrow \mathbf{Z} \cup \{X\}$
21:             **end if**
22:         **end if**
23:         $\mathbf{P} \leftarrow \mathbf{P} \cup \{X\}$
24:         $k \leftarrow k - 1$
25:     **end while**
26:     **return** $\mathbf{S} \setminus \mathbf{Z}$
27: **end procedure**

---

# Conclusions

- **Estimand-Agnostic** Causal Inference as a **practical approach** for do-SHAP.

- **FRA** as an efficient algorithm to significantly **speed up do-SHAP**.

# Questions?

- **Poster**: Wed 3 Dec, 11 a.m. — 2 p.m. PST

- **Email**: alvaro.parafita@bsc.es