# Procedural Fairness Through Decoupling Objectionable Data Generating Components

**Zeyu Tang[1], Jialu Wang[2], Yang Liu[2 4], Peter Spirtes[1], Kun Zhang[1 3]**

zeyutang@cmu.edu    {faldict, yangliu}@ucsc.edu    {ps7z@andrew., kunz1@}cmu.edu

[1]Carnegie Mellon University            [3]MBZUAI
[2]University of California, Santa Cruz   [4]ByteDance Research

Carnegie Mellon University

AI INSTITUTE FOR SOCIETAL DECISION MAKING

UNIVERSITY OF CALIFORNIA SANTA CRUZ

MBZUAI

ByteDance

# In a nutshell

- We reveal the frequently over-looked issue of ***disguised procedural unfairness***

- We propose a framework to address procedural fairness
  - **value instantiation rule**
  - **reference point**

- We determine the value of reference point per Rawls's ***pure procedural justice*** principles

# What is procedural fairness or justice?

- Three kinds of procedural justice (Rawls 1971;2001)

  - perfect procedural justice

  - imperfect procedural justice

  - pure procedural justice

# Principles of *pure procedural justice*

- **Requirement I: Fair Equality of Opportunity**
  - *The opportunity should be open and attainable, with the same prospects of success, for those who are at the same level of talent and ability, and have the same willingness to use them. Such equality of opportunity should not be influenced by arbitrary contingencies.*

- **Requirement II: The Difference Principle**
  - *The (social and economic) inequalities are to be arranged so that they are to the greatest benefit to the least advantaged members of the society.*

# Our framework (part 1) value instantiation rule

**Algorithm 1:** The Value Instantiation Rule for Local Causal Modules

**Input** : The $d_{\text{in}}(V; \mathcal{G})$-ary function $h_V(\text{Parents}(V); \hat{\theta}_V)$ modeling the causal mechanism between the node $V$ and its direct parents, where $d_{\text{in}}(V; \mathcal{G})$ is the the number of direct parents (in-degree) of $V$ in the graph $\mathcal{G}$. The configuration function $\text{ReferencePoint}(\cdot)$, which maps a directed edge corresponding to an objectionable component $\rho \in \mathcal{E}_{\text{Obj}}$ to a reference point (Definition 4.1) with the domain of value of the tail node of the edge.

**Output** : The derivation of the predicted outcome $\widehat{V}$ in the local causal module.

1 **If** *there is additional assumption on the functional form $\widetilde{h}_V(\cdot)$ and/or parameters $\widetilde{\theta}_V$* **Then**

2 $\quad \hat{\theta}_V \leftarrow \widetilde{\theta}_V, h_V \leftarrow \widetilde{h}_V$ ;   `// direct correction of the causal mechanism`

3 **Else**

4 $\quad$ **ForEach** *parent node $W_j$ in* $\text{Parents}(V) = (W_1, W_2, \ldots, W_{d_{\text{in}}(V;\mathcal{G})})$ **Do**

5 $\quad\quad$ **If** *the edge $\rho_j = (W_j, V) \in \mathcal{E}_{\text{Obj}}$, i.e.,* $W_j \to V$ *is an objectionable component* **Then**

6 $\quad\quad\quad$ $w_j$ gets the value $\text{ReferencePoint}(\rho_j)$, because $W_j = \text{Tail}(\rho_j)$;

7 $\quad\quad$ **Else If** *there is at least one ancestor nodes of $W_j$ was set to a reference point* **Then**

8 $\quad\quad\quad$ $w_j$ gets the value that $W_j$ would have taken as a downstream of its ancestor nodes, to which reference points, if any, have been assigned;

9 $\quad\quad$ **Else**

10 $\quad\quad\quad$ $w_j$ gets the value of variable $W_j$ for the record in the data set;

11 $\widehat{v} \leftarrow h_V(w_1, w_2, \ldots, w_{d_{\text{in}}(V;\mathcal{G})}; \hat{\theta}_V)$.

# Our framework (part 2) reference point values

**Algorithm 2:** Aggregating Local Causal Modules while Decoupling Objectionable Components

**Input** : The data set $\mathcal{D}$, the hypothesis class $\mathcal{H}$ and the parameter space $\Theta$, the causal graph $\mathcal{G} = (\mathbf{V}, \mathcal{E})$, the list of index $\mathcal{I}$ for all nodes $\mathbf{V}$. The set of edges $\mathcal{E}_{\text{Obj}}$ where each edge corresponds to an objectionable component. The ReferencePoint($\cdot$) configuration.

**Output** : The derivation of the predicted outcome $\widehat{Y}$ that *decouples* objectionable components from the data generating process, and *only* makes use of neutral components.

1 Sort the list of index $\mathcal{I}$ such that parent nodes, if any, appear before the node itself;

2 **ForEach** *node index* $i \in \mathcal{I}$ **Do**                      // learn model parameters

3     **If** *the number of direct parents of node* $V_i$*, i.e., the in-degree,* $d_{\text{in}}(V_i; \mathcal{G}) > 0$ **Then**

4         Fit model parameters in the local causal module between $V_i$ and its direct parent nodes Parents($V_i$), without any fairness constraint:
$$h_{V_i}, \hat{\theta}_{V_i} \leftarrow \underset{\theta \in \Theta, h \in \mathcal{H}}{\arg\min} \mathcal{L}_{V_i}\big(h(\text{Parents}(V_i); \theta), V_i; \mathcal{D}\big), \mathcal{L}_{V_i} \text{ is the loss function for } V_i;$$

5 According to the sorted list of node index $\mathcal{I}$, apply the value instantiation rule (Algorithm 1) to each local causal module in sequence, and then derive prediction $\widehat{Y}$

$$\widehat{y} = \underset{i \in \mathcal{I}}{\circ} \big(h_{V_i} \circ \text{ReferencePoint}\big)(\mathbf{z}; \mathcal{E}_{\text{Obj}}, \hat{\theta}_{V_1}, \ldots, \hat{\theta}_{V_{|\mathcal{I}|}})$$

# Thank you!

OUR PAPER

OUR CODE