

PC-Fairness: A Unified Framework for Measuring Causality-based Fairness



TL;DR

- o We define a unified fairness notion that covers most previous causality-based fairness notions, namely **Path-specific Counterfactual fairness (PC fairness)**.
- We develop a constrained optimization-based approach for tightly bounding any unidentifiable PC fairness.
- The experiments show the creativeness and effectiveness of the proposed method under any unidentifiable situation or combinations.

Backgrounds for Causal Inference

Structural Causal Model (SEM) and Causal Graph

- Notations
 - Exogenous Variables $U = \{U_X, U_Y, U_Z\}$, and Distribution P(U)

Table 1: Connection between previous fairness notions and PC fairness		
Description	References	Relating to PC fairness
Total causal effect	[10, 7]	$\mathbf{O} = \emptyset$ and $\pi = \Pi$
(System) Direct discrimination	[10, 3, 7]	$\mathbf{O} = \emptyset$ or $\{S\}$ and $\pi = \pi_d = \{S \to \hat{Y}\}$
(System) Indirect discrimination	[10, 3, 7]	$\mathbf{O} = \emptyset$ or $\{S\}$ and $\pi = \pi_i \subset \Pi$
Individual direct discrimination	[8]	$\mathbf{O} = \{S, \mathbf{X}\} \text{ and } \pi = \pi_d = \{S \to \hat{Y}\}$
Path-specific Individual discrimination	[1]	$\mathbf{O} = \{S, \mathbf{X}\}$ and $\pi \in \Pi$
Group direct discrimination	[9]	$\mathbf{O} = \mathbf{Q} = PA_Y \setminus \{S\} \text{ and } \pi = \pi_d = \{S \to \hat{Y}\}$
Counterfactual fairness	[2, 4, 5]	$\mathbf{O} = \{S, \mathbf{X}\} \text{ and } \pi = \Pi$
Counterfactual error rate	[6]	$\mathbf{O} = \{S, Y\}$ and $\pi = \pi_d$ or π_i

Note 1: please refer to our paper for the references.

Measuring Path-specific Counterfactual Fairness

Note 2: please refer to Section 4 for more technical details.



Causal Effects

- Causal effects are responses of the outcome of interest to an external intervention. ✤ Intervention is facilitated with *do*-operator.
 - Intervention do(x) is defined as the substitution of structural equation X = \bullet $f_X(u_X)$ with X = x.
 - The response of Y after intervention do(x) is denoted by Y_x . \bullet
- *** Total Causal Effect** of $X: x_0 \to x_1$ on Y = y:

 $TCE(s^{-}, s^{+}) = P(y_{s^{+}}) - P(y_{s^{-}})$

* Path-specific Effect of $X: x_0 \to x_1$ on Y = y through a path set $\pi = \{X \to Z \to Y\}$: $PE(x_1, x_0) = P(y_{x_1|\pi, x_0|\overline{\pi}}) - P(y_{x_0})$

Counterfactual Effect of $X: x_0 \to x_1$ on Y = y for a condition X = x', O = o': $CE(x_1, x_0) = P(y_{x_1} | x', o') - P(y_{x_0} | x', o')$

The key barrier in PC fairness is the unidentifiable issues. Examples of typical unidentifiable situations are shown below







Figure 2: The "bow graph".

Figure 3: The "kite graph". Figure 4: The "w graph".

• We formulate the bounding problem of PC fairness as a constrained problem by using "response-function variables".

Response-function Variables

- \clubsuit Response-function variables r are used to parameterize the causal model.
 - Categorize the unknown domain of **U** into limited number of equivalent regions, each of which is denoted by a value of a response-function variable.
 - Represent unknown functions **F** by limited number of mappings from parents to dependent variables. Each corresponds to one equivalent region of **U**.
 - As a result, all uncertainties in the causal model parameterized by P(r) U_V U_X

$$\begin{array}{c} & & \\ & &$$

Backgrounds for Causal Fairness Learning

Path-specific Fairness and Counterfactual Fairness

- Three types of path-specific effects:
- The green-path effect is considered as explainable.
- The red-path is considered as indirect discrimination.
- The purple-path is considered as direct discrimination.
- Counterfactual Fairness: the total effect for a group specified by the factual observation.
 - The blue box denotes the real/factual world.
 - The purple box denotes the counterfactual world.
 - The condition $O = \{A, B, C\}$ (in green).
 - The counterfactual quantity is $P(\widehat{Y}_s)$ (in yellow).

Path-specific Counterfactual Fairness

Path-specific Counterfactual Effect

♣ Path-specific Counterfactual Effect of $X: x_0 \rightarrow x_1$ on Y = y for a given condition $\boldsymbol{O} = \boldsymbol{O}$ through a pathset π is

$$PCE_{\pi}(x_1, x_0 | \boldsymbol{o}) = P(y_{x_1 | \pi, x_0 | \overline{\pi}} | \boldsymbol{o}) - P(y_{x_0} | \boldsymbol{o})$$

Path-specific Counterfactual Fairness (PC Fairness)



Bounding Path-specific Counterfactual Fairness

Based on response-function variables, finding the bounds of PC fairness is formulated as follows

min/max $P(\hat{y}_{s_1|\pi,s_0|\bar{\pi}}|\mathbf{o}) - P(\hat{y}_{s_0}|\mathbf{o}),$

s.t.
$$P(\mathbf{V}) = P(\mathcal{D}), \quad \sum P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \ge 0,$$



• Given a predictor \hat{Y} , a condition \boldsymbol{o} , a path set π ,

1) \hat{Y} achieves PC fairness if $PCE_{\pi}(s_1, s_0 | \boldsymbol{o}) = 0$;

2) \hat{Y} achieves τ -PC fairness if $|PCE_{\pi}(s_1, s_0 | \boldsymbol{o})| \leq \tau$.

where s_0, s_1 are two values of the sensitive attribute.

Thanks to the flexibility and capability of the path-specific counterfactual effect, the PC fairness generalize the previous causality-based fairness notions, as summarized in Table 1.

where the objective function is the path-specific counterfactual effect, P(V) is the parameterized distribution, which agrees with the observational distribution $P(\mathbf{D})$.

Acknowledgement



This work was supported in part by NSF 1646654, 1920920, and 1940093. Paper is available at **NIPS Proceedings**. Code is available at <u>http://tiny.cc/pc-fairness-code</u>.

Email: {yw009, lz006, xintaowu}@uark.edu; htong@illinois.edu

2019 Conference on Neural Information Processing Systems (NeurIPS 2019) Vancouver, Canada, Dec. 8-14, 2019