
CausalFairness: An Open Source Python Library for Causal Fairness Analysis

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 As machine learning (ML) systems are increasingly deployed in high-stakes do-
2 mains, the need for robust methods to assess fairness has become more critical.
3 While statistical fairness metrics are widely used due to their simplicity, they are
4 limited in their ability to explain why disparities occur, as they rely on associative
5 relationships in the data. In contrast, causal fairness metrics aim to uncover the
6 underlying data-generating mechanisms that lead to observed disparities, enabling
7 a deeper understanding of the influence of sensitive attributes and their proxies.
8 Despite their promise, causal fairness metrics have seen limited adoption due to
9 their technical and computational complexity. To address this gap, we present
10 CausalFairness, the first open-source Python package designed to compute a di-
11 verse set of causal fairness metrics at both the group and individual levels. The
12 metrics implemented are broadly applicable across classification and regression
13 tasks (with easy extensions for intersectional analysis) and were selected for their
14 significance in the fairness literature. We also demonstrate how standard statistical
15 fairness metrics can be decomposed into their causal components, providing a
16 complementary view of fairness grounded in causal reasoning.

17 1 Introduction

18 Statistical fairness metrics are easy to compute but only capture associations (i.e. conditional
19 probabilities), not causality — limiting their ability to identify whether observed statistical disparities
20 are truly caused by protected attributes or not. Causal fairness metrics, based on Structural Causal
21 Models (SCMs) [6], overcome this by attributing observed disparities to specific sources (protected
22 attributes, mediators or confounders). They also enable causal decompositions of statistical fairness
23 metrics, thus delivering deeper insights into fairness audits. Despite their value, causal metrics
24 are rarely used in practice due to technical challenges: they are harder to compute, require do-
25 interventions, and face identifiability constraints. Each causal fairness metric often needs a custom
26 architecture, and disagreement over the causal graph adds complexity. To address this, we introduce
27 **CausalFairness**—the first open-source Python package to implement generalizable algorithms for
28 key, established causal fairness metrics in the literature, including *Counterfactual Effects* [7]/[12],
29 *Counterfactual Equalized Odds* [11], and *Counterfactual Fairness* [3] (see Table 1 and Appendix
30 A.1 for an overview). The package is designed to work with minimal identifiability assumptions¹,
31 does not (necessarily) require a fully specified SCM and supports both group and individual-level
32 fairness metrics. We demonstrate CausalFairness on three datasets—*Adult Income*, *COMPAS*, and

¹This is not equivalent to us making the claim that we address the challenge of identifiability or choice of cause of graph. Instead, the paper focuses on implementing metrics for which identifiability constraints are not very strong and thus can be implemented for a wide variety of problems and domains without running into problems of identifiability; For causal model discovery, see packages like *CausalNex*, *DoWhy*, or *CausalML*.

33 *LSAC*—and provide code for replication.² The paper positions itself as being a novel contribution to
 34 the causal fairness literature by implementing novel computational algorithms for computing three
 35 existing causal fairness metrics in the *CausalFairness* package, thus filling a critical gap in enabling
 36 practical use of causal fairness metrics³. This work contributes to the counterfactual measurement
 37 branch of causal fairness literature [14] (see Appendix A.2 for a brief literature review).

38 2 Methodology

39 **Notation and Preliminaries** : Causal fairness is typically formalized using Structural Causal Models
 40 (SCMs)[6], where a Directed Acyclic Graph (DAG) represents observed variables (nodes) and their
 41 causal relationships (edges). Y is the true outcome, \hat{Y} the predicted outcome, and y the favorable
 42 outcome (e.g., $y = 1$ in the Adult Income dataset). A is the set of protected attributes (e.g., race,
 43 gender), X contains all features excluding A , and a_0, a_1 denote advantaged and disadvantaged groups.
 44 The causal effect of an intervention $do(X = x)$ is expressed via the counterfactual distribution
 45 $P(Y_x = y)$, where Y_x is the outcome had X been set to x . Often, $P(y|x) = P(Y = y|X = x)$ is
 46 used interchangeably. The *Standard Fairness Model*[12] (see Appendix A.3) outlines three causal
 47 pathways from A to Y and \hat{Y} : **Direct Path** ($A \rightarrow Y$): Captures direct discrimination (e.g., gender or
 48 race directly influencing income or recidivism), interpreted as *disparate treatment*; **Indirect Path**
 49 ($A \rightarrow M \rightarrow Y$): Effects mediated by variables like education or prior offenses, indicating *disparate*
 50 *impact*; **Spurious Path** ($A \leftarrow C \rightarrow Y$): Non-causal associations due to confounders (C), such as
 51 country of residence or age/gender, also contributing to *disparate impact*. Discrimination is assessed
 52 on the DAG using counterfactuals. By applying different *do*-interventions, specific paths can be
 53 isolated. The key idea is *ceteris paribus*, how does changing the protected attribute A affect Y or \hat{Y} ?

54 2.1 CausalFairness : Definitions and Computational Algorithms

55 The core of the *CausalFairness* package is the *CausalFairnessDecomposition* class (see Ta-
 56 ble 1), built on the standard fairness model [12]. It accepts X, Y, \hat{Y} , lists for A, M , optionally C ,
 57 derived from an SCM (algorithmically discovered or expert-curated) or a DAG, and a task-type flag
 58 (regression/classification). The class provides three main methods: `analyse_mean_difference`
 59 – group-level analysis using the Counterfactual Effects framework; `analyse_equalized_odds`
 60 – causal decomposition of error rates; `analyse_counterfactual_fairness` – individual-level
 61 fairness analysis. Each method compares the outcome (acceptance rate, error rates, or predicted
 62 outcome \hat{Y}_i) in two counterfactual worlds: one with observed A , and one with counterfactual A . The
 63 first two rely on (estimated) conditional probabilities using Gaussian Mixture Models for scalability;
 64 the third requires a DAG and fits a graphical causal model.

Table 1: Pseudo-Algorithm for Causal Fairness Metrics

A. Counterfactual Effects (Mean Difference)	B. Counterfactual Equalized Odds (EO)	C. Counterfactual Fairness
Inputs: D, A, M, C, a_0, a_1, y	Inputs: $D, A, C, a_0, a_1, y, \hat{f}$	Inputs: $A, M, C, a_0, a_1, \text{DAG}$
<ol style="list-style-type: none"> For each $(m, c) \in D$: <ul style="list-style-type: none"> Compute: $\mathbb{E}(Y = y a_0, m, c)$ Compute: $\mathbb{E}(Y = y a_1, m, c)$ Estimate via GMM: $P(m a_0, c), P(m a_1, c)$ $P(c a_0), P(c a_1)$ Combine expectations and probabilities to compute the counterfactual effects 	<ol style="list-style-type: none"> For each $c_j \in D$: <ul style="list-style-type: none"> Predict: $\hat{f}(c_j, a_0), \hat{f}(c_j, a_1)$ Obtain: $P(y_{a_0, c_j}), P(y_{a_1, c_j})$ Estimate via GMM: $P(c a_0), P(c a_1)$ Combine predictions and probabilities to compute the Cft-EO 	<ol style="list-style-type: none"> Fit SCM using DAG and dataset D For each individual $i \in D$: <ul style="list-style-type: none"> Get A_{obs} (observed) and A_{cft} (counterfactual) Sample from SCM under: $do(A = A_{obs}) \Rightarrow D_{obs}$ $do(A = A_{cft}) \Rightarrow D_{cft}$ Predict: $\hat{f}(D_{obs}), \hat{f}(D_{cft})$ Check: $Y_{obs} \neq Y_{cft}$

65 2.2 Counterfactual Effects : How does the protected attribute affect the predicted outcome? 66 Calculating Disparate Treatment, Disparate Impact and Explaining the Causal 67 Mechanism Behind Observed Statistical Parity

68 Counterfactual effects [13] is a family of three causal measures of discrimination related to statistical
 69 parity, namely: **Counterfactual Direct Effect (Ctf-DE)**: measures direct discrimination along

²Code will be released; cleared on 10/9/25

³The implementation has internationally been optimized so as to not require any special hardware requirements

70 $A \rightarrow \hat{Y}$ by holding M and C constant, isolating the effect of A on \hat{Y} . [7] define symmetric Ctf-DE
71 as: $DE_a^{\text{sym}}(y|a) = \frac{1}{2} (DE_{a_0, a_1}(y|a) - DE_{a_1, a_0}(y|a))$ i.e. the net treatment, which is the difference
72 between the positive and negative effect of protected group membership. Direct discrimination
73 exists if $DE_a^{\text{sym}}(y|a) > 0$ i.e. the negative effect is greater than the positive effect; **Counterfactual**
74 **Indirect Effect (Ctf-IE)**: measures indirect discrimination along $A \rightarrow M \rightarrow \hat{Y}$ by holding A
75 and C fixed, capturing the effect of M on \hat{Y} . [7] define symmetric Ctf-IE as: $IE_a^{\text{sym}}(y|a) =$
76 $\frac{1}{2} (IE_{a_0, a_1}(y|a) - IE_{a_1, a_0}(y|a))$. Indirect discrimination exists if $IE_a^{\text{sym}}(y|a) > 0$; **Counterfactual**
77 **Spurious Effect (Ctf-SE)**: measures confounding impact along $A \leftarrow C \rightarrow \hat{Y}$, varying C while
78 fixing A and M ⁴. Its is given by: $SE_{a_0, a_1}(y) = P(y_{a_0}|a_1) - P(y|a_0)$. As shown in [7] **disparate**
79 **treatment** (direct discrimination) exists if the symmetric difference due to A , $DE_a^{\text{sym}}(y|a)$, differs
80 from zero. **Disparate impact** (indirect discrimination) exists if either the symmetric indirect effect,
81 $IE_a^{\text{sym}}(y|a)$, or the spurious effect, $SE_{a_0, a_1}(y|a)$, is non-zero. Statistical disparity decomposes as:
82 Mean Difference _{a_0, a_1} (y) = $DE_a^{\text{sym}}(y|a) + IE_a^{\text{sym}}(y|a) + SE_{a_0, a_1}(y|a)$. Without confounders, Ctf-DE
83 and Ctf-IE reduce to the natural direct and indirect effects, respectively.

84 **Algorithmic Procedure**: [12] provide empirical formulas to estimate these effects from observed
85 data using conditional probabilities, avoiding the need for a fully specified SCM, thus aiding ease of
86 application (see Appendix A.4 for the empirical formulations). For each combination $m \in M$ and
87 each combination $c \in C$, we get a subset of D defined by (m, c) . For each subset (m, c) : we calculate
88 condition probability / expectation of the outcome of interest Y_y for a_0 and a_1 i.e. $\mathbb{E}(y|a_0, m, c)$
89 and $\mathbb{E}(y|a_1, m, c)$ respectively. These are the differences in the realisation of Y_y when M and C
90 are the same but A is different. Then for each m , we get the probability of m given c for a_0 and
91 a_1 i.e. $P(m|a_0, c)$ and $P(m|a_1, c)$ i.e. the difference in probability of m when c is the same but A
92 is varied. Then lastly, for each c , we get it's probability for a_0 and a_1 i.e. $P(c|a_0)$ and $P(c|a_1)$
93 respectively. Each of these computed quantities is then combined as per (1) to (11) (see Appendix
94 A.4) to get $DE_a^{\text{sym}}(y|a)$, $IE_a^{\text{sym}}(y|a)$, $SE_{a_0, a_1}(y|a)$ and Mean Difference _{a_0, a_1} (y).

95 2.3 Counterfactual Equalised Odds (Cft-EO) : How does the protected attribute affect the 96 model error rate?

97 Like counterfactual effects, [11] use the standard fairness model to define three causal counterfactual
98 metrics based on equalized error rates: Counterfactual Direct Error Rate ($ER_{a_0, a_1}^d(\hat{y} | a, y) =$
99 $P(\hat{y}_{a_1, y}, (\hat{P}A \setminus X)_{a_0, y} | a, y) - P(\hat{y}_{a_0, y} | a, y)$), Counterfactual Indirect Error Rate ($ER_{a_0, a_1}^i(\hat{y} |$
100 $a, y) = P(\hat{y}_{a_0, y}, (\hat{P}A \setminus X)_{a_1, y} | a, y) - P(\hat{y}_{a_0, y} | a, y)$), and Counterfactual Spurious Error Rate
101 ($ER_{a_0, a_1}^s(\hat{y} | y) = P(\hat{y}_{a_0, y} | a_1, y) - P(\hat{y}_{a_0, y} | a_0, y)$). These measure how error rates would
102 change if the disadvantaged (advantaged) group had the identity, mediators, or confounders of the
103 advantaged (disadvantaged) group. They conclude error rates driven by ER^d indicate bias, while
104 those due to ER^i or ER^s are not discriminatory.⁵ Using these three counterfactual error metrics,
105 [11] show that equalized odds can be broken down into direct, indirect and spurious components as
106 follows: $ER_{a_0, a_1}(\hat{y} | y) = ER_{a_0, a_1}^d(\hat{y} | a_0, y) - ER_{a_1, a_0}^i(\hat{y} | a_0, y) - ER_{a_1, a_0}^s(\hat{y} | y)$.

107 **Algorithmic Procedure**: Unlike counterfactual effects, these metrics face a key limitation: Cft-EO
108 cannot reliably estimate direct, indirect, and spurious effects in the presence of mediators due to
109 identifiability issues from conditioning on both Y and \hat{Y} (whereas Counterfactual Effects condition
110 only on \hat{Y}). The common fix is excluding M from features, using only protected attributes A and
111 confounders C , enabling accurate estimation of ER^d and ER^s . This solution allows for the accurate
112 identification and estimation of Counterfactual Direct Error Rate and Counterfactual Spurious Error
113 Rate. However, this remedial strategy is undesirable in real world applications because the exclusion
114 of M is likely to negatively impact the predictive performance of the model. Thus, to ensure that the
115 metric is used correctly, we remove M from consideration, and modify [11] to estimate the simplified
116 counterfactual error rates as follows: $ER_{a_0, a_1}^d(\hat{y} | a, y) = \sum_c (P(\hat{y}_{a_1, c}) - P(\hat{y}_{a_0, c})) P(c | a, y)$ and
117 $ER_{a_0, a_1}^s(\hat{y} | y) = \sum_c P(\hat{y}_{a_1, c}) (P(c | a_1, y) - P(c | a_0, y))$. To compute these from the observed
118 data D , we use the following procedure: For each $c \in C$ we use the fitted estimator \hat{f} as $\hat{f}(a_1, c)$

⁴Ctf-SE has no symmetric form since confounders C are non-descendants of A and remain unchanged under interventions.

⁵Definitions overlap with Ctf-DE, Ctf-IE, and Ctf-SE from Section 2.1, differing by focusing on error rates instead of mean difference.

119 and $\hat{f}(a_0, c)$ to get $P(\hat{y}_{a_1, c})$ and $P(\hat{y}_{a_0, c})$ respectively. These quantities are the differences in the
 120 realisation of Y_y when C is the same but A is different. Then for each c , we get it’s probability for a_0
 121 and a_1 i.e. $P(c|a_0)$ and $P(c|a_1)$ respectively. Each of these computed quantities is then combined as
 122 per (12) to (14) (see Appendix A.5) to get $ER_{a_0, a_1}^d(\hat{y} | a, y)$, $ER_{a_0, a_1}^s(\hat{y} | y)$ and $ER_{a_0, a_1}(\hat{y} | y)$.

123 2.4 Counterfactual Fairness

124 Unlike Counterfactual Equalised Odds and Counterfactual Effects, Counterfactual fairness[3] is
 125 an individual level causal fairness metric which is achieved if changing an individual i ’s protected
 126 attributes doesn’t change the predicted outcome \hat{Y}_i i.e. $P(\hat{Y}_{A \leftarrow a}(U) = y | X = x, A = a) =$
 127 $P(\hat{Y}_{A \leftarrow a'}(U) = y | X = x, A = a)$

128 **Algorithmic Procedure:** To empirically test for counterfactual fairness first, a fully specified SCM
 129 must be defined using a specified Directed Acyclic Graph (DAG) and dataset D . For each instance i in
 130 D , we retrieve the observed value A_{obs} and compute a counterfactual value A_{cf} . We then generate sam-
 131 ples from the SCM through two do-interventions using the standard "abduction,action,prediction"[6]
 132 procedure. First, we perform an intervention $\text{do}(A = A_{\text{obs}})$ to produce samples for the observed
 133 state $D_{\text{do=observed}}$. Second, we perform an intervention $\text{do}(A = A_{\text{counterfactual}})$ to produce samples for
 134 the counterfactual state $D_{\text{do=cf}}$. Using these samples, we fit functions $\hat{f}(D_{\text{do=obs}})$ and $\hat{f}(D_{\text{do=cf}})$
 135 to obtain predicted outcomes $Y_{D_{\text{do=obs}}}$ and $Y_{D_{\text{do=cf}}}$. To assess counterfactual fairness, we compare the
 136 observed and counterfactual predictions. If $Y_{D_{\text{do=obs}}} \neq Y_{D_{\text{do=cf}}}$, then the prediction function \hat{f} is not
 137 counterfactually fair.

138 **2.5 Scalability :** To address scalability, the algorithms include the following optimizations: **A. GMMs**
 139 **for Conditional Probability Estimation.** For Counterfactual Effects, estimating probabilities via
 140 conditional expectations on a 50,000-sample dataset takes approximately 2 seconds; using GMMs
 141 reduces this to 300–500 ms, depending on the number of features. For Counterfactual Equalized
 142 Odds, estimation takes around 1 second, with GMMs reducing latency to 300–500 ms, depending
 143 on feature count and model complexity. **B. Parallelization of Interventions.** For Counterfactual
 144 Fairness, computing for a single instance takes approximately 10 seconds without parallelization,
 145 and about 1 second with it. Actual times vary with the number of features, interventions, and the
 146 predictive model used.

147 3 Results: Application of CausalFairness to Benchmark Datasets

Dataset	Protected Attribute	Mean Difference	FNR	FPR	$DE_a^{\text{sym}}(y a)$	$IE_a^{\text{sym}}(y a)$	$SE_{a_0, a_1}(y a)$	ER^d	ER^s	ER^*	Counterfactual Fairness
Adult Income	Gender	0.203	0.410	-0.104	0.165	0.039	0.000	0.000	0.000	0.000	-0.031
Adult Income	Intersectional	0.221	0.445	-0.115	0.152	0.069	0.000	0.000	0.000	0.000	-0.068
COMPAS	Race (Black)	0.326	-0.310 (-42)	-0.253 (-0.41)	0.154	0.071	0.101	FPR: 0.297, FNR: 0.265	0	FPR: 0.113, FNR: 0.162	0.055
COMPAS	Intersectional	0.620	-0.620	-0.518	0.513	0.081	0.027	-	-	-	0.640
LSAC	Race (Black)	0.978	-	-	0.554	0.429	0.000	-	-	-	0.001
LSAC	Intersectional	0.990	-	-	0.531	0.458	0.000	-	-	-	-0.007

Table 2: Statistical and Causal Fairness Metrics

148 **3.1 Adult Income Dataset:** We fit a logistic regression using the structure and features in Appendix
 149 A.3, Fig.1.a⁶ to predict $P(\text{Income} > \$50k)$. **Counterfactual Effects:** On average, women are 20.3%
 150 less likely than men to be predicted as earning above \$50k. Most of this disparity (16.5%) is due to
 151 *disparate treatment* ($DE_a^{\text{sym}}(y | a)$), meaning that simply having the social identify of a woman lowers
 152 $P(\text{Income} > \$50k)$. The remaining 3.9% is due to *disparate impact* ($IE_a^{\text{sym}}(y | a)$) via M (years of
 153 education and occupation typical of women); **Counterfactual Equalized Odds:** When refitting the
 154 model without M , the predictor always outputs 0, making the equalized odds and its decomposition
 155 uninformative. We thus cannot determine whether observed disparities stem from disparate treatment
 156 or impact; **Counterfactual Fairness:** As shown in Appendix Fig.2.A, the observed and counterfactual
 157 distributions do not overlap—changing a woman’s gender to male shifts \hat{Y} rightward, increasing
 158 $P(\text{Income} > \$50k)$. Hence, the fitted logistic regression is not counterfactually fair.

159 **3.2 COMPAS Recidivism Dataset:** We fit a logistic regression using the structure and features in
 160 Appendix A.3, Fig.1.b **Counterfactual Effects:** Black individuals are 32.6% more likely than white
 161 individuals to be predicted as high-risk for recidivism. Most of this is due to *disparate treatment*

⁶Country of residence is included as a cause of gender to replicate Zhang and Bareinboim, 2018

162 (15.4%), meaning that being Black alone increases $P(\text{Recidivism})$. *Disparate impact* comes from
163 both M and C : confounders like age and gender raise risk by $\sim 10\%$ (spurious effect), and M
164 contributes an additional 7.1%. **Counterfactual Equalized Odds:** Excluding M does not make the
165 model naive, though it increases error rates. Decomposing FPR/FNR shows most of the disparity
166 stems from direct discrimination: 29.7% of the 41% FPR and 26.5% of the 42% FNR. **Counterfactual**
167 **Fairness:** The COMPAS model is not counterfactually fair (see Appendix Fig. 2.C)—changing race
168 from Black to white shifts the distribution of \hat{Y} leftward, decreasing $P(\text{Recidivism})$.

169 **3.3 Law School Admission Council (LSAC) Dataset:** We fit a Random Forest regressor on GPA,
170 LSAT, Race, and Gender to predict average grade, where $A = (\text{Race}, \text{Gender})$, $M = (\text{GPA}, \text{LSAT})$,
171 and no C . **Counterfactual Effects:** The predicted average grade for the white subgroup is 0.978
172 higher than for the Black subgroup. Both *disparate treatment* and *disparate impact* are significant,
173 with most of the gap (0.55) due to direct discrimination. Since there are no confounders, the direct
174 and indirect effects correspond to the natural direct and indirect effects. **Counterfactual Equalized**
175 **Odds:** Not applicable since this is a regression task. **Counterfactual Fairness:** The model is
176 counterfactually fair, illustrating that fairness can differ at the individual vs. group level.⁷

177 **3.4 Intersectional Causal Fairness:** This package supports intersectional analysis to detect potential
178 “double” or higher-order discrimination. **Adult Income:** Black women are 22.1% less likely than
179 white men to have $P(\text{Income} > \$50k)$ —2% more than the non-intersectional gender gap—with direct
180 effect being the largest contributor. The model is also more counterfactually unfair: changing a Black
181 woman’s identity to a white man increases $P(\text{Income} > \$50k)$ by 6% (vs. 2% non-intersectionally)
182 (see Appendix Fig. 2.B). **COMPAS:** The mean difference in $P(\text{Recidivism})$ between Black men
183 and white women (60%) exceeds the non-intersectional racial gap, with direct discrimination as the
184 main driver. Counterfactual unfairness also rises by $\sim 11\%$: changing a Black man’s identity to a
185 white woman lowers predicted recidivism by 64% (vs. 55%) (see Appendix Fig. 2.D). **LSAC:** The
186 mean difference between Black women and white men is slightly higher than the non-intersectional
187 comparison (0.99 vs. 0.978), again mainly due to direct discrimination. As in the non-intersectional
188 case, the model remains counterfactually fair.

189 4 Limitations of CausalFairness

190 Generally, 1) deciding the right causal model from competing models of bias or achieving causal
191 fairness simultaneously across multiple competing models remains an active area of research and 2)
192 defining a hypothetical intervention on protected attributes remains a fraught process. The example
193 application to benchmark datasets highlighted how 1) lack of identifiability can limit analysis [4]
194 and 2) lack of methods for falsifying DAGs in the presence of competing causal models can lead
195 to disagreements about the validity of the conclusions. For example, in the Adult Income dataset,
196 identifiability issues prevented the causal decomposition of equalized odds. For counterfactual effects,
197 the DAG must be Markovian; otherwise, counterfactual probabilities cannot be empirically estimated⁸.
198 Extending the three discussed metrics to path-specific discrimination [6] is also limited by stricter
199 identifiability constraints. Hence, causal fairness metrics should be applied cautiously,

200 5 Conclusion

201 This paper introduced CausalFairness - the first open source generalizable implementation for
202 calculating key causal fairness metrics and applied it to 3 fairness benchmarking datasets. The
203 application to benchmark datasets demonstrated how CausalFairness provides practitioners with the
204 actionable insight - for example, at the very least the Adult Income model must eliminate at least
205 16.5% difference in statistical parity, while the COMPAS model needs to address 15.4% disparity in
206 statistical parity and 29.7-26.5% in error rates (all of which can be attributed to direct discrimination),
207 but this varies across intersectionally. Future work will expand the metrics available and extend the
208 package to include methods for bias reduction in the causal fairness literature (see Appendix A.7).

⁷Our experiments also show that “fairness through unawareness”—excluding A from training—can worsen fairness. For example, excluding Race in the LSAC dataset leads to a counterfactual fairness score of -0.50, meaning changing a Black individual’s race to white increases the predicted average grade by 0.50.

⁸In the presence of unobserved confounding, counterfactual effects may be estimated using counterfactual randomization ?, which is not implemented here.

209 **References**

- 210 [1] Barocas, S., Hardt, M. & Narayanan, A., Fairness and Machine Learning. Fairness and machine
211 learning. Available at: <https://fairmlbook.org/> [Accessed September 17, 2022]. Castelnovo, A. et al.,
212 2022. A clarification of the nuances in the fairness metrics landscape. Nature News. Available at:
213 <https://www.nature.com/articles/s41598-022-07939-1> [Accessed September 17, 2022].
- 214 [2] Kilbertus, N. et al., 2018. Avoiding discrimination through causal reasoning. arXiv.org. Available at:
215 <https://arxiv.org/abs/1706.02744> [Accessed September 17, 2022].
- 216 [3] Kusner, M.J. et al., 2018. Counterfactual fairness. arXiv.org. Available at: <https://arxiv.org/abs/1703.06856>
217 [Accessed September 17, 2022].
- 218 [4] Makhoulouf, K., Zhioua, S. & Palamidessi, C., 2022. Survey on causal-based machine learning fairness
219 notions. arXiv.org. Available at: <https://arxiv.org/abs/2010.09553> [Accessed September 17, 2022].
- 220 [5] PARK, K.E.V.I.N.A. & QUERCIA, R.O.B.E.R.T.O.G., Who lends beyond the Red Line? - university of
221 Pennsylvania. Available at: https://penniuir.upenn.edu/uploads/media/Park_Quercia.pdf [Accessed September
222 16, 2022].
- 223 [6] Pearl, J., Causality, 2nd edition, 2009. Available at: <http://bayes.cs.ucla.edu/BOOK-2K/> [Accessed September
224 17, 2022].
- 225 [7] Plecko, D. & Bareinboim, E., Causal fairness analysis. Available at: <https://causalai.net/r90.pdf> [Accessed
226 September 16, 2022].
- 227 [8] Russell, C. et al., 2017. When worlds collide: Integrating different counterfactual as-
228 sumptions in fairness. Advances in Neural Information Processing Systems. Available at:
229 <https://papers.nips.cc/paper/2017/hash/1271a7029c9df08643b631b02cf9e116-Abstract.html> [Accessed Septem-
230 ber 17, 2022].
- 231 [9] Spielkamp, M., 2020. Inspecting algorithms for bias. MIT Technology Review. Available at:
232 <https://www.technologyreview.com/2017/06/12/105804/inspecting-algorithms-for-bias/> [Accessed September
233 17, 2022].
- 234 [10] Williams, D.R., Priest, N. & Anderson, N.B., 2016. Understanding associations among
235 race, socioeconomic status, and Health: Patterns and prospects. Health psychology : official
236 journal of the Division of Health Psychology, American Psychological Association. Available at:
237 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4817358/> [Accessed September 17, 2022].
- 238 [11] Zhang, J. & Bareinboim, E., 1970. Equality of opportunity in classification: A
239 Causal approach. Advances in Neural Information Processing Systems. Available at:
240 <https://proceedings.neurips.cc/paper/2018/hash/ff1418e8cc993fe8abcfe3ce2003e5c5-Abstract.html> [Ac-
241 cessed September 17, 2022].
- 242 [12] Zhang, J. & Bareinboim, E., 2018. Fairness in decision-making - the causal explanation formula: Proceed-
243 ings of the Thirty-second AAAI Conference on Artificial Intelligence and thirtieth innovative applications of
244 Artificial Intelligence Conference and eighth AAAI symposium on educational advances in artificial intelligence.
245 Guide Proceedings. Available at: <https://dl.acm.org/doi/abs/10.5555/3504035.3504283> [Accessed September 17,
246 2022].
- 247 [13] Zhang, L., Wu, Y. & Wu, X., 2016. A causal framework for discovering and removing direct and indirect
248 discrimination. arXiv.org. Available at: <https://arxiv.org/abs/1611.07509> [Accessed September 17, 2022].
- 249 [14] Kusner, M.J. & Loftus, J.R., 2020. The Long Road to fairer algorithms. Nature News. Available at:
250 <https://www.nature.com/articles/d41586-020-00274-3> [Accessed September 17, 2022].

251 **A Appendix**

252 **A.1 Overview of Metrics Implemented in CausalFairness**

Table 3: Overview of Metrics Implemented in CausalFairness

Metric	Query Addressed	Level	Supported Metric Decomposition	Counterfactual Estimation Procedure
Counterfactual Effects (for Statistical Parity)	What would the disadvantaged (advantaged) group's acceptance rate be if they had the identity (A), mediating characteristics (M), or confounding characteristics (C) of the advantaged (disadvantaged) group?	Aggregate	Direct, Indirect, Spurious	Conditional probabilities (computed or estimated using GMM)
Counterfactual Equalized Odds	What would the disadvantaged (advantaged) group's error rate be if they had the identity (A), mediating characteristics (M), or confounding characteristics (C) of the advantaged (disadvantaged) group?	Aggregate	Direct, Spurious	Conditional probabilities (computed or estimated using GMM)
Counterfactual Fairness	What would the disadvantaged (advantaged) individual's predicted Y be if they had the identity (A), mediating characteristics (M), and confounding characteristics (C) of the advantaged (disadvantaged) group?	Individual	N/A	Predictions from functional relationships of fitted SCM

253 **A.2 A Brief Literature Review**

254 There has been considerable interest in the use of causal mechanisms to better understand black-box machine
 255 learning systems, and literature on causal fairness situates itself within the same. The causal fairness literature
 256 has three primary approaches for aiding algorithmic fairness assessment [14]: 1) Counterfactual measurement:
 257 aids the answer of what if cause-effect questions without running randomized control trials. For instance, *ceritus*
 258 *paribus*, if a woman's gender was changed to male, would her expected income be higher?; 2) Sensitivity analysis:
 259 how sensitive a model is to latent / confounding variables (which is often the status of protected attributes). For
 260 instance, sensitivity analysis can be used to "explore how sensitive our estimate of the causal link between legal
 261 representation and guilty verdict is to different levels of jury racism" and give recommendations for altering jury
 262 selection to minimize bias [14] and 3) Impact evaluation: to measure the long term consequences of automated
 263 decision making systems through the use of interventions. Following the principle of what gets measured gets
 264 managed, we recognize that causal identification of discrimination is crucial before moving on to remedial
 265 actions and impact analysis. Thus, this paper - and package - focus on addressing the gap in practical, broad
 266 adoption of causal (counterfactual) fairness metrics by providing implementations of [12]/[7], [11] and [3].

267 **A.3 Standard Fairness Model: Examples**

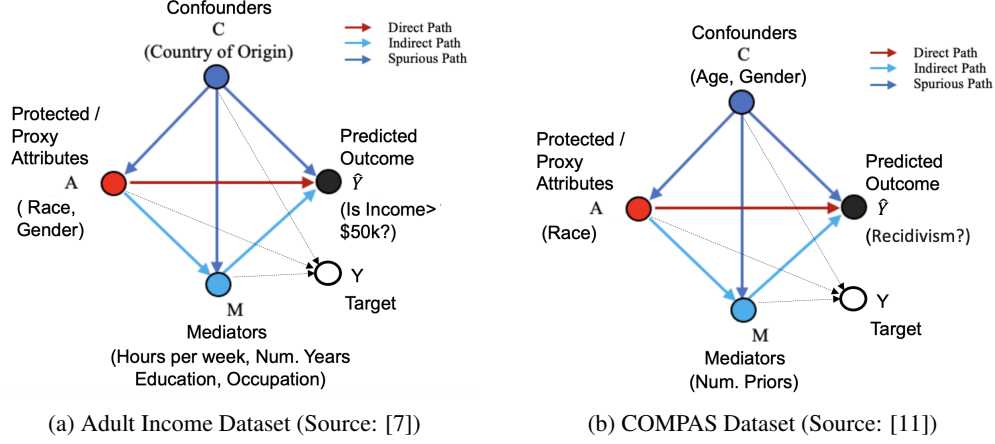


Figure 1: Standard Fairness Model

268 A.4 Counterfactual effects: Empirical Formulations

269 Given that protected group membership can have positive and negative impacts, the direct disadvantage due to
 270 protected group membership is given by:

$$DE_{a_0, a_1}(y|a) = P(y_{a_1}, m_{a_0}, c_a|a) - P(y_{a_0}, m_{a_0}, c_a|a) \quad (1)$$

271 while direct advantage is given by

$$DE_{a_1, a_0}(y|a) = P(y_{a_0}, m_{a_1}, c_a|a) - P(y_{a_1}, m_{a_1}, c_a|a) \quad (2)$$

272 The disadvantage due to the impact of A mediated through M is given by:

$$IE_{a_0, a_1}(y|a) = P(y_{a_0}, m_{a_1}, c_a|a) - P(y_{a_0}, m_{a_0}, c_a|a) \quad (3)$$

273 while the advantage is given by

$$IE_{a_1, a_0}(y|a) = P(y_{a_1}, m_{a_0}, c_a|a) - P(y_{a_1}, m_{a_1}, c_a|a) \quad (4)$$

274 Lastly, spurious effect if given by:

$$SE_{a_0, a_1}(y) = P(y_{a_0}|a_1) - P(y|a_0) \quad (5)$$

275 The corresponding empirical formulations for (1) to (5) are as follows:

$$DE_{a_0, a_1}(y|a) = \sum_{c, m} (\mathbb{E}(y|a_1, m, c) - \mathbb{E}(y|a_0, m, c)) P(m|a_0, c) P(c|a) \quad (6)$$

$$DE_{a_1, a_0}(y|a) = \sum_{c, m} (\mathbb{E}(y|a_0, m, c) - \mathbb{E}(y|a_1, m, c)) P(m|a_1, c) P(c|a) \quad (7)$$

$$IE_{a_0, a_1}(y|a) = \sum_{c, m} \mathbb{E}(y|a_0, m, c) (P(m|a_1, c) - P(m|a_0, c)) P(c|a) \quad (8)$$

$$IE_{a_1, a_0}(y|a) = \sum_{c, m} \mathbb{E}(y|a_1, m, c) (P(m|a_0, c) - P(m|a_1, c)) P(c|a) \quad (9)$$

$$SE_{a_0, a_1}(y|a) = \sum_{c, m} \mathbb{E}(y|a_0, m, c) P(m|a_0, c) (P(c|a_1) - P(c|a_0)) \quad (10)$$

$$\text{Mean Difference}_{a_0, a_1}(y) = DE_a^{\text{sym}}(y|a) + IE_a^{\text{sym}}(y|a) + SE_{a_0, a_1}(y|a) \quad (11)$$

276 **A.5 Counterfactual Equalized Odds: Empirical Formulations**

$$ER_{a_0, a_1}^d(\hat{y} | a, y) = \sum_c (P(\hat{y}_{a_1, c}) - P(\hat{y}_{a_0, c})) P(c | a, y) \quad (12)$$

$$ER_{a_0, a_1}^s(\hat{y} | y) = \sum_c P(\hat{y}_{a_1, c}) (P(c | a_1, y) - P(c | a_0, y)) \quad (13)$$

277 Using these two counterfactual error metrics, equalized odds can be broken down into direct and spurious
 278 components:

$$ER_{a_0, a_1}(\hat{y} | y) = ER_{a_0, a_1}^d(\hat{y} | a_0, y) - ER_{a_1, a_0}^s(\hat{y} | y) \quad (14)$$

279 **A.6 Counterfactual Fairness Plots**

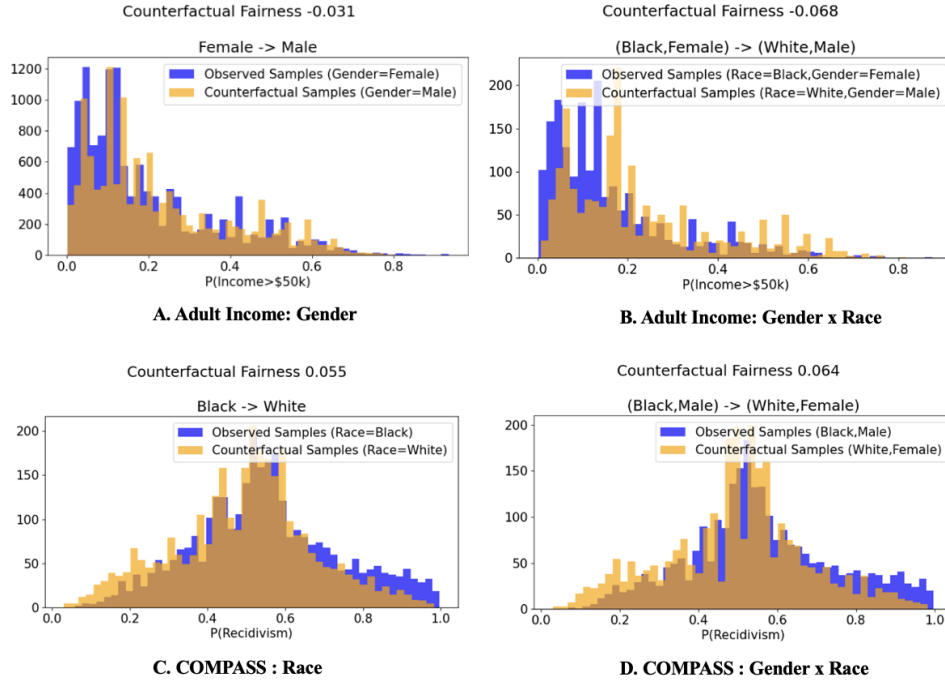


Figure 2: Counterfactual Fairness

280 A.7 Bias Reduction with CausalFairness and Upcoming Future Work

281 But now that the causal bias has been detected, what remedial steps can practitioners take to achieve causal
 282 fairness? The easiest way to achieve causal fairness is to train an estimator using only the observable non-
 283 descendants of A [3]. However, as most observable features are likely to be descendants of A, this strategy
 284 is unsuitable.⁹ Latent variables which are non-descendants of A but affect X and Y can be used to train
 285 counterfactually fair estimators [3]. Counterfactual effects [12] or counterfactual error rates [11] can be used
 286 for feature and sample selection to minimize direct, indirect and spurious discrimination. Using counterfactual
 287 fairness, a multi-world causal fairness penalty can be created to achieve counterfactual fairness under competing
 288 SCMs [8]. While addressing causal bias correction algorithms is out of scope for the current paper, this is an
 289 active area of research in the causal fairness literature which we aim to incorporate into forthcoming versions of
 290 the package along side sensitivity and impact analysis.

⁹Linear regression which includes the protected attributes is guaranteed to be counterfactually fair [3]

291 **NeurIPS Paper Checklist**

292 **1. Claims**

293 Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s
294 contributions and scope?

295 Answer: [Yes]

296 Justification: The abstract states that a novel causal fairness package has been introduced and the paper
297 elaborates on that.

298 Guidelines:

- 299 • The answer NA means that the abstract and introduction do not include the claims made in the
300 paper.
- 301 • The abstract and/or introduction should clearly state the claims made, including the contributions
302 made in the paper and important assumptions and limitations. A No or NA answer to this
303 question will not be perceived well by the reviewers.
- 304 • The claims made should match theoretical and experimental results, and reflect how much the
305 results can be expected to generalize to other settings.
- 306 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not
307 attained by the paper.

308 **2. Limitations**

309 Question: Does the paper discuss the limitations of the work performed by the authors?

310 Answer: [Yes]

311 Justification: The limitation section and the conclusion discuss the limitations of causal fairness
312 metrics

313 Guidelines:

- 314 • The answer NA means that the paper has no limitation while the answer No means that the paper
315 has limitations, but those are not discussed in the paper.
- 316 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 317 • The paper should point out any strong assumptions and how robust the results are to violations of
318 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,
319 asymptotic approximations only holding locally). The authors should reflect on how these
320 assumptions might be violated in practice and what the implications would be.
- 321 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested
322 on a few datasets or with a few runs. In general, empirical results often depend on implicit
323 assumptions, which should be articulated.
- 324 • The authors should reflect on the factors that influence the performance of the approach. For
325 example, a facial recognition algorithm may perform poorly when image resolution is low or
326 images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide
327 closed captions for online lectures because it fails to handle technical jargon.
- 328 • The authors should discuss the computational efficiency of the proposed algorithms and how
329 they scale with dataset size.
- 330 • If applicable, the authors should discuss possible limitations of their approach to address problems
331 of privacy and fairness.
- 332 • While the authors might fear that complete honesty about limitations might be used by reviewers
333 as grounds for rejection, a worse outcome might be that reviewers discover limitations that
334 aren’t acknowledged in the paper. The authors should use their best judgment and recognize
335 that individual actions in favor of transparency play an important role in developing norms that
336 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize
337 honesty concerning limitations.

338 **3. Theory assumptions and proofs**

339 Question: For each theoretical result, does the paper provide the full set of assumptions and a complete
340 (and correct) proof?

341 Answer: [NA]

342 Justification: The paper implements algorithms whose theoretical grantees have been proven in their
343 parent papers.

344 Guidelines:

- 345 • The answer NA means that the paper does not include theoretical results.

- 346 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 347 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 348 • The proofs can either appear in the main paper or the supplemental material, but if they appear in
- 349 the supplemental material, the authors are encouraged to provide a short proof sketch to provide
- 350 intuition.
- 351 • Inversely, any informal proof provided in the core of the paper should be complemented by
- 352 formal proofs provided in appendix or supplemental material.
- 353 • Theorems and Lemmas that the proof relies upon should be properly referenced.

354 4. Experimental result reproducibility

355 Question: Does the paper fully disclose all the information needed to reproduce the main experimental
356 results of the paper to the extent that it affects the main claims and/or conclusions of the paper
357 (regardless of whether the code and data are provided or not)?

358 Answer: [Yes]

359 Justification: The paper describes the algorithms and datasets required to reproduce the same.

360 Guidelines:

- 361 • The answer NA means that the paper does not include experiments.
- 362 • If the paper includes experiments, a No answer to this question will not be perceived well by the
- 363 reviewers: Making the paper reproducible is important, regardless of whether the code and data
- 364 are provided or not.
- 365 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make
- 366 their results reproducible or verifiable.
- 367 • Depending on the contribution, reproducibility can be accomplished in various ways. For
- 368 example, if the contribution is a novel architecture, describing the architecture fully might suffice,
- 369 or if the contribution is a specific model and empirical evaluation, it may be necessary to either
- 370 make it possible for others to replicate the model with the same dataset, or provide access to
- 371 the model. In general, releasing code and data is often one good way to accomplish this, but
- 372 reproducibility can also be provided via detailed instructions for how to replicate the results,
- 373 access to a hosted model (e.g., in the case of a large language model), releasing of a model
- 374 checkpoint, or other means that are appropriate to the research performed.
- 375 • While NeurIPS does not require releasing code, the conference does require all submissions
- 376 to provide some reasonable avenue for reproducibility, which may depend on the nature of the
- 377 contribution. For example
- 378 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to
- 379 reproduce that algorithm.
- 380 (b) If the contribution is primarily a new model architecture, the paper should describe the
- 381 architecture clearly and fully.
- 382 (c) If the contribution is a new model (e.g., a large language model), then there should either be
- 383 a way to access this model for reproducing the results or a way to reproduce the model (e.g.,
- 384 with an open-source dataset or instructions for how to construct the dataset).
- 385 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are
- 386 welcome to describe the particular way they provide for reproducibility. In the case of
- 387 closed-source models, it may be that access to the model is limited in some way (e.g.,
- 388 to registered users), but it should be possible for other researchers to have some path to
- 389 reproducing or verifying the results.

390 5. Open access to data and code

391 Question: Does the paper provide open access to the data and code, with sufficient instructions to
392 faithfully reproduce the main experimental results, as described in supplemental material?

393 Answer: [No]

394 Justification: The package was cleared by IP review very close to the submission date. If accepted, the
395 paper will be updated with a link to the public package repo and relevant replication scripts for this
396 paper.

397 Guidelines:

- 398 • The answer NA means that paper does not include experiments requiring code.
- 399 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 400
- 401 • While we encourage the release of code and data, we understand that this might not be possible,
- 402 so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless
- 403 this is central to the contribution (e.g., for a new open-source benchmark).

- 404 • The instructions should contain the exact command and environment needed to run to reproduce
405 the results. See the NeurIPS code and data submission guidelines ([https://nips.cc/public/
406 guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 407 • The authors should provide instructions on data access and preparation, including how to access
408 the raw data, preprocessed data, intermediate data, and generated data, etc.
- 409 • The authors should provide scripts to reproduce all experimental results for the new proposed
410 method and baselines. If only a subset of experiments are reproducible, they should state which
411 ones are omitted from the script and why.
- 412 • At submission time, to preserve anonymity, the authors should release anonymized versions (if
413 applicable).
- 414 • Providing as much information as possible in supplemental material (appended to the paper) is
415 recommended, but including URLs to data and code is permitted.

416 6. Experimental setting/details

417 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,
418 how they were chosen, type of optimizer, etc.) necessary to understand the results?

419 Answer: [Yes]

420 Justification: Sections 2 and 3 elaborate on this

421 Guidelines:

- 422 • The answer NA means that the paper does not include experiments.
- 423 • The experimental setting should be presented in the core of the paper to a level of detail that is
424 necessary to appreciate the results and make sense of them.
- 425 • The full details can be provided either with the code, in appendix, or as supplemental material.

426 7. Experiment statistical significance

427 Question: Does the paper report error bars suitably and correctly defined or other appropriate informa-
428 tion about the statistical significance of the experiments?

429 Answer: [NA]

430 Justification: Because the implemented methods didn't include statistical tests in them.

431 Guidelines:

- 432 • The answer NA means that the paper does not include experiments.
- 433 • The authors should answer "Yes" if the results are accompanied by error bars, confidence
434 intervals, or statistical significance tests, at least for the experiments that support the main claims
435 of the paper.
- 436 • The factors of variability that the error bars are capturing should be clearly stated (for example,
437 train/test split, initialization, random drawing of some parameter, or overall run with given
438 experimental conditions).
- 439 • The method for calculating the error bars should be explained (closed form formula, call to a
440 library function, bootstrap, etc.)
- 441 • The assumptions made should be given (e.g., Normally distributed errors).
- 442 • It should be clear whether the error bar is the standard deviation or the standard error of the
443 mean.
- 444 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report
445 a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is
446 not verified.
- 447 • For asymmetric distributions, the authors should be careful not to show in tables or figures
448 symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 449 • If error bars are reported in tables or plots, The authors should explain in the text how they were
450 calculated and reference the corresponding figures or tables in the text.

451 8. Experiments compute resources

452 Question: For each experiment, does the paper provide sufficient information on the computer
453 resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

454 Answer: [Yes]

455 Justification: Yes, the lack of specialized hardware for using this package is mentioned

456 Guidelines:

- 457 • The answer NA means that the paper does not include experiments.

- 458
- 459
- 460
- 461
- 462
- 463
- 464
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

465 **9. Code of ethics**

466 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code
467 of Ethics <https://neurips.cc/public/EthicsGuidelines?>

468 Answer: [Yes]

469 Justification: Reviewed and complied with NeurIPS Code of Ethics

470 Guidelines:

- 471
- 472
- 473
- 474
- 475
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

476 **10. Broader impacts**

477 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts
478 of the work performed?

479 Answer: [Yes]

480 Justification: The limitations mention the ill effects of using causal fairness when inappropriate.

481 Guidelines:

- 482
- 483
- 484
- 485
- 486
- 487
- 488
- 489
- 490
- 491
- 492
- 493
- 494
- 495
- 496
- 497
- 498
- 499
- 500
- 501
- 502
- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

503 **11. Safeguards**

504 Question: Does the paper describe safeguards that have been put in place for responsible release of
505 data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or
506 scraped datasets)?

507 Answer: [NA]

508 Justification: Artifacts with high risk for misuse are not part of this publication.

509 Guidelines:

- 510
- 511
- 512
- 513
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.

- 514 • Datasets that have been scraped from the Internet could pose safety risks. The authors should
515 describe how they avoided releasing unsafe images.
516 • We recognize that providing effective safeguards is challenging, and many papers do not require
517 this, but we encourage authors to take this into account and make a best faith effort.

518 12. Licenses for existing assets

519 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper,
520 properly credited and are the license and terms of use explicitly mentioned and properly respected?

521 Answer: [Yes]

522 Justification: All authors have been notified and mentioned in the paper.

523 Guidelines:

- 524 • The answer NA means that the paper does not use existing assets.
525 • The authors should cite the original paper that produced the code package or dataset.
526 • The authors should state which version of the asset is used and, if possible, include a URL.
527 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
528 • For scraped data from a particular source (e.g., website), the copyright and terms of service of
529 that source should be provided.
530 • If assets are released, the license, copyright information, and terms of use in the package should
531 be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for
532 some datasets. Their licensing guide can help determine the license of a dataset.
533 • For existing datasets that are re-packaged, both the original license and the license of the derived
534 asset (if it has changed) should be provided.
535 • If this information is not available online, the authors are encouraged to reach out to the asset's
536 creators.

537 13. New assets

538 Question: Are new assets introduced in the paper well documented and is the documentation provided
539 alongside the assets?

540 Answer: [No]

541 Justification: The package was cleared by IP review very close to the submission date. If accepted, the
542 paper will be updated with a link to the public package repo and relevant replication scripts for this
543 paper.

544 Guidelines:

- 545 • The answer NA means that the paper does not release new assets.
546 • Researchers should communicate the details of the dataset/code/model as part of their sub-
547 missions via structured templates. This includes details about training, license, limitations,
548 etc.
549 • The paper should discuss whether and how consent was obtained from people whose asset is
550 used.
551 • At submission time, remember to anonymize your assets (if applicable). You can either create an
552 anonymized URL or include an anonymized zip file.

553 14. Crowdsourcing and research with human subjects

554 Question: For crowdsourcing experiments and research with human subjects, does the paper include
555 the full text of instructions given to participants and screenshots, if applicable, as well as details about
556 compensation (if any)?

557 Answer: [NA]

558 Justification: The paper does not have crowdsourcing experiments and research with human subjects

559 Guidelines:

- 560 • The answer NA means that the paper does not involve crowdsourcing nor research with human
561 subjects.
562 • Including this information in the supplemental material is fine, but if the main contribution of the
563 paper involves human subjects, then as much detail as possible should be included in the main
564 paper.
565 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other
566 labor should be paid at least the minimum wage in the country of the data collector.

567 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

568 Question: Does the paper describe potential risks incurred by study participants, whether such
569 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an
570 equivalent approval/review based on the requirements of your country or institution) were obtained?

571 Answer: [NA]

572 Justification: The paper does not have study participants

573 Guidelines:

- 574 • The answer NA means that the paper does not involve crowdsourcing nor research with human
575 subjects.
- 576 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be
577 required for any human subjects research. If you obtained IRB approval, you should clearly state
578 this in the paper.
- 579 • We recognize that the procedures for this may vary significantly between institutions and
580 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for
581 their institution.
- 582 • For initial submissions, do not include any information that would break anonymity (if applica-
583 ble), such as the institution conducting the review.

584 16. **Declaration of LLM usage**

585 Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard
586 component of the core methods in this research? Note that if the LLM is used only for writing,
587 editing, or formatting purposes and does not impact the core methodology, scientific rigor, or
588 originality of the research, declaration is not required.

589 Answer: [No]

590 Justification: [TODO]

591 Guidelines:

- 592 • The answer NA means that the core method development in this research does not involve LLMs
593 as any important, original, or non-standard components.
- 594 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what
595 should or should not be described.