

CausalFairnessInAction : A Novel Open Source Python Library For Causal Fairness Analysis

A Practical Guide to Diagnosing Bias with Causal Fairness

Forthcoming at: <https://github.com/amazon-science/causal-fairness-in-action>

Kriti Mahajan, Amazon

kritimhj@amazon.com

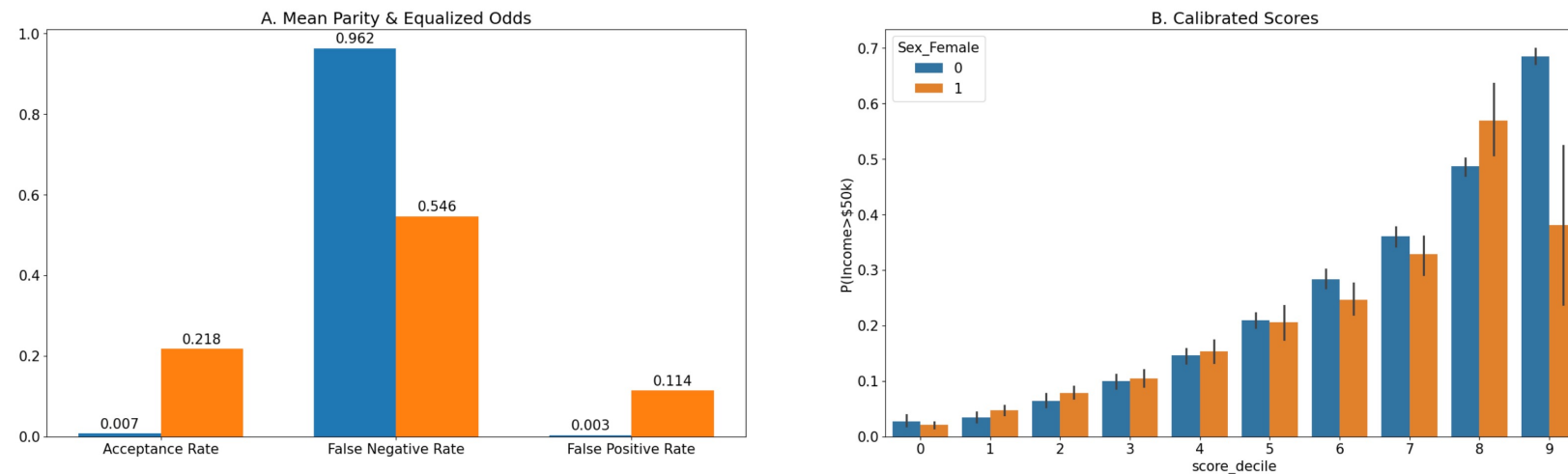


Motivation

Our model shows disparity. But is it discrimination?

- As machine learning enters high-stakes domains, assessing fairness becomes vital
- Statistical fairness metrics are widely used : They tell us *what* disparities exist in a model's outcome

Figure 1: Statistical Fairness Metrics (Ex: Adult Income Dataset)



Commonly used statistical fairness metrics :

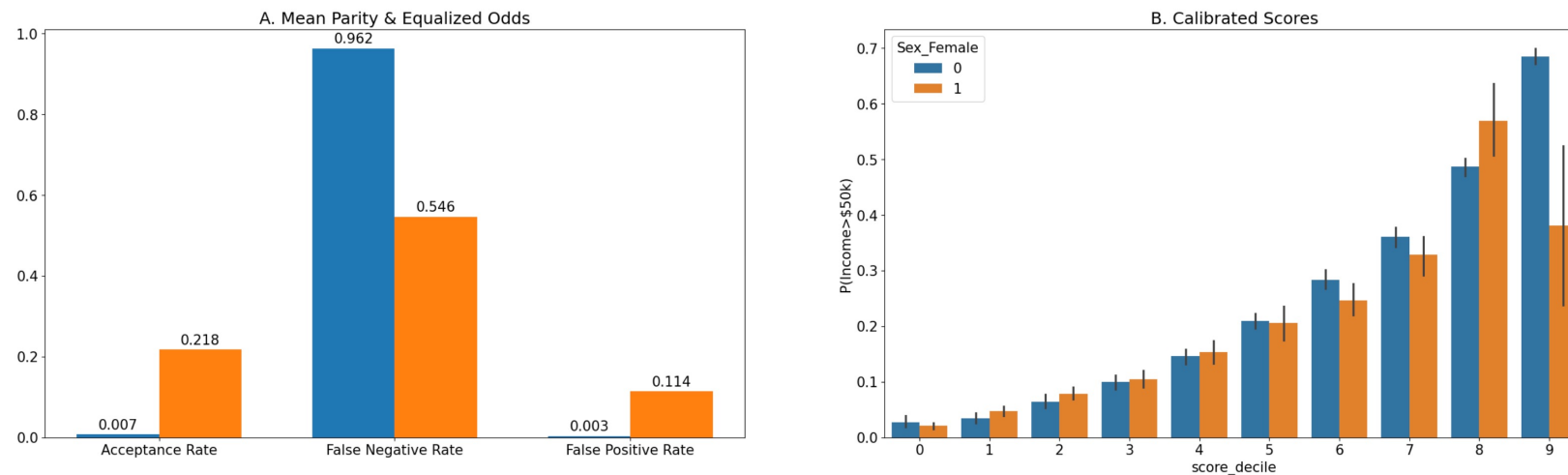
- Mean Parity
- Equalized Odds
- Calibrated Scores

Motivation

Our model shows disparity. But is it discrimination?

- But the typically used statistical fairness metrics have a key limitation: **They are associations - conditional probabilities - thus cannot explain why these disparities occur**
- Need methods that identify *why* disparities occur

Figure 1: Statistical Fairness Metrics (Ex: Adult Income Dataset)



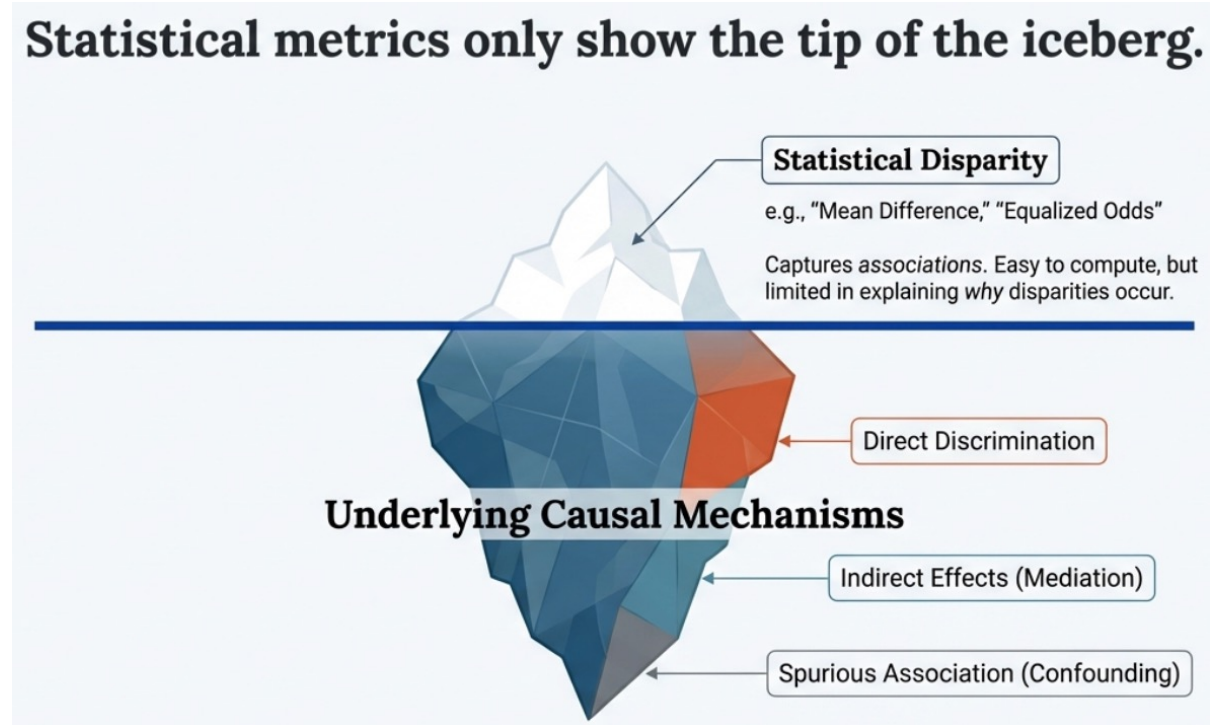
Commonly used statistical fairness metrics :

- Mean Parity
- Equalized Odds
- Calibrated Scores

Motivation

Our model shows disparity. But is it discrimination?

- **Causal Fairness metrics** based on Structural Causal Models (SCMs) attribute observed disparities to specific sources - protected attributes, mediators or confounders
- But they have limited adoption due to **technical & computational complexity**.



Contribution : CausalFairnessInAction

The first open-source Python package for computing causal fairness metrics

CausalFairnessInAction implements generalizable algorithms for three key metrics in the causal fairness literature:

- Counterfactual Effects for Mean Parity

Plecko, D. & Bareinboim, E., 2024. "Causal fairness analysis.". In: Foundations and Trends® in Machine Learning: Vol. 17, No. 3, pp 1–238

- Counterfactual Equalized Odds

Zhang, J. & Bareinboim, E., 2018. "Equality of opportunity in classification: A Causal approach." In: Advances in Neural Information Processing Systems.

- Counterfactual Fairness

Kusner, M.J. et al., 2017. "Counterfactual fairness" In: NIPS'17: Proceedings of the 31st International Conference on Neural Information Processing Systems

Contribution : CausalFairnessInAction

The first open-source Python package for computing causal fairness metrics

- **Practical:**

- Applicable across classification and regression tasks
- Designed to work with minimal identifiability constraints
- Doesn't require fully specified SCMs

- **Comprehensive:**

- Computes metrics at both **group** & **individual** levels
- Supports **intersectional analysis**

- **Efficient:** Optimized for scalability using Gaussian Mixture Models, parallelization to reduce latency

Enables **actionable audits** by decomposing statistical fairness metrics into causal components.

A Brief Literature Review

Broadly, there are three branches in the causal fairness literature:

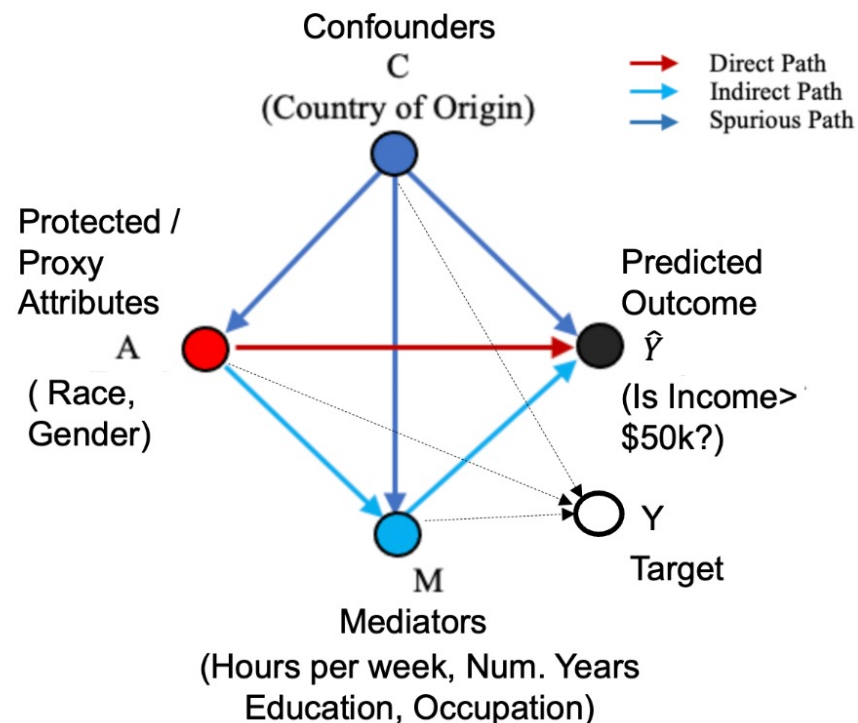
- **Counterfactual measurement** (focus of this package) : helps answering what-if cause-effect questions without running randomized control trials.
- **Sensitivity analysis** : how sensitive a model is to latent / confounding variables
- **Impact evaluation** : measures the long-term consequences of automated decision-making systems through the use of interventions.

What gets measured gets managed - causal identification of discrimination is crucial before moving on to remedial actions and impact analysis.

Methodology & Framework

The Standard Fairness Model

The `CausalFairnessDecomposition` class is built on the standard fairness model (Zhang, J. & Bareinboim, E., 2018):



- DAG includes protected attribute A, mediators M, confounders C, true outcome Y, prediction \hat{Y}

- Three paths from $A \rightarrow Y / \hat{Y}$:

- 1) **Direct (Causal Pathway):** $A \rightarrow Y / \hat{Y}$
(disparate treatment)

- 2) **Indirect (Causal Pathway):** $A \rightarrow M \rightarrow Y / \hat{Y}$
(mediated disparate impact)

- 3) **Spurious (Non-Causal Pathway):** $A \leftarrow C \rightarrow Y / \hat{Y}$
(confounding: spurious correlation)

Figure 2: Standard Fairness Model (Ex: Adult Income Dataset;
source - Zhang, J. & Bareinboim, E., 2018)

Implemented Metrics and Methods

`CausalFairnessDecomposition` class has 3 methods, each corresponding to a different causal fairness metric in the literature

Counterfactual Effects for Mean Parity	Counterfactual Equalized Odds	Counterfactual Fairness
<p>↓</p> <p>Group Level</p> <p>1. <code>`analyse_mean_difference`</code></p>	<p>↓</p> <p>Group Level</p> <p>2. <code>`analyse_equalized_odds`</code></p>	<p>↓</p> <p>Individual Level</p> <p>3. <code>`analyse_counterfactual_fairness`</code></p>
<p>Question Addressed →</p> <p>What would a group's acceptance rate be if they had the identity, mediators, or confounders of another group?</p>	<p>Question Addressed →</p> <p>How would a group's error rate (FPR/FNR) change under counterfactual conditions?</p>	<p>Question Addressed →</p> <p>If we changed an <i>individual's</i> protected attribute, would their predicted outcome change?</p>
<p>Decomposes</p> <p>Statistical Parity into Counterfactual Direct Effect (Ctf-DE), Indirect Effect (Ctf-IE), and Spurious Effect (Ctf-SE).</p>	<p>Decomposes</p> <p>Equalized Odds into Counterfactual Direct and Spurious Error Rates.</p>	<p>Measures</p> <p>Whether a prediction is counterfactually fair for each person in the dataset.</p>

Implemented Metrics and Methods

Counterfactual Effects : Calculating Disparate Treatment, Disparate Impact and Explaining the Causal Mechanism Behind Statistical Parity

- **Counterfactual Direct Effect (Ctf-DE):** measures direct discrimination along $A \rightarrow \hat{Y}$ by holding M and C constant. **Symmetric Ctf-DE** is the difference between the positive and negative effect of protected group membership.
Direct discrimination exists if Symmetric Ctf-DE > 0
→ **disparate treatment**
- **Counterfactual Indirect Effect (Ctf-IE):** measures indirect discrimination along $A \rightarrow M \rightarrow \hat{Y}$ by holding A and C fixed. **Symmetric Ctf-IE** is the difference between the positive and negative effect of having the group's mediating characteristics. **Indirect discrimination exists if Symmetric Ctf-IE > 0**
→ **disparate impact**
- **Counterfactual Spurious Effect (Ctf-SE):** measures confounding impact along $A \leftarrow C \rightarrow \hat{Y}$, varying C while fixing A and M
→ **disparate impact**

Implemented Metrics and Methods

Counterfactual Effects : Calculating Disparate Treatment, Disparate Impact and Explaining the Causal Mechanism Behind Statistical Parity

Mean Difference Causal Decomposition: Using these three counterfactual metrics, (Plecko, D. & Bareinboim, E.,2024) show that mean difference can be broken down into direct, indirect and spurious components as follows

$$\text{Mean Difference} = \text{Symmetric Ctf-DE} + \text{Symmetric Ctf-IE} + \text{Ctf-SE}$$

Implemented Metrics and Methods

Counterfactual Effects : Pseudo-algorithm for `analyse_mean_difference`

(Zhang, J. & Bareinboim, E., 2018) provide empirical formulas for estimation from observed data using conditional probabilities

→ **fully specified SCM not needed**

For each combination $m \in M$ and $c \in C$, we get a subset of D defined by (m, c) .

- For each (m, c) compute the expected outcome $\mathbf{E}(y \mid \mathbf{a}, m, c)$ for $\mathbf{a}_0, \mathbf{a}_1$.
- For each \mathbf{m} , calculate probability of \mathbf{m} when c is fixed under $\mathbf{a}_0, \mathbf{a}_1$: $\mathbf{P}(\mathbf{m} \mid \mathbf{a}_0, c), \mathbf{P}(\mathbf{m} \mid \mathbf{a}_1, c)$
- For each \mathbf{c} , calculate probability of \mathbf{c} under $\mathbf{a}_0, \mathbf{a}_1$: $\mathbf{P}(\mathbf{c} \mid \mathbf{a}_0), \mathbf{P}(\mathbf{c} \mid \mathbf{a}_1)$.

Inputs: D, A, M, C, a_0, a_1, y

1. For each $(m, c) \in D$:
 - Compute: $\mathbb{E}(Y = y \mid a_0, m, c)$
 - Compute: $\mathbb{E}(Y = y \mid a_1, m, c)$
2. Estimate via GMM:
 - $P(m \mid a_0, c), P(m \mid a_1, c)$
 - $P(c \mid a_0), P(c \mid a_1)$
3. Combine expectations and probabilities to compute the counterfactual effects

Implemented Metrics and Methods

Counterfactual Equalized Odds

- **Counterfactual Direct Error Rate**
- **Counterfactual Indirect Error Rate**
- **Counterfactual Spurious Error Rate**

Equalized Odds Causal Decomposition: Using these three counterfactual error metrics, (Zhang, J. & Bareinboim, E., 2018) show that equalized odds can be broken down into direct, indirect and spurious components as follows:

$$\text{Equalized Odds} = \text{Counterfactual Direct Error Rate} + \text{Counterfactual Indirect Error Rate} + \text{Counterfactual Spurious Error Rate}$$

Implemented Metrics and Methods

Counterfactual Equalized Odds : Pseudo-algorithm for `analyse_equalized_odds`

Limitation: cannot reliably estimate direct, indirect, and spurious effects in the presence of mediators due to lack of identifiability from conditioning on both Y and \hat{Y}

Solution : Refit estimator without M to get direct & spurious effects but no indirect effect

- For each $c \in C$ use the fitted estimator to compute predicted outcomes under \mathbf{a}_0 and \mathbf{a}_1 : $\hat{y}_{\mathbf{a}_0, \mathbf{c}}$ and $\hat{y}_{\mathbf{a}_1, \mathbf{c}}$.
- For each \mathbf{c} , compute how its probability under \mathbf{a}_0 , \mathbf{a}_1 : $P(\mathbf{c} | \mathbf{a}_0, \mathbf{y})$, $P(\mathbf{c} | \mathbf{a}_1, \mathbf{y})$.

Inputs: $D, A, C, a_0, a_1, y, \hat{f}$

1. For each $c_j \in D$:
 - Predict: $\hat{f}(c_j, a_0), \hat{f}(c_j, a_1)$
 - Obtain: $P(\hat{y}_{a_0, c_j}), P(\hat{y}_{a_1, c_j})$
2. Estimate via GMM:
 $P(c | a_0), P(c | a_1)$
3. Combine predictions and probabilities to compute the Cft-EO

Implemented Metrics and Methods

Counterfactual Fairness

Counterfactual Fairness is an individual level metric which is achieved if changing an individual i 's protected attributes doesn't change the individual's predicted outcome

Counterfactual Fairness Decomposition: Not supported

Implemented Metrics and Methods

Counterfactual Fairness : Pseudo-algorithm in `analyse_counterfactual_fairness`

Define a structural causal model (SCM) and get observed \mathbf{A}_{obs} and counterfactual \mathbf{A}_{cf} for $i \in D$

- Generate two datasets by intervening on \mathbf{A} :
 $\text{do}(\mathbf{A} = \mathbf{A}_{\text{obs}}) \rightarrow$ observed world
 $\text{do}(\mathbf{A} = \mathbf{A}_{\text{cf}}) \rightarrow$ counterfactual world
- Fit the model on each dataset to get predictions $\hat{\mathbf{Y}}_{\text{obs}}$ and $\hat{\mathbf{Y}}_{\text{cf}}$
- The model is **not counterfactually fair** if $\hat{\mathbf{Y}}_{\text{obs}} \neq \hat{\mathbf{Y}}_{\text{cf}}$

Inputs: $A, M, C, a_0, a_1, \text{DAG}$

1. Fit SCM using DAG and dataset D
2. For each individual $i \in D$:
 - Get A_{obs} (observed) and A_{cf} (counterfactual)
 - Sample from SCM under:
 $\text{do}(A = A_{\text{obs}}) \Rightarrow D_{\text{obs}}$
 $\text{do}(A = A_{\text{cf}}) \Rightarrow D_{\text{cf}}$
 - Predict: $\hat{f}(D_{\text{obs}}), \hat{f}(D_{\text{cf}})$
 - Check: $Y_{\text{obs}} \neq Y_{\text{cf}}$

Implemented Metrics and Methods

Example API Call Applied To The Adult Income Dataset

```
cf = CausalFairnessDecomposition(**{"X": X_train,
                                     "y_true" : y_train.values,
                                     "y_pred": X_train["prediction"],
                                     "model": trained_logistic_regression_classifier,
                                     "protected_attr": ["Sex_Female"],
                                     "mediators": ([x for x in X_train.columns if "Occupation_" in x]
                                                  + ["EducationNum"]
                                                  + ["HoursPerWeek"]),
                                     "confounders": [x for x in X_train.columns if "Country_" in x],
                                     "yi": 1,
                                     "advantage_group": 0,
                                     "disadvantage_group": 1,
                                     "continuous": False})

mean_diff_decomposition = cf.analyse_mean_difference(how='decompose')
fig_md, ax_md = plot_mean_diff_waterfall(mean_diff_decomposition, mean_difference)
```

Implemented Metrics and Methods (Contd.)

Example API Call Applied To The Adult Income Dataset

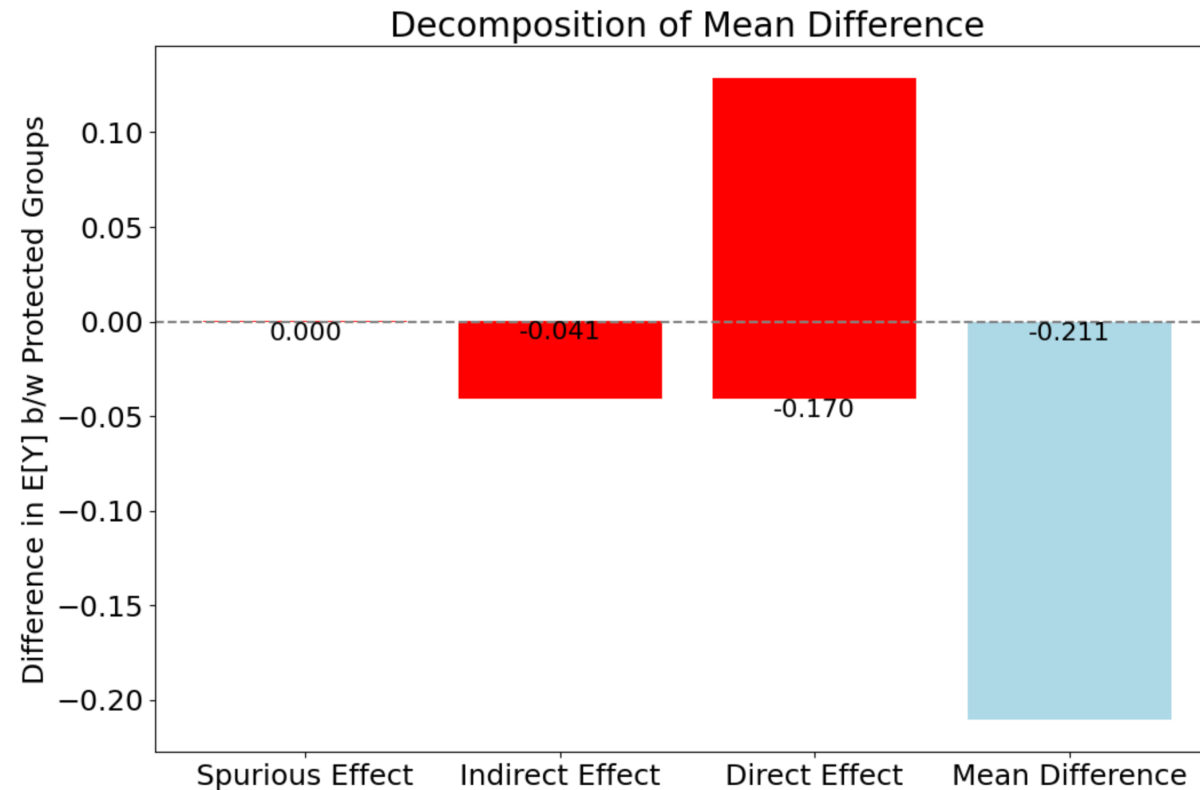


Figure 3: Counterfactual Effects Decomposition of Mean Difference for Adult Income Dataset

$$\text{Mean Difference} = \text{Symmetric Ctf-DE} + \text{Symmetric Ctf-IE} + \text{Ctf-SE}$$

Implemented Metrics and Methods

Computational Optimizations

- GMMs for conditional probability estimation
- Parallelization
- No specialized hardware requirement

Application To Benchmark Datasets

Overview of Findings

We benchmarked the library on 3 datasets : **Adult Income, COMPAS, LSAC**

- **Direct discrimination is the primary contributor** to mean difference and equalized odds across all 3 datasets
- The classifier for Adult Income, COMPAS is not counterfactually fair but is counterfactually fair for LSAC i.e. **group fairness can differ from individual fairness**
- **Intersectional Analysis (Race x Sex) worsens direct discrimination** across all three datasets

Table 1: CausalFairnessInAction Benchmarking Results

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^i	ER^s	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

Application To Benchmark Datasets

Adult Income Dataset

Task : logistic regression fit to predict $P(\text{Income} > \$50k)$

- **Counterfactual Effects:** On average, women are 20.3% less likely than men to be predicted as earning above \$50k. Most of this disparity (16.5%) is due to disparate treatment (direct effect)
- **Counterfactual Equalized Odds:** Not identifiable due to mediator issues
- **Counterfactual Fairness:** Not counterfactually fair—changing female→male increases predicted $P(\text{Income} > \$50k)$

DAG & Variables used for Fitting `CausalFairnessDecomposition`

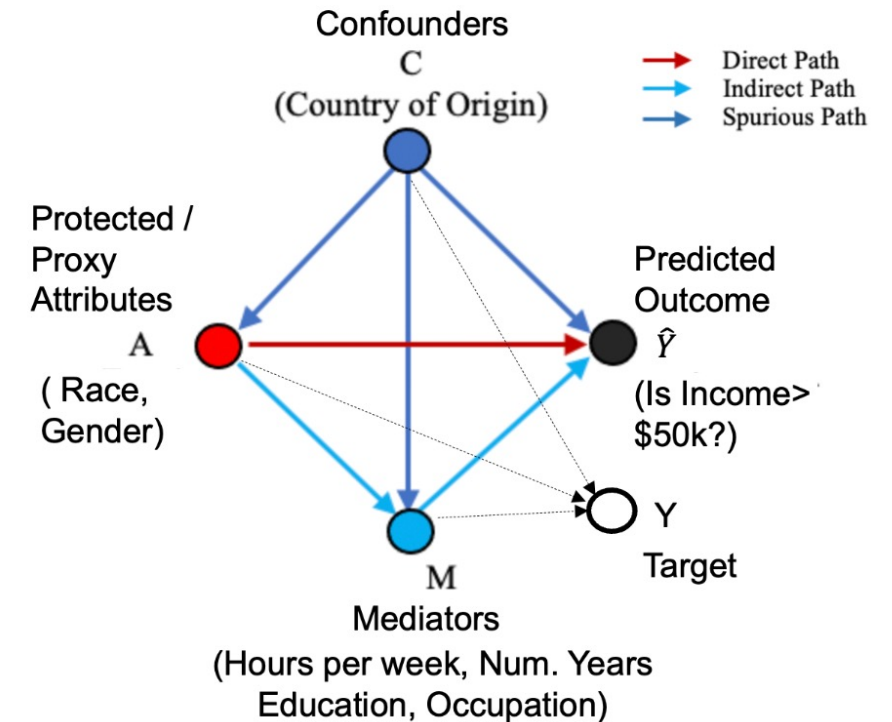


Figure 2: Adult Income Standard Fairness Model (Zhang, J. & Bareinboim, E., 2018)

Application To Benchmark Datasets

COMPAS Dataset

Task : logistic regression fit to predict $P(\text{Recidivism})$

- **Counterfactual Effects:** Black individuals are 32.6% more likely than white individuals to be predicted as high-risk for recidivism. Majority is attributed to disparate treatment (15.4%). Disparate impact comes from both M and C: confounders raise risk by ~10% (spurious effect); M contributes an additional 7.1% (indirect effect)
- **Counterfactual Equalized Odds:** Excluding M does not make the model naive, though it increases error rates. Decomposing FPR/FNR shows most of the disparity stems from direct discrimination: 29.7% of the 41% FPR and 26.5% of the 42% FNR
- **Counterfactual Fairness:** Not counterfactually fair—changing black \rightarrow white decreases predicted $P(\text{Recidivism})$

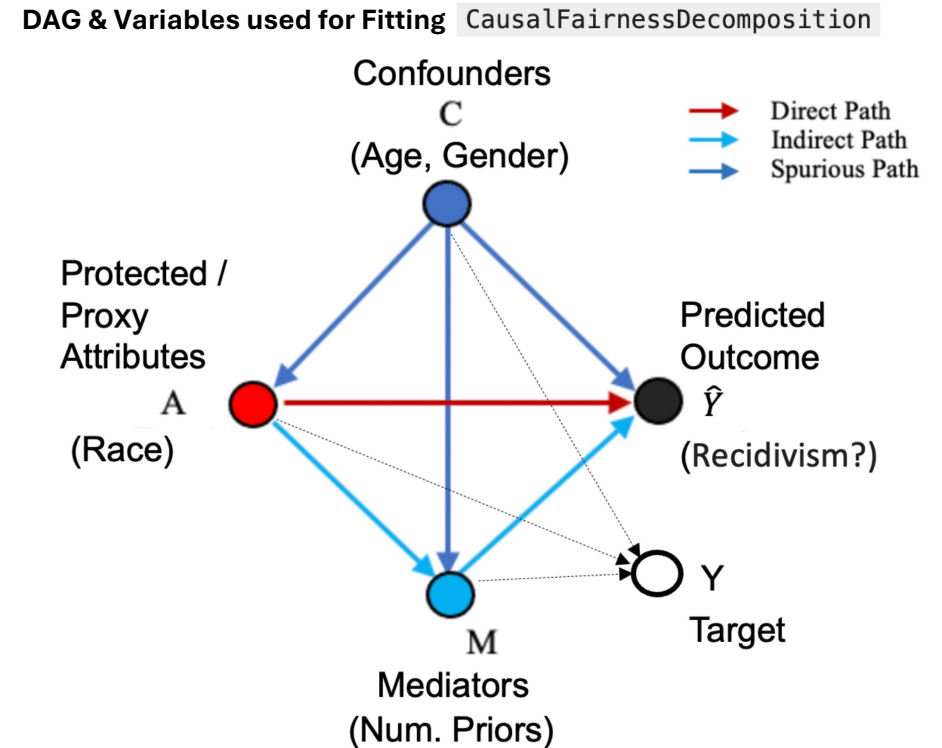


Figure 4: COMPAS Standard Fairness Model (Zhang, J. & Bareinboim, E., 2018)

Application To Benchmark Datasets

LSAC Dataset

Task : fit Random Forest regressor to **predict** average grade

- **Counterfactual Effects:** The predicted average grade for the white subgroup is 0.978 higher than for the Black subgroup with majority of the gap (0.55) due to direct discrimination.
- **Counterfactual Equalized Odds:** Not applicable since this is a regression task
- **Counterfactual Fairness:** The model is counterfactually fair, illustrating that fairness can differ at the individual vs. group level

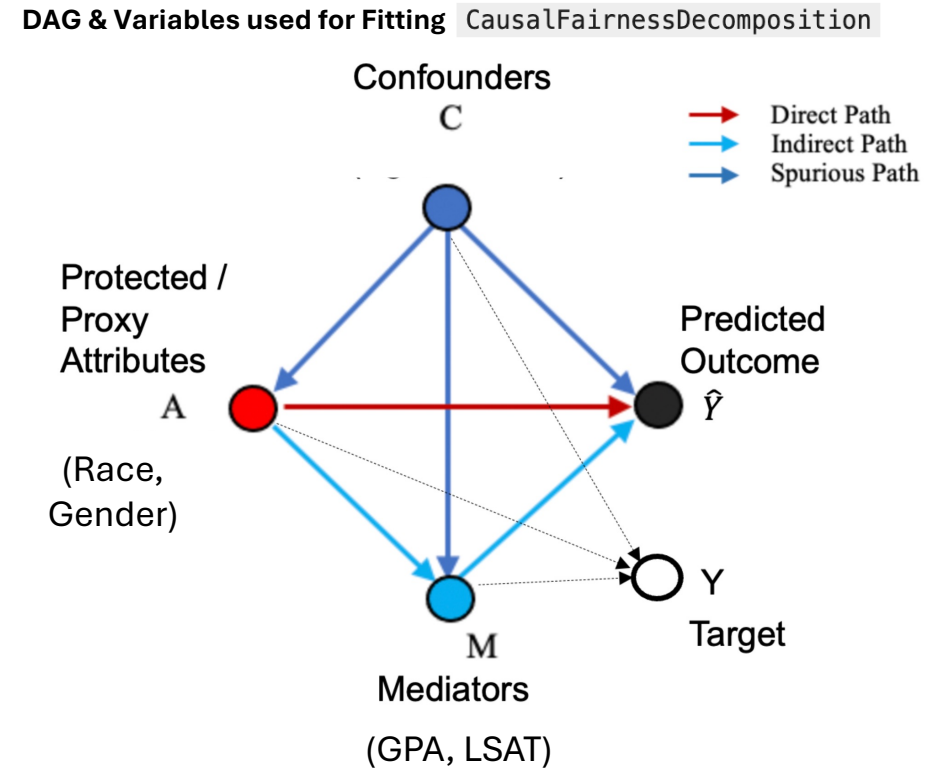


Figure 5: LSAC Standard Fairness Model (Kusner, M.J. et al., 2017.)

Application To Benchmark Datasets

Intersectional Analysis

- **Adult Income:** Black women are 22.1% less likely than white men to have $P(\text{Income} > \$50k)$ —2% more than the non-intersectional gender gap— with direct effect being the largest contributor. The model is also more counterfactually unfair: changing a Black woman's identity to a white man increases $P(\text{Income} > \$50k)$ by 6% (vs. 2% non-intersectionally)
- **COMPAS:** The mean difference in $P(\text{Recidivism})$ between Black men and white women (60%) exceeds the non-intersectional racial gap, with direct discrimination as the main driver. Counterfactual unfairness also rises by ~11%: changing a Black man's identity to a white woman lowers predicted recidivism by 64% (vs. 55%)
- **LSAC:** The mean difference between Black women and white men is slightly higher than the non-intersectional comparison (0.99 vs. 0.978), mainly due to direct discrimination. As in the non-intersectional case, the model remains counterfactually fair.

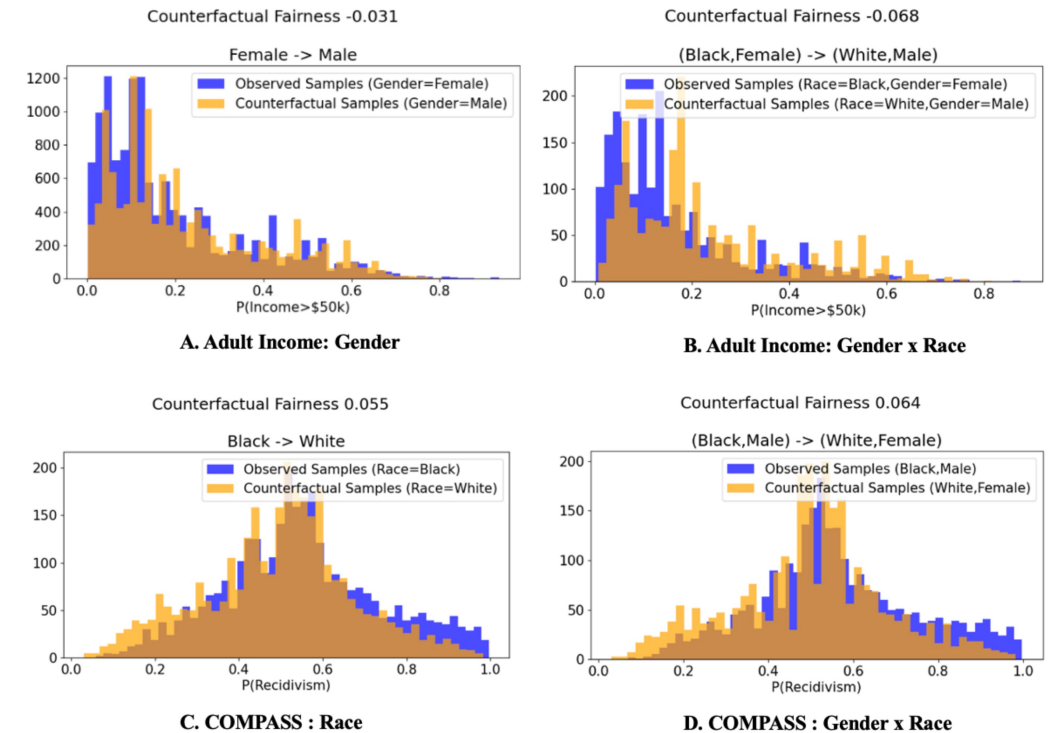


Figure 6: Counterfactual Fairness

Limitations of CausalFairnessInAction

- **Lack of identifiability can limit analysis**
 - Ex: in the Adult Income dataset, identifiability issues prevented the causal decomposition of equalized odds
- **Lack of methods for falsifying DAGs** in the presence of competing causal models can lead to disagreements about result validity
- **Defining a hypothetical intervention on protected attributes** remains a fraught process

Future Work

- Use package results to guide fairness interventions:
 - Feature selection
 - Sample-level reweighting
 - Multi-world regularization
- Integrate bias-reduction algorithms into package

Conclusion

We introduced ***CausalFairnessInAction*** - the first open-source generalizable implementation for calculating key causal fairness metrics

- Applied it to 3 fairness benchmarking datasets

Demonstrated how CausalFairnessInAction provides practitioners with the actionable insight

- Ex : at the very least the Adult Income model must eliminate at least 16.5% difference in statistical parity, while the COMPAS model needs to address 15.4% disparity in statistical parity and 29.7-26.5% in error rates (all of which can be attributed to direct discrimination)

Thank you!

Explore the library, replicate these findings, and apply causal analysis to your own models.



github.com/amazon-science/causal-fairness-in-action

**Forthcoming*

Appendix

CausalFairnessInAction: An Open Source Python Library for Causal Fairness Analysis

Kriti Mahajan, Amazon, kritimhj@amazon.com
Forthcoming at:
<https://github.com/amazon-science/causal-fairness-in-action>



Motivation & Contribution

The Problem: As machine learning enters high-stakes domains, assessing fairness becomes vital—but the typically used statistical fairness metrics have a key limitation: They are **associations(conditional probabilities)** thus, they can state **what the observed disparity is but not why it exists**.

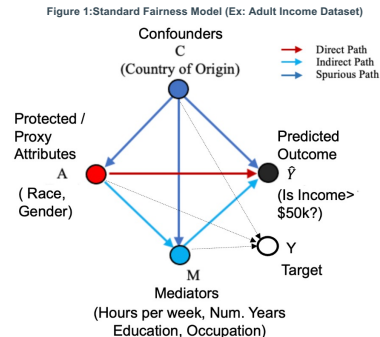
Causal Fairness metrics solve this by using Structural Causal Models (SCMs) to uncover generating mechanisms, but have limited adoption due to **technical & computational complexity**.

The Solution: **CausalFairnessInAction**, the **first open-source Python package** for computing diverse causal fairness metrics, enabling actionable audits by decomposing statistical disparities into causal components.

- Practical:** applicable across classification and regression tasks ; designed to work with minimal identifiability constraints; doesn't require fully specified SCMs.
- Comprehensive:** Computes metrics at both **group & individual** levels; supports **intersectional analysis**.
- Efficient:** Optimized for scalability using Gaussian Mixture Models, parallelization to reduce latency

Methodology & Framework

The **CausalFairnessDecomposition** class is built on the standard fairness model [1]:



Three Implemented Metrics and Methods

analyse_mean_difference → **Implements Counterfactual Effects**¹

Query : 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?

Supported Decompositions: Direct, Indirect, Spurious

analyse_equalized_odds → **Implements Counterfactual Equalized Odds**²

Query : What would the disadvantaged (advantaged) **group's error rate** be if they had A, M or C of the advantaged (disadvantaged) group?

Supported Decompositions: Direct, Spurious

analyse_counterfactual_fairness → **Implements Counterfactual Fairness**³

Query : What would the disadvantaged (advantaged) **individual's predicted Y** be if they had the A, M, and C of the advantaged (disadvantaged) group?

Supported Decompositions: N/A

Table 1: Pseudo-Algorithms for Causal Fairness Metrics

analyse_mean_difference	analyse_equalized_odds	analyse_counterfactual_fairness
Inputs: D, A, M, C, a_0, a_1, y	Inputs: $D, A, C, a_0, a_1, y, \hat{f}$	Inputs: $A, M, C, a_0, a_1, \text{DAG}$
1. For each $(m, c) \in D$: - Compute: $E(Y = y a_0, m, c)$ - Compute: $E(Y = y a_1, m, c)$ 2. Estimate via GMM: $P(m a_0, c), P(m a_1, c)$ $P(c a_0), P(c a_1)$ 3. Combine expectations and probabilities to compute the counterfactual effects	1. For each $c_j \in D$: - Predict: $\hat{f}(c_j, a_0), \hat{f}(c_j, a_1)$ - Obtain: $P(y_{a_0, c_j}), P(y_{a_1, c_j})$ 2. Estimate via GMM: $P(c a_0), P(c a_1)$ 3. Combine predictions and probabilities to compute the Cft-EO	1. Fit SCM using DAG and dataset D 2. For each individual $i \in D$: - 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: $Y_{obs} \neq Y_{cft}$

Code block 1: Example API Call Applied To The Adult Income Dataset

```
cfd = CausalFairnessDecomposition(**{"X": X_train,
                                     "y_true": y_train.values,
                                     "y_pred": X_train["prediction"],
                                     "model": trained_logistic_regression_classifier,
                                     "protected_attr": ["Sex_Female"],
                                     "mediators": [x for x in X_train.columns if "Occupation" in x]
                                     + ["EducationNum"]
                                     + ["HoursPerWeek"]},
                                     "confounders": [x for x in X_train.columns if "Country_" in x],
                                     "yi": 1,
                                     "advantage_group": 0,
                                     "disadvantage_group": 1,
                                     "continuous": False})

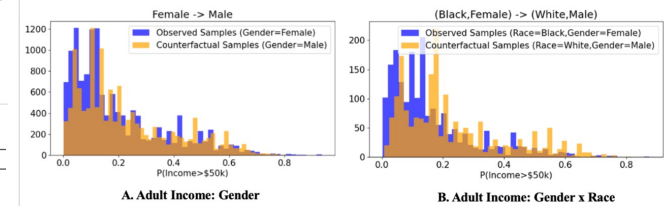
mean_diff_decomposition = cfd.analyse_mean_difference(how='decompose')
fig_md, ax_md = plot_mean_diff_waterfall(mean_diff_decomposition, mean_difference)
```

Application To Benchmark Datasets

We benchmarked the library on 3 datasets : **Adult Income**, **COMPAS**, and **LSAC**

- Direct discrimination is the primary contributor** to mean difference and equalized odds across all 3 datasets
- The classifier for Adult Income, COMPAS is not counterfactually fair but is counterfactually fair for LSAC i.e. **group fairness can differ from individual fairness**
- Intersectional Analysis (Race x Sex) worsens direct discrimination** across all three datasets

Figure 2: Counterfactual Fairness Plots



Limitations

- Lack of identifiability can limit analysis:** Ex - in the Adult Income dataset, identifiability issues prevent the causal decomposition of equalized odds

Conclusion & Future Work

- Actionable:** Provides specific targets for bias mitigation (e.g., fixing the 16.5% direct effect in Adult Income)
- Future:** Extending the package to include remediation algorithms and sensitivity analysis.

References

- Plecko, D. & Bareinboim, E., 2024. "Causal fairness analysis." In: Foundations and Trends® in Machine Learning: Vol. 17, No. 3, pp 1–238
- Zhang, J. & Bareinboim, E., "Equality of opportunity in classification: A Causal approach." In: Advances in Neural Information Processing Systems.
- Kusner, M.J. et al., 2017. "Counterfactual fairness" In: NIPS'17: Proceedings of the 31st International Conference on Neural Information Processing Systems

Table 2: CausalFairnessInAction Benchmarking											
Dataset	Protected Attribute	Mean Difference	FNR	FPR	$DE_a^{sym}(y/a)$	$IE_a^{sym}(y/a)$	$SE_{a_0, a_1}(y/a)$	ER^d	ER^i	ER^s	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 (-.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