

A GENERIC FRAMEWORK FOR CONFORMAL FAIRNESS

Aditya T. Vadlamani^{1,*}, Anutam Srinivasan^{1,*}, Pranav Maneriker^{2,†},
Ali Payani³, Srinivasan Parthasarathy¹

¹The Ohio State University ²Dolby Laboratories ³Cisco Systems

ABSTRACT

Conformal Prediction (CP) is a popular method for uncertainty quantification with machine learning models. While conformal prediction provides probabilistic guarantees regarding the coverage of the true label, these guarantees are agnostic to the presence of sensitive attributes within the dataset. In this work, we formalize *Conformal Fairness*, a notion of fairness using conformal predictors, and provide a theoretically well-founded algorithm and associated framework to control for the gaps in coverage between different sensitive groups. Our framework leverages the exchangeability assumption (implicit to CP) rather than the typical IID assumption, allowing us to apply the notion of Conformal Fairness to data types and tasks that are not IID, such as graph data. Experiments were conducted on graph and tabular datasets to demonstrate that the algorithm can control fairness-related gaps in addition to coverage aligned with theoretical expectations.

1 INTRODUCTION

Machine learning (ML) models are increasingly used to make critical decisions in many fields of human endeavor, making it essential to quantify the uncertainty associated with their predictions. Conformal Prediction (CP) is a distribution-free framework (Vovk et al., 2005) which produces confidence sets with rigorous theoretical guarantees and has become popular in real-world applications (Cherian & Bronner, 2020). Post-hoc CP allows for facile integration into ML pipelines and applies to a wide variety of data types, including graph data (H. Zargarbashi et al., 2023; Huang et al., 2024), because of its weaker requirement of a *statistical exchangeability*.

Relatedly, ensuring the fairness of machine learning models is vital for their high-stakes deployments in critical decision-making. Biases affect ML models at different stages - from data collection to algorithmic learning stages (Mehrabi et al., 2021). During the data collection stage, measurement and representation biases can skew how each feature is interpreted, leading to inaccurate determinations by learning models. Algorithmic bias, caused by model design choices and prioritization of specific metrics while learning the model, can also lead to unfair outcomes. Many models inherit biases from historical outcomes (Kallus & Zhou, 2018; Dwork et al., 2012) and inadvertently skew decisions towards members of certain advantaged groups (Mehrabi et al., 2021). These biases have led to several global actors proposing and requiring practitioners to adhere to certain *fairness* standards (Hirsch et al., 2023). To facilitate ML pipeline and model adherence to socio-cultural or regulatory fairness standards, researchers have proposed methods to either construct fair-predictors (Alghamdi et al., 2022; Creager et al., 2019; Zhao et al., 2023) or audit fairness claims made by deployed machine learning models (Ghosh et al., 2021; Maneriker et al., 2023; Yan & Zhang, 2022).

However, these efforts on fairness (predictors, auditing, and uncertainty quantification) primarily focus on binary classification, often implicitly relying on the independent and identically distributed (IID) assumption, and largely do not bridge fairness and uncertainty quantification. The need to both quantify uncertainty and ensure fairness considerations are met is critical. A few researchers have started to examine how to assess (and possibly improve) the prediction quality of unreliable models (Wang & Wang, 2024) while meeting socio-cultural or regulatory standards of fairness. However, these efforts are limited in that they either require knowledge of group membership at inference time (a somewhat impractical assumption) (Lu et al., 2022) or are model-specific (Wang & Wang, 2024).

*Equal Contribution, †This work was done while the author was a student at The Ohio State University

Key Contributions: We propose a novel and comprehensive Conformal Fairness (CF) Framework to redress these concerns.

First, we develop the theoretical insights that facilitate how our framework leverages CP’s distribution-free approach to build and construct fair uncertainty sets according to user-specified notions of fairness. Our framework is not only comprehensive but also highly flexible, as it can be adapted to bespoke user-specified fairness criteria. This adaptability ensures that the framework can be customized to meet the specific needs of different users, enhancing its practicality and usability.

Second, the weaker (exchangeability) assumptions required by CP allow us to extend the utility of our framework to fairness problems in graph models. Graph models, in particular, suffer from the *homophily effect*, which exacerbates inherent segregation due to node linkages and causes further biases in predictions (Current et al., 2022; Dong et al., 2023; He et al., 2023).

Third, we discuss how our approach serves as a fairness auditing tool for conformal predictors. This function is important as it allows one to verify model fairness, ensuring that fairness is not just a theoretical concept but a practical reality in predictive modeling.

Finally, we demonstrate the effectiveness of our CF Framework by evaluating fairness using multiple popular fairness metrics for multiple different conformal predictors on both real-world graph and tabular fairness datasets.

2 BACKGROUND

2.1 CONFORMAL PREDICTION

Conformal Prediction (Vovk et al., 2005) is a framework for quantifying the uncertainty of a model by constructing prediction sets that satisfy a *coverage* guarantee. For expository simplicity, we will focus on split (or inductive) conformal prediction (CP) in the classification setting. Given a calibration dataset, $\mathcal{D}_{\text{calib}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and a test point $(\mathbf{x}_{n+1}, y_{n+1})$, where $\mathbf{x}_i \in \mathcal{X} = \mathbb{R}^d$ and $y_i \in \mathcal{Y} = \{0, \dots, K - 1\}$, CP is used to construct a prediction set $\mathcal{C}(\mathbf{x}_{n+1})$ such that:

$$1 - \alpha \leq \Pr[y_{n+1} \in \mathcal{C}(\mathbf{x}_{n+1})] \leq 1 - \alpha + \frac{1}{n + 1}, \quad (1)$$

where $1 - \alpha \in (0, 1)$ is the coverage bound. Concretely, given a non-conformity score function $s : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, let $\hat{q}(\alpha) = \text{Quantile}\left(\frac{\lfloor (n+1)(1-\alpha) \rfloor}{n}; \{s(\mathbf{x}_i, y_i)\}_{i=1}^n\right)$ be the *conformal quantile*. Then $\mathcal{C}_{\hat{q}(\alpha)}(\mathbf{x}_{n+1}) = \{y \in \mathcal{Y} : s(\mathbf{x}_{n+1}, y) \leq \hat{q}(\alpha)\}$ is a prediction set that satisfies Equation 1.

Evaluating CP: *Coverage* quantifies the true test time probability $\Pr[y_{n+1} \in \mathcal{C}_{\hat{q}(\alpha)}(\mathbf{x}_{n+1})]$ while *efficiency* is the average test prediction set size, $|\mathcal{C}_{\hat{q}(\alpha)}(\mathbf{x}_{n+1})|$. Intuitively, there is an inverse relationship between coverage and efficiency, as a higher desired coverage is harder to achieve, the method may produce larger prediction sets to satisfy the guarantee. In CP, the only assumption made about the data is that $\mathcal{D}_{\text{calib}} \cup \{(\mathbf{x}_{n+1}, y_{n+1})\}$ is *exchangeable* – a weaker notion than iid, enabling its use on non-iid data, including graph data.

Graph CP: In this work, we focus on the node classification task. Given an attributed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X})$, where \mathcal{V} is the set of nodes, \mathcal{E} is the set of edges, and \mathbf{X} is the set of node attributes. Let \mathbf{A} be the adjacency matrix for the graph. Further, let $\mathcal{Y} = \{0, \dots, K - 1\}$ denote the set of classes associated with the nodes. For $v \in \mathcal{V}$, $\mathbf{x}_v \in \mathbb{R}^d$ denotes its features and $y_v \in \mathcal{Y}$ denotes its true class. The task of node classification is to learn a model that predicts the label for each node given all the node features and the adjacency matrix, i.e., $(\mathbf{X}, \mathbf{A}, v) \mapsto y_v$. In the transductive setting, the entire graph, including test points, is accessible during the base model training. In this scenario, for any trained permutation-equivariant function (e.g., GNN) trained on a set of training/validation nodes, the scores produced on the calibration set and test set are exchangeable, thus enabling CP to be applied (H. Zargarbashi et al., 2023; Huang et al., 2024).

2.2 FAIRNESS METRICS

Group (or statistical) fairness requires individuals from different sensitive groups to be treated equally. Sensitive groups are subpopulations characterized by sensitive attributes, including gen-

der, race, and/or ethnicity. Group fairness metrics aim to observe bias in the predictions of a model between the different groups in a dataset. This work considers several popular fairness metrics, including equal opportunity, equalized odds, demographic parity, predictive equality, and predictive parity. For generality, we define the metrics for the multiclass setting with an n -ary sensitive attribute. Let \mathcal{Y}^+ denote the set of advantaged labels (e.g., “is_approved” in a loan approval task), Y be the true label, and \hat{Y} be the predicted label from a classifier. Let \mathcal{G} be the set of all groups for the sensitive attribute(s). Table A1 discusses the formal definitions of different fairness metrics considered in this work.

Achieving *exact fairness* (i.e., the equality in Table A1) can be challenging or, in some cases, impossible (Barocas et al., 2023). Often, regulatory requirements focus on the difference (or ratio) in probabilities between groups for any given positive label. This is achieved by ensuring the difference (or ratio) meets a prespecified *closeness criterion*. For example, many regulatory bodies consider the **Four-Fifths Rule** (EEOC, 1979; Feldman et al., 2015), which asserts that the ratio of the selection probabilities between groups is at least 0.8.

3 CONFORMAL FAIRNESS (CF) FRAMEWORK

In this section, we propose a theoretically well-founded framework using conformal predictors to control for fairness disparity between different sensitive groups. The framework is motivated by adapting the standard CP algorithm to determine conditional coverage *given* a score threshold, λ , for the prediction sets (i.e. $\mathcal{C}_\lambda(\mathbf{x}_{n+1}) = \{y \in \mathcal{Y} \mid s(\mathbf{x}_{n+1}, y) \leq \lambda\}$). Depending on the fairness metric, fairness disparity refers to gaps in group-conditional or group-and-class-conditional coverages between groups and advantaged labels. The conditional coverages are leveraged to evaluate if fairness is achieved for some closeness criterion c for different fairness metrics. This is achieved by searching a *threshold space* Λ for an optimal threshold λ_{opt} that achieves the closeness criteria. The framework also handles user-defined metrics as discussed in Section 3.4, thus controlling for quantities, potentially orthogonal to conditional coverage.

3.1 EXEMPLAR CONFORMAL FAIRNESS (CF) METRICS

For conformal fairness, we adapt popular fairness metrics defined for *multiclass classification* (shown in Table A1). For standard point-wise predictions, fairness measures are concerned with the probability a prediction is a specific label (i.e., $\tilde{y} = \hat{Y}$), given a condition, i.e., $X \in g_a, Y = \tilde{y}$ for Equal Opportunity, for a particular covariate (X, Y) . We replace equivalence to the predicted value with set membership ($\tilde{y} \in \mathcal{C}_\lambda(X)$) to adapt these notions for prediction sets. The adapted conformal fairness metrics are in Table 1.

Table 1: Conformal Fairness Metrics.

Metric	Definition
Demographic (or Statistical) Parity	$ \Pr[\tilde{y} \in \mathcal{C}_\lambda(X) \mid X \in g_a] - \Pr[\tilde{y} \in \mathcal{C}_\lambda(X) \mid X \in g_b] < c, \forall g_a, g_b \in \mathcal{G}, \forall \tilde{y} \in \mathcal{Y}^+$
Equal Opportunity	$ \Pr[\tilde{y} \in \mathcal{C}_\lambda(X) \mid Y = \tilde{y}, X \in g_a] - \Pr[\tilde{y} \in \mathcal{C}_\lambda(X) \mid Y = \tilde{y}, X \in g_b] < c, \forall g_a, g_b \in \mathcal{G}, \forall \tilde{y} \in \mathcal{Y}^+$
Predictive Equality	$ \Pr[\tilde{y} \in \mathcal{C}_\lambda(X) \mid Y \neq \tilde{y}, X \in g_a] - \Pr[\tilde{y} \in \mathcal{C}_\lambda(X) \mid Y \neq \tilde{y}, X \in g_b] < c, \forall g_a, g_b \in \mathcal{G}, \forall \tilde{y} \in \mathcal{Y}^+$
Equalized Odds	Equal Opp. and Pred. Equality
Predictive Parity	$\Pr[Y = \tilde{y} \mid \tilde{y} \in \mathcal{C}_\lambda(X), X \in g_a] = \Pr[Y = \tilde{y} \mid \tilde{y} \in \mathcal{C}_\lambda(X), X \in g_b], \forall g_a, g_b \in \mathcal{G}, \forall \tilde{y} \in \mathcal{Y}^+$

3.2 CONFORMAL FAIRNESS (CF) THEORY

Before presenting our framework, we first lay out the necessary theoretical groundwork. Detailed proofs are in Appendix B.

Filtering $\mathcal{D}_{\text{calib}}$: Group fairness metrics are evaluated on a subset of the population, defined by a condition on the data (i.e., membership in a group, true label value). For example, Demographic Parity is evaluated *per group* ($X \in g_a$ in definition), while Equal Opportunity is evaluated *per group and true label* ($Y = y, X \in g_a$ in definition). To formalize this notion, let M denote

a fairness metric (e.g. Equal Opportunity) and define $F_M : \mathcal{X} \times \mathcal{Y} \times \mathcal{G} \times \mathcal{Y}^+ \rightarrow \{0, 1\}$ be a **filter function** which maps a calibration point along with a group and positive label, $(\mathbf{x}_i, y_i, g, \tilde{y})$, to 0 or 1 depending on whether the condition for the fairness metric, M , is satisfied. For Equal Opportunity, F_M would instantiate to $F_{EO}(\mathbf{x}_i, y_i, g, \tilde{y}) := \mathbf{1}[\mathbf{x}_i \in g \cap y_i = \tilde{y}]$. We can filter $\mathcal{D}_{\text{calib}}$ to be $\mathcal{D}_{\text{calib}(g, \tilde{y})} = \{(\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{calib}} \mid F_M(\mathbf{x}_i, y_i, g, \tilde{y}) = 1\}$. By doing so, we provide guarantees regarding the conditional coverages as stated in Lemma 3.1

Lemma 3.1. *For any $(g, \tilde{y}) \in \mathcal{G} \times \mathcal{Y}^+$, calibrating on $\mathcal{D}_{\text{calib}(g, \tilde{y})} = \{(\mathbf{x}_i, y_i) \mid F_M(\mathbf{x}_i, y_i, g, \tilde{y}) = 1\}$ guarantees the following about the conditional coverage:*

$$1 - \alpha \leq \Pr[y_{n+1} \in \mathcal{C}_\lambda(\mathbf{x}_{n+1}) \mid F_M(\mathbf{x}_{n+1}, y_{n+1}, g, \tilde{y}) = 1] \leq 1 - \alpha + \frac{1}{|\mathcal{D}_{\text{calib}(g, \tilde{y})}| + 1} \quad (2)$$

The *interval width* is $\frac{1}{|\mathcal{D}_{\text{calib}(g, \tilde{y})}| + 1}$.

Prior work (Ding et al., 2024; Vovk et al., 2005; Lei et al., 2016) focused on the upper bound; however, the lower bound is also necessary for our framework.

Inverse Quantile: Given an $(1 - \alpha)$ -coverage level, we have that the $(1 - \alpha)$ -quantile of the calibration non-conformity scores is the appropriate threshold to achieve Equation 1. For our framework, given a threshold, λ , we recover the coverage level. This can be done using the **inverse λ -quantile**. Formally, if $(\mathbf{x}_{n+1}, y_{n+1})$ is a test point and $\mathcal{S}_{\text{calib}} = \{s(\mathbf{x}_i, y_i) \mid (\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{calib}}\}$, the inverse λ -quantile is given by

$$Q^{-1}(\lambda, \mathcal{S}_{\text{calib}}) := \Pr[s(\mathbf{x}_{n+1}, y_{n+1}) \leq \lambda] = \Pr[y_{n+1} \in \mathcal{C}_\lambda(\mathbf{x}_{n+1})].$$

Moreover, $Q^{-1}(\lambda, \mathcal{S}_{\text{calib}})$ is the coverage level for the label y_{n+1} . Lemma 3.2 asserts that the coverage level is within a bounded interval of length $\frac{1}{|\mathcal{D}_{\text{calib}}| + 1}$.

Lemma 3.2. *For $\lambda \in [0, 1]$ and $n = |\mathcal{D}_{\text{calib}}|$,*

$$\frac{\sum_{i=1}^n \mathbf{1}[s(\mathbf{x}_i, y_i) \leq \lambda]}{n + 1} \leq \Pr[y_{n+1} \in \mathcal{C}_\lambda(\mathbf{x}_{n+1})] \leq \frac{\sum_{i=1}^n \mathbf{1}[s(\mathbf{x}_i, y_i) \leq \lambda] + 1}{n + 1}, \quad (3)$$

CF for a Fixed Label: In standard CP, coverage is only evaluated for the true label, y_i . However, for fairness evaluation, it is essential to balance disparity between groups for **all** positive labels (see Table A1). So for conformal fairness evaluation, coverage needs to be balanced between groups for any given $\tilde{y} \in \mathcal{Y}^+$, as seen in Table 1. Lemma 3.3 asserts that we can perform CP using a fixed label and get the same coverage guarantees.

Lemma 3.3. *Equation 1 holds if we replace $\{(\mathbf{x}_i, y_i)\}$ with $\{(\mathbf{x}_i, \tilde{y})\}$ for a fixed $\tilde{y} \in \mathcal{Y}$.*

Connecting Theory to the Framework: For a particular fairness metric, we filter the calibration set based on the conditional from Table 1 and achieve bounds on the conditional coverage with Lemma 3.1. By Lemma 3.3, the bounds continue to hold when considering the conditional coverage for a fixed positive label. We use Lemma 3.2 to perform an inverse quantile to compute the coverage under various λ thresholds. With the coverages for a fixed positive label and each sensitive group, we compute the worst pairwise coverage gap across the groups using the bounds given by Lemma 3.3 to evaluate and control fairness at the desired closeness criterion.

3.3 CORE CONFORMAL FAIRNESS (CF) ALGORITHM

Input: The input to the core CF algorithm (Algorithm 1), include the calibration set, $\mathcal{D}_{\text{calib}}$, the set of labels (\mathcal{Y}) and positive labels (\mathcal{Y}^+), the set of sensitive groups, \mathcal{G} , a closeness criterion, c , the threshold search space, Λ , a fairness metric, M , and a corresponding filter function, F_M .

Specifying c : In practice, the choice of closeness criterion, c , may not be in the hands of the practitioner but instead defined by a (external) regulatory framework. For example, for Demographic Parity and c , we want that $\forall g_a, g_b \in \mathcal{G}$,

$$|\Pr[y_{n+1} \in \mathcal{C}_\lambda(\mathbf{x}_{n+1}) \mid \mathbf{x}_{n+1} \in g_a] - \Pr[y_{n+1} \in \mathcal{C}_\lambda(\mathbf{x}_{n+1}) \mid \mathbf{x}_{n+1} \in g_b]| < c.$$

Choosing Λ : The algorithm accepts a user-provided search space, Λ , which avoids degenerate thresholds and can guarantee desirable conditions. For our experiments, we set $\Lambda =$

$[\hat{q}(\alpha), \max\{\mathcal{S}_{\text{calib}}\}]$, ensuring that the optimal threshold, λ_{opt} , is at least $\hat{q}(\alpha)$. Since $\lambda_{\text{opt}} \geq \hat{q}(\alpha)$, the coverage increases for larger thresholds and still satisfies the $1 - \alpha$ coverage requirement. That is,

$$1 - \alpha \leq \Pr[y_{n+1} \in \mathcal{C}_{\hat{q}(\alpha)}(\mathbf{x}_{n+1})] \leq \Pr[y_{n+1} \in \mathcal{C}_{\lambda_{\text{opt}}}(\mathbf{x}_{n+1})].$$

Procedure: For each $\lambda \in \Lambda$, we want to check if it balances the coverage between groups for all positive labels. So, for each $(g, \tilde{y}) \in \mathcal{G} \times \mathcal{Y}^+$, we use F_M to filter $\mathcal{D}_{\text{calib}}$ (Line 11 in Algorithm 1) and then compute the non-conformity scores, $\mathcal{S}_{\text{calib}(g, \tilde{y})}$ (Line 12). With the inverse quantile, the coverage level is computed at the λ threshold on the scores (Line 14). We then compare the coverages for a fixed $y \in \mathcal{Y}^+$ between groups and check if the worst-case disparity satisfies the desired closeness criterion (Lines 16-21), forming the set Λ_M (Line 2). We choose $\lambda_{\text{opt}} = \min_{\lambda} \Lambda_M$ (Line 3) to minimize the final prediction set size (i.e., get the best efficiency). When evaluating multiple fairness metrics simultaneously, for example with Equalized Odds, the framework can be used to construct the set of satisfying lambdas for Equal Opportunity and Predictive Equality, Λ_{EO} and Λ_{PE} respectively. Then, $\lambda_{\text{opt}} = \min_{\lambda} \{\Lambda_{EO} \cap \Lambda_{PE}\}$.

Algorithm 1 Conformal Fairness Framework

```

1: procedure CONFORMAL_FAIRNESS( $\mathcal{D}_{\text{calib}}, \mathcal{Y}, \mathcal{Y}^+, \mathcal{G}, c, \Lambda, F_M$ )
2:    $\Lambda_M = \{\lambda \in \Lambda \mid \text{SATISFY\_LAMBDA}(\mathcal{D}_{\text{calib}}, \mathcal{Y}, \mathcal{Y}^+, \mathcal{G}, c, \lambda, F_M)\}$  ▷ Optimize for fairness
3:    $\lambda_{\text{opt}} = \min_{\lambda} \Lambda_M$  ▷ Optimize for efficiency
4:   return  $\lambda_{\text{opt}}$ 
5: end procedure
6:
7: procedure SATISFY_LAMBDA( $\mathcal{D}_{\text{calib}}, \mathcal{Y}, \mathcal{Y}^+, \mathcal{G}, c, \lambda, F_M$ )
8:   label_coverages =  $[0]_{(g_i, y) \in \mathcal{G} \times \mathcal{Y}}$ 
9:   interval_widths =  $[0]_{(g_i, y) \in \mathcal{G} \times \mathcal{Y}}$ 
10:  for  $(g, \tilde{y}) \in \mathcal{G} \times \mathcal{Y}^+$  do
11:     $\mathcal{D}_{\text{calib}(g, \tilde{y})} = \{(\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{calib}} \mid F_M(\mathbf{x}_i, y_i, g, \tilde{y}) = 1\}$ 
12:     $\mathcal{S}_{\text{calib}(g, \tilde{y})} = \{s(\mathbf{x}_i, y_i) \mid (\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{calib}(g, \tilde{y})}\}$ 
13:    interval_widths $[(g, \tilde{y})] = \frac{1}{|\mathcal{D}_{\text{calib}(g, \tilde{y})}| + 1}$  ▷ Uses Lemma 3.1
14:    label_coverages $[(g, \tilde{y})] = Q^{-1}(\lambda, \mathcal{S}_{\text{calib}(g, \tilde{y})})$  ▷ Uses Lemma 3.2
15:  end for
16:  for  $\tilde{y} \in \mathcal{Y}^+$  do ▷ Uses Lemma 3.3
17:     $\alpha_{\min} = \min(\text{label\_coverages}[(\cdot, \tilde{y})] - \text{interval\_widths}[(\cdot, \tilde{y})])$ 
18:     $\alpha_{\max} = \max(\text{label\_coverages}[(\cdot, y)])$ 
19:    if  $\alpha_{\max} - \alpha_{\min} > c$  then return False
20:  end if
21: end for
22: return True
23: end procedure

```

Using multiple λ thresholds: We also consider a classwise approach where we choose a $[\lambda_{\text{opt}}^0, \dots, \lambda_{\text{opt}}^{k-1}] = \boldsymbol{\lambda}_{\text{opt}} \in [0, 1]^K$ for each of the K classes. λ_{opt}^i is only required to satisfy the closeness criterion for the i^{th} class. One can achieve this by setting $\mathcal{Y}^+ = \{\tilde{y}\}$ and repeating Lines 2 and 3 in Algorithm 1 for each $\tilde{y} \in \mathcal{Y}^+$. This allows for smaller λ_{opt}^i to be chosen for most classes as they are no longer impacted by minority classes, which require a larger threshold to meet the closeness criterion.

A distinguishing feature of the CF framework is that it does not require group information at inference time. Though one can choose a different λ for each $(g, y) \in \mathcal{G} \times \mathcal{Y}$ pair, in streaming (or online) settings, the sensitive attribute may be unavailable. For example, loan applications may be race or gender-blind to enforce fairer judgment. In these settings, the CF Framework is not limited and provides group conditional coverage when group information is absent at inference time.

3.4 FRAMEWORK EXTENSIBILITY

Algorithm 1 directly applies to Demographic Parity, Equal Opportunity, Predictive Equality, and Equalized Odds. The following modifications are necessary to accommodate Disparate Impact, Predictive Parity, and some user-defined metrics.

Disparate Impact: The standard criterion for Disparate Impact is the *Four-Fifths Rule* applied to Demographic Parity. To control the conditional coverages for the Four-Fifths Rule, we only change Line 19 in Algorithm 1 to check if $(1 - \alpha_{\max})/(1 - \alpha_{\min}) < c$ for $c = 0.8$.

Predictive Parity: Predictive Parity seeks to balance the Positive Predictive Value (PPV) between groups (Verma & Rubin, 2018). It differs from the other fairness metrics in Table 1 as it is *conditioned on membership in the prediction set*. Given the objective of balancing conditional coverage, the conformal definition of Predictive Parity, and Bayes’ Theorem, we get

$$\Pr[Y = \tilde{y} \mid \tilde{y} \in \mathcal{C}_\lambda(X), X \in g_i] = \underbrace{\frac{\Pr[\tilde{y} \in \mathcal{C}_\lambda(X) \mid Y = \tilde{y}, X \in g_i]}{\Pr[\tilde{y} \in \mathcal{C}_\lambda(X) \mid X \in g_i]}}_{\text{Equal Opportunity over Demographic Parity}} \cdot \underbrace{\Pr[Y = \tilde{y} \mid X \in g_i]}_{\text{Conditional Label Probability}} \quad (4)$$

for $\tilde{y} \in \mathcal{Y}^+$ and $g_i \in \mathcal{G}$. A threshold, λ , is guaranteed to exist for any $\mathcal{Y}^+ \subseteq \mathcal{Y}$ if c is greater than the maximum pairwise total variation distance of the group-conditioned label distribution. This is formalized in Theorem 3.4.

Theorem 3.4. *Let W be a random variable for a label distribution over \mathcal{Y} . Let $W_i \sim W \mid (X \in g_i)$ – the label distribution conditioned on group membership. Then there exists λ such that for $c \geq \max\{D_{TV}(W_i, W_j) \mid i, j \in \{1, \dots, |\mathcal{G}|\}\}$, where D_{TV} is the total variation distance¹, the difference in Predictive Parity between groups is within c .*

In Equation 4, the Equal Opportunity, Demographic Parity, and Conditional Label Probability terms lie within finite intervals. This allows us to compute an interval where conformal Predictive Parity holds and use the CF framework to identify values of λ that meet the coverage closeness criterion. Further theoretical details and the proof of Theorem 3.4 are provided in Appendix C.

To control for arbitrarily small values of c , we use the *Predictive Parity Proxy*—an example of a user-defined metric—defined in Equation 5. For all $g_i \in \mathcal{G}, \tilde{y} \in \mathcal{Y}^+$,

$$\Pr(Y = \tilde{y} \mid \tilde{y} \in \mathcal{C}_\lambda(X), X \in g_i) - \Pr(Y = \tilde{y} \mid X \in g_i). \quad (5)$$

In cases where it is possible to assume the label distribution is independent of group membership, Equation 4 can be directly controlled for an arbitrarily small closeness criterion, c . Proofs and technical details on these modifications can be found in Appendix C.

3.5 LEVERAGING THE CF FRAMEWORK FOR FAIRNESS AUDITING

Using the Conformal Fairness Framework, one can audit if the disparity of a conformal predictor between multiple groups violates a user-specified fairness criterion. Specifically, we have thus far focused on fairness criteria concerning bounding the disparity between groups using the fairness metrics described in Table 1 by some closeness criterion, c . It is straightforward to support user-defined fairness metrics concerning label coverage. While Algorithm 1, as presented, gives a method of finding an optimal λ threshold which satisfies the fairness guarantees using Lemmas 3.1, 3.2, and 3.3, the same SATISFY_LAMBDA procedure can be leveraged to check if a *given* λ used by a conformal predictor satisfies the same fairness guarantees. Notably, the CF framework can also be leveraged even if the conformal predictor is treated as a black-box model. In this case, we construct an $\mathcal{D}_{\text{audit}}$ set exchangeable with the calibration data used for the conformal predictor. Using $\mathcal{D}_{\text{audit}}$, we can determine if the conformal predictor satisfies the corresponding fairness guarantee given the fairness metric and the λ threshold used.

3.6 NON-CONFORMITY SCORES

There are several choices for the non-conformity score for performing fair conformal prediction with classification tasks. We currently implement TPS (Sadinle et al., 2019), APS (Romano et al., 2020), RAPS (Angelopoulos et al., 2022), DAPS (H. Zargarbashi et al., 2023), and CFGNN (Huang et al., 2024) in the CF framework, though any non-conformity score can be used. More details on the specifics of each non-conformity score can be found in Appendix D.2.

¹A *modified total variation distance*, $D_{TV}^+(W_i, W_j) := \sup_{k \in \mathcal{Y}^+} |\Pr[W_i = k] - \Pr[W_j = k]|$, can be used in place of D_{TV} in Theorem 3.4 for a weaker assumption about c , which still gives a satisfying λ .

4 EXPERIMENTS

4.1 SETUP

Datasets: To evaluate the CF Framework, we used five multi-class datasets: Pokec-n (Takac & Zabovsky, 2012), Pokec-z (Takac & Zabovsky, 2012), Credit (Agarwal et al., 2021), ACSIncome (Ding et al., 2021), and ACSEducation (Ding et al., 2021) (see Table 2 for details). For each dataset, we use a 30%/20%/25%/25% stratified split of the labeled points for $\mathcal{D}_{\text{train}}/\mathcal{D}_{\text{valid}}/\mathcal{D}_{\text{calib}}/\mathcal{D}_{\text{test}}$.

Table 2: Dataset Statistics. T refers to Tabular, and G refers to Graph.

Name	Type	Size	# Labeled	# Groups	# Classes
ACSIncome	T	1,664,500	ALL	race(9)	4
ACSEducation	T	1,664,500	ALL	race(9)	6
Name	Type	($ \mathcal{V} , \mathcal{E} $)	# Labeled	# Groups	# Classes
Credit	T/G	(30,000, 1,436,858)	ALL	age(2)	4
Pokec-n	G	(66,569, 729,129)	8,797	region(2), gender(2)	4
Pokec-z	G	(66,569, 729,129)	8,797	region(2), gender(2)	4

Models: For the graph datasets, we evaluated with GCN (Kipf & Welling, 2017), GraphSAGE (Hamilton et al., 2017), or GAT (Veličković et al., 2018) as the base model (results reported are for the highest performing base model). For Credit, we additionally evaluated XGBoost (Chen & Guestrin, 2016) (i.e., ignoring the graph structure) as we empirically observed this approach to outperform the graph neural network baselines in terms of efficiency for this dataset. The choice of ignoring edge information while training Credit on XGBoost does not prohibit us from using CFGNN or DAPS, which utilize the edge information. The conformal predictor requires the softmax logits from the base model (i.e., XGBoost) but is otherwise model agnostic. For ACSIncome and ACSEducation, we used an XGBoost model. Each model’s hyperparameters were tuned as discussed in Appendix D.3.

Baseline: For each dataset and CP non-conformity score, we built a conformal predictor to achieve a coverage level of $1 - \alpha = 0.9$. Then, we assess fairness according to the specific fairness metric using the SATISFY_LAMBDA from Algorithm 1 for $\lambda = \hat{q}(\alpha)$.

Evaluation Metrics: We report the *worst fairness disparity* and *efficiency*. For Disparate Impact, the worst fairness disparity is the *minimum* $(1 - \alpha_{\max})/(1 - \alpha_{\min})$ across the positive labels. For the remaining metrics, we record the *maximum* $\alpha_{\max} - \alpha_{\min}$ across the positive labels.

4.2 RESULTS

For each figure, we use a line to indicate the base conformal predictor’s *average* worst fairness disparity across different thresholds, the bar plot for the worst fairness disparity using the CF Framework, and a dot to denote the desired fairness disparity. We report the average base performance for simplicity and readability of the figures. In every experiment, except for Figure 2, the CF framework was **better** than the average base conformal predictor. We provide a more granular version of Figure 2 with Figure E4, where it is clear that the framework performs better for every closeness threshold.

Controlling for Fairness Disparity: For different closeness thresholds, our CF Framework effectively controls the fairness disparity for several metrics compared to the base conformal predictor. In Figure 1 and 2, we can observe that in terms of fairness disparity, our CF Framework **precisely** (note step-wise change with c on violations) improves upon the baseline conformal predictor. As with algorithmic fairness, a trade-off is involved in that there is a slightly worse efficiency. From Figure 2, we continue to observe this for both standard and graph-based conformal predictors. Furthermore, if the base conformal predictor is already “fair” according to our fairness disparity criterion, then the CF Framework will report the results accordingly. This phenomenon is observed with the CFGNN results in Figure 2, where the CF Framework matches the baseline regarding both evaluation metrics. This behavior of the CF Framework makes it suitable to leverage for black box fairness auditing (as noted previously). We present additional results, for example, the disparity results for the CF

Framework without classwise lambdas in Appendix E. Notably, the prediction set sizes are more prominent due to selecting a larger λ than the classwise approach (see Figure E3 vs Figure 1).

Controlling for Disparate Impact: For Disparate Impact, we present results for the standard *Four-Fifths Rule*. In Table 3, we see that using the CF Framework can significantly improve the base conformal predictor for the *Four-Fifths Rule*. The disparate impact value is far below the desired 0.8 for the base conformal predictor, sometimes even less than 0.4, as with Credit with TPS and ACSIncome dataset. Our framework, however, is close to the 0.8 value and in some cases surpasses it, like in Credit with CFGNN, with minor effects on the efficiency for both datasets.

Table 3: *Four-Fifths Rule* for Credit and ACSIncome. Our framework surpasses the base conformal predictor and achieves close to or exceeds the disparate impact value of 0.80. The - means N/A.

		APS		RAPS		TPS		CFGNN		DAPS	
		Base	CF	Base	CF	Base	CF	Base	CF	Base	CF
Credit	Disp. Impact	0.646	0.821	0.586	0.768	0.252	0.793	0.922	0.922	0.539	0.809
	Efficiency	2.326	2.513	2.326	2.509	2.268	2.558	2.202	2.202	2.254	2.526
ACSIncome	Disp. Impact	0.397	0.797	0.387	0.790	0.356	0.798	-	-	-	-
	Efficiency	2.212	2.674	2.169	2.752	2.109	2.679	-	-	-	-

Agnostic to Non-Conformity Score: As discussed earlier, the CF Framework can support a variety of non-conformity scores, emphasizing the agnostic nature of our framework. We achieved effective results for conformal predictors with different underlying non-conformity score functions for all the experiments. Further results can be found in Appendix E.

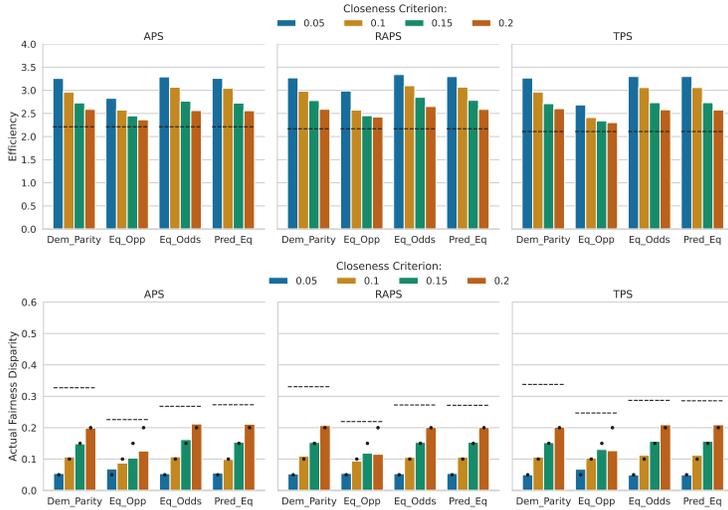


Figure 1: **ACSIncome**. The left two plots are efficiency results, while the right two are the fairness disparities for (a) APS, (b) RAPS, and (c) TPS. In all cases, our framework gives results at or better than the desired threshold and better than the baseline.

Intersectional Fairness: When characterizing data points into groups, we are not limited to a single sensitive attribute. In many applications, there can be multiple sensitive attributes (e.g., race and gender) that need to be considered. Our CF Framework is not limited to analyzing a single sensitive attribute. To demonstrate this, we experiment with the Pokec-n dataset. Pokec-n has two sensitive attributes, namely *region* and *gender*. We treat each combination of region and gender as a separate sensitive group and apply the CF framework to control for fairness disparities. Figure 3 shows that the CF framework improves upon the base conformal predictor regarding fairness disparity. This improvement is starker with the graph-based conformal predictors, CFGNN, and DAPS, as seen in Figure 3 plots (b) and (c).

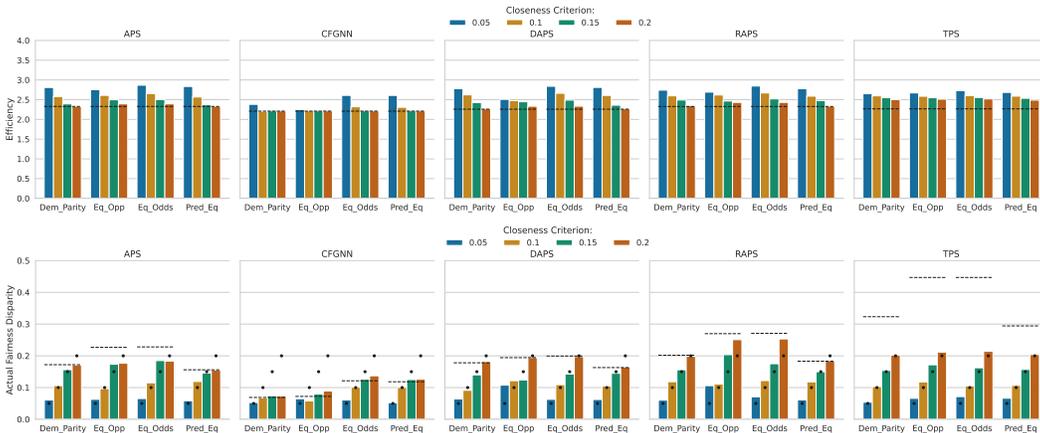


Figure 2: **Credit**. The top four plots are efficiency results, while the bottom four are the fairness disparities for (a) APS, (b) CFGNN, (c) DAPS, (d) RAPS, and (e) TPS. In all cases, our framework achieves the desired coverage gap better than the baseline, with a minor impact on efficiency.

A key challenge with intersectional fairness is the multiplicative increase in the number of groups (i.e., combinations of sensitive attributes and classes) that must be calibrated and evaluated. This increases the data requirements needed to satisfy the coverage guarantees discussed in Section 3.2, as these guarantees become harder to achieve when the size of $\mathcal{D}_{(g,y)}$ decreases. This problem is exacerbated (in empirical results) for datasets with only a few labeled points, such as Pokec-n. For Pokec-n, using a standard data split, the calibration set has around 2200 data points. The calibration set is then further split to get the conditional positive label coverage for each positive label and group pair. This results in the calibration being done with sets of fewer than a few hundred points, which is much lower than the suggested 1000 points in the literature (Angelopoulos & Bates, 2021). In Figure 3, the effect of this challenge is seen with the fairness disparity given by the CF Framework being slightly above the desired closeness threshold for $c = 0.1$. Nevertheless, the guarantees still hold, even under intersectional fairness constraints.

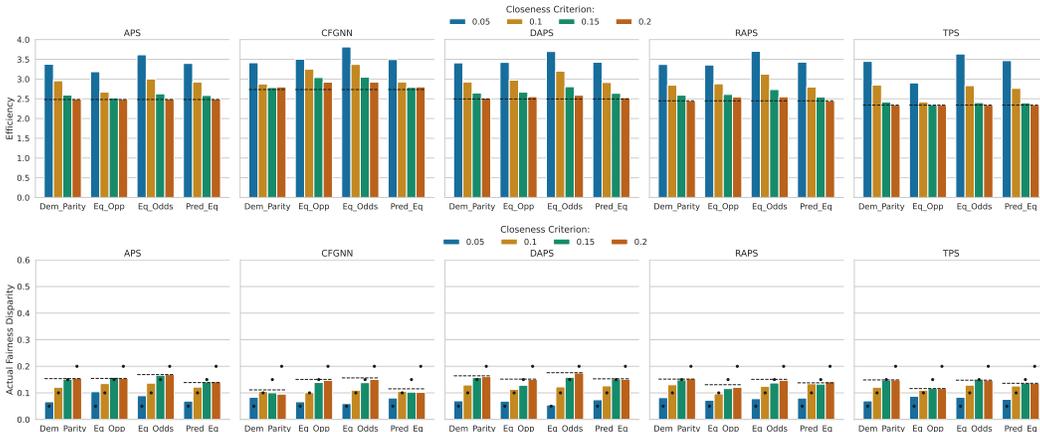


Figure 3: **Pokec-n** using **both** sensitive attributes. The top four plots are the efficiency results, while the bottom four are the fairness disparities for (a) APS, (b) CFGNN, (c) DAPS, (d) RAPS, and (e) TPS. CFGNN (b) and DAPS (c) achieve the desired fairness coverage thresholds better than standard CP methods.

Predictive Parity Proxy: As discussed, the CF framework is extensible to user-defined fairness notions. We consider the Predictive Parity Proxy in Equation 5 as an example of a user’s ability to provide a reasonable fairness measure (Disparate Impact, above, is another example). An experiment on ACSEducation in Table 4 demonstrates we can control for arbitrarily small values of c , unlike the standard notion of Predictive Parity. Additionally, it empirically illustrates that we can control

for disparities of probabilities conditioned on the prediction set. This metric can also be applied in the graph setting, as seen in Appendix E.

Table 4: **ACSEducation**. The worst-case fairness disparity, based on the Predictive Parity Proxy, with our method is below the desired c threshold, while the *average* baseline disparity is much higher (> 0.30) than all of the c thresholds we consider.

		Closeness Threshold (c)					
		0.05	0.10	0.15	0.20	Base (Average)	
APS	Max Fairness Disparity	0.044	0.093	0.152	0.166	0.411	
	Efficiency	3.662	3.236	3.049	3.008	2.982	
RAPS	Max Fairness Disparity	0.043	0.094	0.153	0.172	0.268	
	Efficiency	3.948	3.339	3.102	3.063	3.030	
TPS	Max Fairness Disparity	0.038	0.091	0.167	0.199	0.319	
	Efficiency	3.662	3.061	2.880	2.845	2.828	

4.3 DISCUSSION AND RELATED WORK

Few prior efforts study fairness and conformal prediction (Wang et al., 2024; Lu et al., 2022; Liu et al., 2022). One line of work has focused on applying fairness notions toward CP problems for regression tasks, explicitly focusing on Demographic Parity (Liu et al., 2022) and Equal Opportunity (Wang et al., 2024). Another line of work focuses on applying the notion of Overall Accuracy Equality for CP (Lu et al., 2022). This effort considers a specific medical application of detecting malignant skin conditions and applies group-balanced CP (Vovk, 2012).

An orthogonal direction is on (group) conditional CP. Foygel Barber et al. (2021) provides a theoretical grounding for conditional CP, while Gibbs et al. (2023) considers the impact of covariate shift for conditional coverage under the I.I.D assumption. Others Bastani et al. (2022); Jung et al. (2023) look at multivald CP, which requires (1) group-conditional and (2) threshold-calibrated coverage guarantees – a distinct notion from Conformal Fairness. Deng et al. (2023) also introduces a generalization for multi-calibration and how it relates to algorithmic fairness for conformal prediction, focusing on equalized coverage for regression.

Our work differs in its breadth and flexibility (i.e., support for several fairness metrics and non-conformity scores) and focus on classification. Some existing works represent specific instantiations in our framework (e.g., Wang et al. (2024)). Others provide baselines for comparison (e.g., BatchGCP (Jung et al., 2023)), without our framework’s theoretical guarantees for fairness. Appendix E.6 contains more details on BatchGCP and results. The CF framework generalizes group-balanced CP to consider the notion of coverage for specific labels, thus allowing us to evaluate disparity based on classical fairness metrics in a manner that does not require *a priori* knowledge of group membership at inference time, unlike many approaches listed above.

Lastly, CF can be used in fairness-critical domains where conditional conformal prediction is infeasible, such as finance, which can have strict fairness requirements (Agarwal et al., 2021), and health care (Wang et al., 2024), where privileged information may be unavailable at inference time.

5 CONCLUSION

In this work, we formalize the notion of Conformal Fairness (CF) for conformal predictors and propose a novel and comprehensive CF Framework. We provide a theoretically grounded algorithm that can be used to control for the gaps in conditional coverage, defined based on different fairness metrics, across sensitive groups. We conduct experiments on tabular and graph datasets, leveraging the exchangeability assumption of conformal prediction. We present results for CF based on various classical and user-defined fairness metrics on conformal predictors with various non-conformity score functions, including results on the framework’s effectiveness in evaluating intersectional fairness with conformal predictors. We further describe how the CF framework can be practically leveraged for applications, including fairness auditing of conformal predictors. Future work could extend the framework to regression tasks and strengthen the theory by relaxing assumptions and exploring non-exchangeable settings.

ACKNOWLEDGMENTS

The authors acknowledge support from National Science Foundation (NSF) grant #2112471 (AI-EDGE) and a grant from Cisco Research (US202581249). Any opinions and findings are those of the author(s) and do not necessarily reflect the views of the granting agencies. The authors also thank the anonymous reviewers for their constructive feedback on this work.

REPRODUCIBILITY STATEMENT

In the spirit of reproducibility, the proofs for the theoretical aspects of our work can be found in Appendices B and C. Details of the experiments, including datasets, non-conformity scores, and hyperparameters, are provided in Appendix D. The source code is available at <https://github.com/AdityaVadlamani/conformal-fairness>.

REFERENCES

- Chirag Agarwal, Himabindu Lakkaraju, and Marinka Zitnik. Towards a unified framework for fair and stable graph representation learning. In *Uncertainty in Artificial Intelligence*, pp. 2114–2124. PMLR, 2021.
- Wael Alghamdi, Hsiang Hsu, Haewon Jeong, Hao Wang, Peter Michalak, Shahab Asodeh, and Flavio Calmon. Beyond adult and compas: Fair multi-class prediction via information projection. *Advances in Neural Information Processing Systems*, 35:38747–38760, 2022.
- Anastasios Angelopoulos, Stephen Bates, Jitendra Malik, and Michael I. Jordan. Uncertainty sets for image classifiers using conformal prediction, 2022. URL <https://arxiv.org/abs/2009.14193>.
- Anastasios N. Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *CoRR*, abs/2107.07511, 2021. URL <https://arxiv.org/abs/2107.07511>.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning: Limitations and Opportunities*. MIT Press, 2023.
- Osbert Bastani, Varun Gupta, Christopher Jung, Georgy Noarov, Ramya Ramalingam, and Aaron Roth. Practical adversarial multivald conformal prediction. *Advances in Neural Information Processing Systems*, 35:29362–29373, 2022.
- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’16, pp. 785–794, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450342322. doi: 10.1145/2939672.2939785. URL <https://doi.org/10.1145/2939672.2939785>.
- John Cherian and Lenny Bronner. How the washington post estimates outstanding votes for the 2020 presidential election. *Retrieved September*, 13:2023, 2020.
- Elliot Creager, David Madras, Jörn-Henrik Jacobsen, Marissa Weis, Kevin Swersky, Toniann Pitassi, and Richard Zemel. Flexibly fair representation learning by disentanglement. In *International conference on machine learning*, pp. 1436–1445. PMLR, 2019.
- Sean Current, Yuntian He, Saket Gurukar, and Srinivasan Parthasarathy. Fairegm: Fair link prediction and recommendation via emulated graph modification. In *Proceedings of the 2nd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, EAAMO ’22, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450394772. doi: 10.1145/3551624.3555287. URL <https://doi.org/10.1145/3551624.3555287>.
- Zhun Deng, Cynthia Dwork, and Linjun Zhang. Happymap: A generalized multicalibration method. In *14th Innovations in Theoretical Computer Science Conference (ITCS 2023)*, 2023.
- Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair machine learning. *Advances in neural information processing systems*, 34:6478–6490, 2021.

- Tiffany Ding, Anastasios N. Angelopoulos, Stephen Bates, Michael I. Jordan, and Ryan J. Tibshirani. Class-conditional conformal prediction with many classes. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA, 2024. Curran Associates Inc.
- Yushun Dong, Jing Ma, Song Wang, Chen Chen, and Jundong Li. Fairness in graph mining: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 35(10):10583–10602, 2023.
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pp. 214–226, 2012.
- The U.S. EEOC. Uniform guidelines on employee selection procedures. March 1979.
- Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '15*, pp. 259–268, New York, NY, USA, 2015. Association for Computing Machinery. ISBN 9781450336642. doi: 10.1145/2783258.2783311. URL <https://doi.org/10.1145/2783258.2783311>.
- Rina Foygel Barber, Emmanuel J Candes, Aaditya Ramdas, and Ryan J Tibshirani. The limits of distribution-free conditional predictive inference. *Information and Inference: A Journal of the IMA*, 10(2):455–482, 2021.
- Bishwamitra Ghosh, Debabrota Basu, and Kuldeep S. Meel. Justicia: A stochastic sat approach to formally verify fairness. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(9):7554–7563, May 2021. doi: 10.1609/aaai.v35i9.16925. URL <https://ojs.aaai.org/index.php/AAAI/article/view/16925>.
- Isaac Gibbs, John J Cherian, and Emmanuel J Candès. Conformal prediction with conditional guarantees. *arXiv preprint arXiv:2305.12616*, 2023.
- Soroush H. Zargarbashi, Simone Antonelli, and Aleksandar Bojchevski. Conformal prediction sets for graph neural networks. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 12292–12318. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/h-zargarbashi23a.html>.
- William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. *CoRR*, abs/1706.02216, 2017. URL <http://arxiv.org/abs/1706.02216>.
- Yuntian He, Saket Gurukar, and Srinivasan Parthasarathy. Fairmile: Towards an efficient framework for fair graph representation learning. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, pp. 1–10, 2023.
- Dennis Hirsch, Timothy Bartley, Aravind Chandrasekaran, Davon Norris, Srinivasan Parthasarathy, and Piers Norris Turner. *Business Data Ethics: Emerging Models for Governing AI and Advanced Analytics*. Springer Nature, 2023.
- Kexin Huang, Ying Jin, Emmanuel Candès, and Jure Leskovec. Uncertainty quantification over graph with conformalized graph neural networks. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA, 2024. Curran Associates Inc.
- Christopher Jung, Georgy Noarov, Ramya Ramalingam, and Aaron Roth. Batch multivald conformal prediction. *11th International Conference on Learning Representations (ICLR)*, 2023.
- Nathan Kallus and Angela Zhou. Residual unfairness in fair machine learning from prejudiced data. In *International Conference on Machine Learning*, pp. 2439–2448. PMLR, 2018.

- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=SJU4ayYgl>.
- Jing Lei, Max Grazier G’Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry A. Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113:1094–1111, 2016. URL <https://api.semanticscholar.org/CorpusID:13741419>.
- Richard Liaw, Eric Liang, Robert Nishihara, Philipp Moritz, Joseph E Gonzalez, and Ion Stoica. Tune: A research platform for distributed model selection and training. *arXiv preprint arXiv:1807.05118*, 2018.
- Meichen Liu, Lei Ding, Dengdeng Yu, Wulong Liu, Linglong Kong, and Bei Jiang. Conformalized fairness via quantile regression. *Advances in Neural Information Processing Systems*, 35:11561–11572, 2022.
- Tianci Liu, Haoyu Wang, Yaqing Wang, Xiaoqian Wang, Lu Su, and Jing Gao. Simfair: A unified framework for fairness-aware multi-label classification. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(12):14338–14346, Jun. 2023. doi: 10.1609/aaai.v37i12.26677. URL <https://ojs.aaai.org/index.php/AAAI/article/view/26677>.
- Charles Lu, Andréanne Lemay, Ken Chang, Katharina Höbel, and Jayashree Kalpathy-Cramer. Fair conformal predictors for applications in medical imaging. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(11):12008–12016, Jun. 2022. doi: 10.1609/aaai.v36i11.21459. URL <https://ojs.aaai.org/index.php/AAAI/article/view/21459>.
- Pranav Maneriker, Codi Burley, and Srinivasan Parthasarathy. Online fairness auditing through iterative refinement. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD ’23*, pp. 1665–1676, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701030. doi: 10.1145/3580305.3599454. URL <https://doi.org/10.1145/3580305.3599454>.
- Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM computing surveys (CSUR)*, 54(6):1–35, 2021.
- Yaniv Romano, Evan Patterson, and Emmanuel Candes. Conformalized quantile regression. *Advances in neural information processing systems*, 32, 2019.
- Yaniv Romano, Matteo Sesia, and Emmanuel Candes. Classification with valid and adaptive coverage. *Advances in neural information processing systems*, 33:3581–3591, 2020.
- Julien Rouzot, Julien Ferry, and Marie-José Huguet. Learning optimal fair scoring systems for multi-class classification. In *2022 IEEE 34th International Conference on Tools with Artificial Intelligence (ICTAI)*, pp. 197–204, 2022. doi: 10.1109/ICTAI56018.2022.00036.
- Mauricio Sadinle, Jing Lei, and Larry Wasserman. Least ambiguous set-valued classifiers with bounded error levels. *Journal of the American Statistical Association*, 114(525):223–234, 2019.
- Lubos Takac and Michal Zabovsky. Data analysis in public social networks. In *International scientific conference and international workshop present day trends of innovations*, volume 1, 2012.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks, 2018. URL <https://arxiv.org/abs/1710.10903>.
- Sahil Verma and Julia Rubin. Fairness definitions explained. In *Proceedings of the international workshop on software fairness*, pp. 1–7, 2018.
- Vladimir Vovk. Conditional validity of inductive conformal predictors. In Steven C. H. Hoi and Wray Buntine (eds.), *Proceedings of the Asian Conference on Machine Learning*, volume 25 of *Proceedings of Machine Learning Research*, pp. 475–490, Singapore Management University, Singapore, 04–06 Nov 2012. PMLR. URL <https://proceedings.mlr.press/v25/vovk12.html>.

- Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a random world*, volume 29. Springer, 2005.
- Fangxin Wang, Lu Cheng, Ruocheng Guo, Kay Liu, and Philip S. Yu. Equal opportunity of coverage in fair regression. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA, 2024. Curran Associates Inc.
- Ziyan Wang and Hao Wang. Variational imbalanced regression: Fair uncertainty quantification via probabilistic smoothing. *Advances in Neural Information Processing Systems*, 36, 2024.
- Tom Yan and Chicheng Zhang. Active fairness auditing. In *International Conference on Machine Learning*, pp. 24929–24962. PMLR, 2022.
- I-Cheng Yeh and Che-hui Lien. The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. *Expert systems with applications*, 36(2): 2473–2480, 2009.
- Chen Zhao, Le Wu, Pengyang Shao, Kun Zhang, Richang Hong, and Meng Wang. Fair representation learning for recommendation: A mutual information perspective. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(4):4911–4919, Jun. 2023. doi: 10.1609/aaai.v37i4.25617. URL <https://ojs.aaai.org/index.php/AAAI/article/view/25617>.

A FAIRNESS METRICS

As discussed in Section 2.2, exact fairness is difficult to achieve. So, for many metrics, we want to say something about the difference between groups. Formally, for Demographic Parity, we may require that for some (small) $c \in (0, 1]$,

$$\max_{\tilde{y} \in \mathcal{Y}^+} \left\{ \left| \Pr[\hat{Y} = \tilde{y} \mid X \in g_a] - \Pr[\hat{Y} = \tilde{y} \mid X \in g_b] \right| \mid \forall g_a, g_b \in \mathcal{G} \right\} < c.$$

Similar requirements exist for the other fairness metrics.

Table A1: Fairness metrics formulations for multiclass classification (Rouzot et al., 2022).

Metric	Definition
Demographic (or Statistical) Parity	$\Pr[\hat{Y} = y \mid X \in g_a] = \Pr[\hat{Y} = y \mid X \in g_b], \forall g_a, g_b \in \mathcal{G}, \forall y \in \mathcal{Y}^+$
Equal Opportunity	$\Pr[\hat{Y} = y \mid Y = y, X \in g_a] = \Pr[\hat{Y} = y \mid Y = y, X \in g_b], \forall g_a, g_b \in \mathcal{G}, \forall y \in \mathcal{Y}^+$
Predictive Equality	$\Pr[\hat{Y} = y \mid Y \neq y, X \in g_a] = \Pr[\hat{Y} = y \mid Y \neq y, X \in g_b], \forall g_a, g_b \in \mathcal{G}, \forall y \in \mathcal{Y}^+$
Equalized Odds	Equal Opp. and Pred. Equality
Predictive Parity	$\Pr[Y = y \mid \hat{Y} = y, X \in g_a] = \Pr[Y = y \mid \hat{Y} = y, X \in g_b], \forall g_a, g_b \in \mathcal{G}, \forall y \in \mathcal{Y}^+$

B PROOFS

B.1 PROOF OF LEMMA 3.1

Lemma 3.1. *For any $(g, \tilde{y}) \in \mathcal{G} \times \mathcal{Y}^+$, calibrating on $\mathcal{D}_{\text{calib}(g, \tilde{y})} = \{(\mathbf{x}_i, y_i) \mid F_M(\mathbf{x}_i, y_i, g, \tilde{y}) = 1\}$ guarantees the following about the conditional coverage:*

$$1 - \alpha \leq \Pr[y_{n+1} \in \mathcal{C}_\lambda(\mathbf{x}_{n+1}) \mid F_M(\mathbf{x}_{n+1}, y_{n+1}, g, \tilde{y}) = 1] \leq 1 - \alpha + \frac{1}{|\mathcal{D}_{\text{calib}(g, y)}| + 1} \quad (2)$$

The interval width is $\frac{1}{|\mathcal{D}_{\text{calib}(g, y)}| + 1}$.

Proof. The proof for the upper bound follows Ding et al. (2024), where the test score $s(\mathbf{x}_{n+1}, y_{n+1})$ follows the distribution of scores in $\mathcal{D}_{\text{calib}} \cap \mathcal{R}$, thus holding the conditional coverage. The upper-bound guarantee directly follows Romano et al. (2019), assuming distinct non-conformity scores or a suitably random way to tie-break equal scores. \square

B.2 PROOF OF LEMMA 3.2

To prove Lemma 3.2, we use the following Lemma that is stated in the proof of Theorem D.1 from Angelopoulos & Bates (2021):

Lemma B.1. *Suppose you have $n + 1$ exchangeable random variables $Z_1, Z_2, \dots, Z_n, Z_{n+1}$. Then, $\Pr[Z_{n+1} \leq Z_{(k)}] = \frac{k}{n+1}$ for all $k \leq n$, where Z_{n+1} is the test point.*

Proof Sketch. Using Z_1, \dots, Z_n , you can form $n + 1$ intervals, $(L, Z_{(1)}], (Z_{(1)}, Z_{(2)}], \dots, (Z_{(n-1)}, Z_{(n)}], (Z_{(n)}, U]$, where L and U are the lower and upper bounds of the random variables. Exchangeability gives us that Z_{n+1} falls in any of the $n + 1$ intervals with equal probability. Thus, $\Pr[Z_{n+1} \leq Z_{(k)}]$ is the probability that it falls in the first k intervals giving us $\Pr[Z_{n+1} \leq Z_{(k)}] = \frac{k}{n+1}$. \square

Now using Lemma B.1, we will prove Lemma 3.2.

Lemma 3.2. *For $\lambda \in [0, 1]$ and $n = |\mathcal{D}_{\text{calib}}|$,*

$$\frac{\sum_{i=1}^n \mathbf{1}[s(\mathbf{x}_i, y_i) \leq \lambda]}{n+1} \leq \Pr[y_{n+1} \in \mathcal{C}_\lambda(\mathbf{x}_{n+1})] \leq \frac{\sum_{i=1}^n \mathbf{1}[s(\mathbf{x}_i, y_i) \leq \lambda] + 1}{n+1}, \quad (3)$$

Proof. Let $s_i = s(\mathbf{x}_i, y_i)$ for brevity. By Lemma B.1, we have that $\Pr(s_{n+1} \leq s_{(k)}) = \frac{k}{n+1}$ for all $k \leq n$. Let k be s.t. $s_{(k)} \leq \lambda \leq s_{(k+1)}$. We then have that $\Pr(s_{n+1} \leq s_{(k)}) \leq \Pr(s_{n+1} \leq \lambda) \leq \Pr(s_{n+1} \leq s_{(k+1)})$. So,

$$\frac{k}{n+1} \leq \Pr(s_{n+1} \leq \lambda) \leq \frac{k+1}{n+1}.$$

Since s_1, \dots, s_n are known empirically, we have that $k = \sum_{i=1}^n \mathbf{1}[s_i \leq \lambda]$. Hence,

$$\frac{\sum_{i=1}^n \mathbf{1}[s_i \leq \lambda]}{n+1} \leq \Pr(s_{n+1} \leq \lambda) \leq \frac{\sum_{i=1}^n \mathbf{1}[s_i \leq \lambda] + 1}{n+1}.$$

Since $\Pr[s_{n+1} \leq \lambda] = \Pr[y_{n+1} \in \mathcal{C}_\lambda(\mathbf{x}_{n+1})]$, we get Equation 3. □

B.3 PROOF OF LEMMA 3.3

Lemma 3.3. *Equation 1 holds if we replace $\{(\mathbf{x}_i, y_i)\}$ with $\{(\mathbf{x}_i, \tilde{y})\}$ for a fixed $\tilde{y} \in \mathcal{Y}$.*

Proof. By computing $\hat{q}(\alpha) = \text{Quantile}\left(\frac{\lceil (n+1)(1-\alpha) \rceil}{n}; \{s(\mathbf{x}_i, \tilde{y})\}_{i=1}^n\right)$, where \tilde{y} is the fixed label. Using the assumption of exchangeability we have that $s(\mathbf{x}_1, \tilde{y}), s(\mathbf{x}_2, \tilde{y}), \dots, s(\mathbf{x}_{n+1}, \tilde{y})$ is an exchangeable sequence. Thus

$$1 - \alpha \leq \Pr[s(\mathbf{x}_{n+1}, \tilde{y}) \leq \hat{q}(\alpha)] \leq 1 - \alpha + \frac{1}{n+1}. \quad (6)$$

Equivalently,

$$1 - \alpha \leq \Pr[\tilde{y} \in \mathcal{C}_{\hat{q}(\alpha)}(\mathbf{x}_{n+1})] \leq 1 - \alpha + \frac{1}{n+1}. \quad (7)$$

□

C FURTHER DISCUSSION ON PREDICTIVE PARITY

Theorem 3.4. *Let W be a random variable for a label distribution over \mathcal{Y} . Let $W_i \sim W | (X \in g_i)$ – the label distribution conditioned on group membership. Then there exists λ such that for $c \geq \max\{D_{TV}(W_i, W_j) \mid i, j \in \{1, \dots, |\mathcal{G}|\}\}$, where D_{TV} is the total variation distance², the difference in Predictive Parity between groups is within c .*

Proof. Let λ be the maximum value it can take on as the threshold of the prediction set. Let $y \in \mathcal{Y}^+$ and $g_m, g_n \in \mathcal{G}$. Then,

$$\begin{aligned} \Pr[Y = y \mid y \in \mathcal{C}_\lambda(X), X \in g_m] &= \frac{\Pr[y \in \mathcal{C}_\lambda(X) \mid Y=y, X \in g_m]}{\Pr[y \in \mathcal{C}_\lambda(X) \mid X \in g_m]} \Pr[Y = y \mid X \in g_m] \\ &= \Pr[Y = y \mid X \in g_m] \end{aligned}$$

since the numerator and the denominator are 1 due to the selection of λ . Using the same argument for g_n and the definition of D_{TV} the difference in Predictive Parities is:

$$\left| \Pr[Y = y \mid X \in g_m] - \Pr[Y = y \mid X \in g_n] \right| \leq D_{TV}(W_m, W_n)$$

Since, $D_{TV}(W_m, W_n) \leq \max\{D_{TV}(W_i, W_j) \mid i, j \in \{1, \dots, |\mathcal{G}|\}\} \leq c$, the chosen λ value causes the Predictive Parity difference to be less than c . Thus a solution exists for $c \geq \max\{D_{TV}(W_i, W_j) \mid i, j \in \{1, \dots, |\mathcal{G}|\}\}$. □

²A modified total variation distance, $D_{TV}^+(W_i, W_j) := \sup_{k \in \mathcal{Y}^+} |\Pr[W_i = k] - \Pr[W_j = k]|$, can be used in place of D_{TV} in Theorem 3.4 for a weaker assumption about c , which still gives a satisfying λ .

C.1 ACHIEVING ARBITRARY CLOSENESS

We will further discuss the two methods for controlling for arbitrarily small values of c , depending on whether the label distribution is independent of group membership.

C.1.1 LABEL DISTRIBUTION INDEPENDENT OF GROUP MEMBERSHIP

Assuming independence, we have the following corollary of Theorem 3.4:

Corollary C.1. *Given a random variable, W , from a label distribution over \mathcal{Y} that is independent of group membership, then the Conformal Fairness Framework can find a λ such that the disparity of Predictive Parity is within any $c > 0$.*

Proof. Let $W_i \sim W | (X \in g_i)$ be the label distribution condition on group membership. Using the independence assumption we get $W = W_i, \forall i \in \{1, \dots, |\mathcal{G}|\}$, thus $D_{TV}(W_i, W_j) = 0, \forall i, j \in \{1, \dots, |\mathcal{G}|\}$. Hence, using Theorem 3.4, $c > \max\{D_{TV}(W_i, W_j) | i, j \in \{1, \dots, |\mathcal{G}|\}\} = 0$. Thus a λ can be found for any $c > 0$. \square

C.1.2 WITHOUT INDEPENDENCE ASSUMPTION

When independence cannot be assumed, then we propose a proxy for Predictive Parity that the Conformal Fairness Framework can use. We propose a *Predictive Parity Proxy* as the balancing the quantity $\Pr(Y = y | y \in \mathcal{C}_\lambda(X), X \in g_a) - \Pr(Y = y | X \in g_a)$ across all groups and positive labels. Thus, for the framework, given a user-specified value for c , we want $\forall g_a, g_b \in \mathcal{G}, \forall y \in \mathcal{Y}^+$:

$$\begin{aligned} & |(\Pr[Y = y | y \in \mathcal{C}_\lambda(X), X \in g_a] - \Pr[Y = y | X \in g_a]) \\ & - (\Pr[Y = y | y \in \mathcal{C}_\lambda(X), X \in g_b] - \Pr[Y = y | X \in g_b])| < c. \end{aligned} \quad (8)$$

Observe that using $\lambda = \sup \Lambda$ will make Equation 8 equal to zero, thus c can be arbitrarily small. Intuitively, this proxy balances the information provided about an outcome if the label is in the prediction set, similar to Predictive Parity. Formally, if the label distribution is independent of group membership, then balancing the proxy would be the same as balancing Predictive Parity.

D ADDITIONAL EXPERIMENT DETAILS

D.1 DATASETS

Credit (G): The Credit dataset is from the UCI repository, and traditionally the binary target is to predict the existence of default payments (Yeh & Lien, 2009). We used a graph version of Credit as considered by Agarwal et al. (2021). To convert the dataset to a multi-class dataset, we used the education level (4 labels) as the target and used gender as the sensitive attribute, as done by Liu et al. (2023).

Pokec-{n,z} (G): The Pokec dataset (Takac & Zabovsky, 2012) is a social-network graph dataset collected from Pokec, a popular social network in Slovakia. Since several rows in the dataset are missing features, two commonly used subgraphs are the Pokec-z and Pokec-n datasets. The graphs have 4 labels, corresponding to the fieldwork. They also have two sensitive attributes, gender (2 groups) and region (2 groups). Our experiments consider each attribute individually and intersectional fairness by creating an attribute with 4 (2x2) groups.

ACSIIncome (T): In the fairness space, the American Community Services (ACS) datasets from the `Folktables` library are widely used (Ding et al., 2021). For ACSIIncome, we used the standard ACSIIncome dataset in `Folktables` however, we divided the targets into 4 classes by evenly dividing the income into 4 brackets. Race is the sensitive attribute and has 9 groups.

ACSEducation (T): Similar to ACSIIncome, we used the ACS data and selected the Education Level as our target. We broke the education level into 6 groups {did not complete high school,

has a high school diploma, has a GED, started an undergrad program, completed an undergrad program, and completed graduate or professional school}. ACSEducation also uses race as a sensitive attribute.

D.2 NON-CONFORMITY SCORES

Let $\hat{\pi}$ be a trained classification model with softmaxed output.

Threshold Prediction Sets (TPS) In TPS (Sadinle et al., 2019), the score function is $s(\mathbf{x}, y) = 1 - \hat{\pi}(\mathbf{x})_y$, where $\hat{\pi}(\mathbf{x})_y$ is the class probability for the correct class. This is the simplest method, which is also shown to be optimal with respect to efficiency (Sadinle et al., 2019).

Adaptive Prediction Sets (APS) The most popular baseline when comparing CP method is APS (Romano et al., 2019). The scoring function works by sorting the softmax logits in descending order and accumulating the class probabilities until the correct class is included. For tighter prediction sets, randomization is introduced through a uniform random variable. Formally, if $\hat{\pi}(\mathbf{x})_{(1)} \geq \hat{\pi}(\mathbf{x})_{(2)} \geq \dots \geq \hat{\pi}(\mathbf{x})_{(K-1)}$, $u \sim U(0, 1)$, and r_y is the rank of the correct label, then

$$s(\mathbf{x}, y) = \left[\sum_{i=1}^{r_y} \hat{\pi}(\mathbf{x})_{(i)} \right] - u \hat{\pi}(\mathbf{x})_y.$$

Regularized Adaptive Prediction Sets (RAPS) One drawback of APS is that it can produce large prediction sets. Angelopoulos et al. (2022) introduces a regularization approach for APS. Given the same setup and notation as APS, define $o(\mathbf{x}, y) = |\{c \in \mathcal{Y} : \hat{\pi}(\mathbf{x})_y \geq \hat{\pi}(\mathbf{x})_c\}|$. Then,

$$s(\mathbf{x}, y) = \left[\sum_{i=1}^{r_y} \hat{\pi}(\mathbf{x})_{(i)} \right] - u \hat{\pi}(\mathbf{x})_y + \nu \cdot \max\{o(\mathbf{x}, y) - k_{reg}, 0\},$$

where ν and $k_{reg} \geq 0$ are regularization hyperparameters.

Diffusion Adaptive Prediction Sets (DAPS) Graphs are rich with neighborhood information, with nodes tending to be homophilous. Intuitively, this suggests that the non-conformity scores of connected nodes are also related. To utilize this observation, DAPS H. Zargarbashi et al. (2023) performs a one-step diffusion update on the non-conformity scores. Formally, if $s(\mathbf{x}, y)$ is a point-wise score function (e.g., APS), then the diffusion step gives a new score function

$$\hat{s}(\mathbf{x}, y) = (1 - \delta)s(\mathbf{x}, y) + \frac{\delta}{|\mathcal{N}_{\mathbf{x}}|} \sum_{\mathbf{u} \in \mathcal{N}_{\mathbf{x}}} s(\mathbf{u}, y),$$

where $\delta \in [0, 1]$ is a diffusion hyperparameter and $\mathcal{N}_{\mathbf{x}}$ is the 1-hop neighborhood of \mathbf{x} .

Conformalized GNN (CFGNN) CFGNN (Huang et al., 2024) is a GNN approach to graph CP. The underlying observation is that the inefficiencies are correlated between nodes with similar neighborhood topologies. Bearing this in mind, using the calibration set, a second GNN is trained to correct the scores from the base model to optimize for efficiency through an inefficiency loss function Huang et al. (2024) propose. The inefficiency loss function definition includes a point-wise score function and can be different for training and validation. For our experiments, we set the score function to be APS and kept it consistent between training and validation.

D.3 HYPERPARAMETER TUNING

Hyperparameter tuning was done using Ray Tune (Liaw et al., 2018). For the Pokec_n and Pokec_z datasets, hyperparameters for the base GNN models were tuned via random search using Table D1 for each model type (i.e., GCN, GAT, and GraphSAGE) and for each choice of the sensitive attribute(s). For the Credit, ACS Income, and ACS Education datasets, the base XGBoost models were tuned via random search using Table D2 for each choice of sensitive attribute(s).

For Credit, Pocec_n, and Pocec_z, we tune the hyperparameters for the CFGNN model via random search using Table D3 for each model type (e.g., GCN, GAT, and GraphSAGE), for each dataset, and choice of sensitive attribute(s). We set $\mathcal{D}_{\text{calib}} = \mathcal{D}_{\text{test}} = (1 - \mathcal{D}_{\text{train}} - \mathcal{D}_{\text{valid}})/2$.

All experiments were run on a single P100 GPU.

In the interest of reproducibility, the source code for the CF Framework is provided in the supplementary material.

Table D1: Hyperparameter search space for the base GNN model for Pocec_n and Pocec_z. The last two rows are layer-type specific for GAT and GraphSAGE, respectively.

Hyperparameter	Search Space
batch_size	64
lr	loguniform(10^{-4} , 10^{-1})
hidden_channels	{16, 32, 64, 128}
layers	{1, 2, 4}
dropout	uniform(0.1, 0.8)
heads	{2, 4, 8}
aggr_fn	{mean, gcn, pool, lstm}

Table D2: Hyperparameter search space for the base XGBoost model for Credit, ACS Education, and ACS Income.

Hyperparameter	Search Space
lr	loguniform(10^{-4} , 10^{-1})
n_estimators	{2, ..., 500}
max_depth	{2, ..., 30}
gamma	uniform(0, 1)
colsample_bytree	uniform(0.25, 1.0)
colsample_bylevel	uniform(0.25, 1.0)
colsample_bynode	uniform(0.25, 1.0)
subsample	uniform(0.5, 1.0)

Table D3: Hyperparameter search space for the CFGNN model for Credit, Pocec_n, and Pocec_z. The last two rows are layer-type specific for GAT and GraphSAGE, respectively.

Hyperparameter	Search Space
batch_size	64
lr	loguniform(10^{-4} , 10^{-1})
hidden_channels	{16, 32, 64, 128}
layers	{1, 2, 3, 4}
dropout	uniform(0.1, 0.8)
τ	loguniform(10^{-3} , 10^1)
heads	{2, 4}
aggr_fn	{mean, gcn, pool, lstm}

E ADDITIONAL RESULTS

Here we provide additional results and discourse for each dataset and experiment we discuss in the main paper.

E.1 ACSEEDUCATION

Below we include results for the ACSEducation dataset, which is unique as it has a greater number of classes and is a custom dataset generated from the ACS datasets.

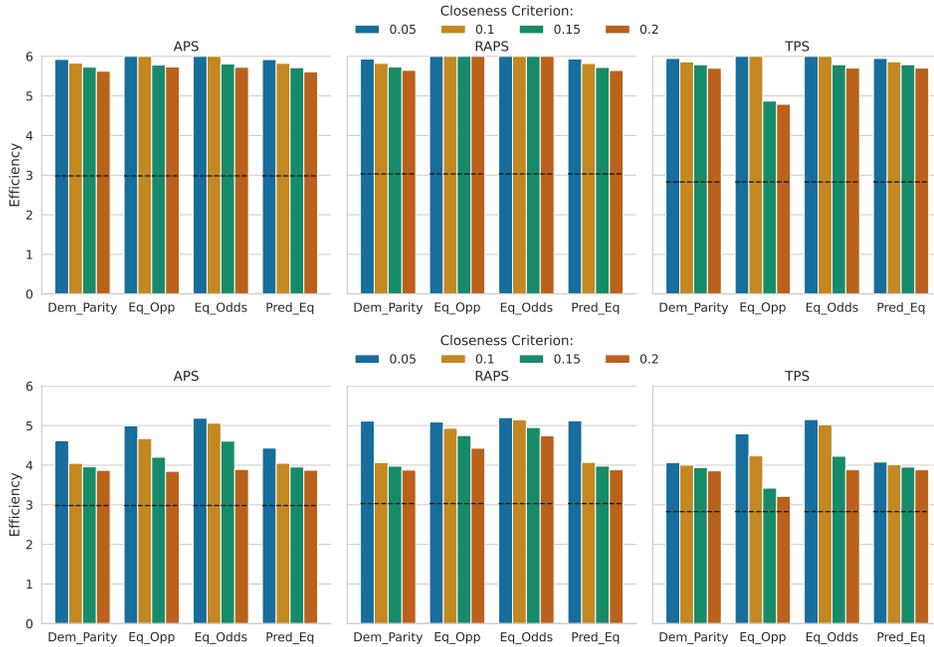


Figure E1: **ACSEducation**. Comparison of efficiencies when using the CF Framework without (top) and with (bottom) classwise lambdas. We observe that the efficiencies are better in the *right* plot. This is because $\forall_{i, \dots, k} \lambda_{\text{non-classwise}} \geq \lambda_{\text{classwise}}^i$ (k is the number of classes), which causes fewer labels to be included in the prediction set, thus improving efficiency with the classwise approach. For some experiments, the fairness disparity is 0 (e.g. APS and RAPS in the no-classwise setting), because the framework is producing the full prediction set—the trivial case—which means the coverage of $\tilde{y} \in \mathcal{Y}^+$ is 1.00, thus causing the disparity to be 0.

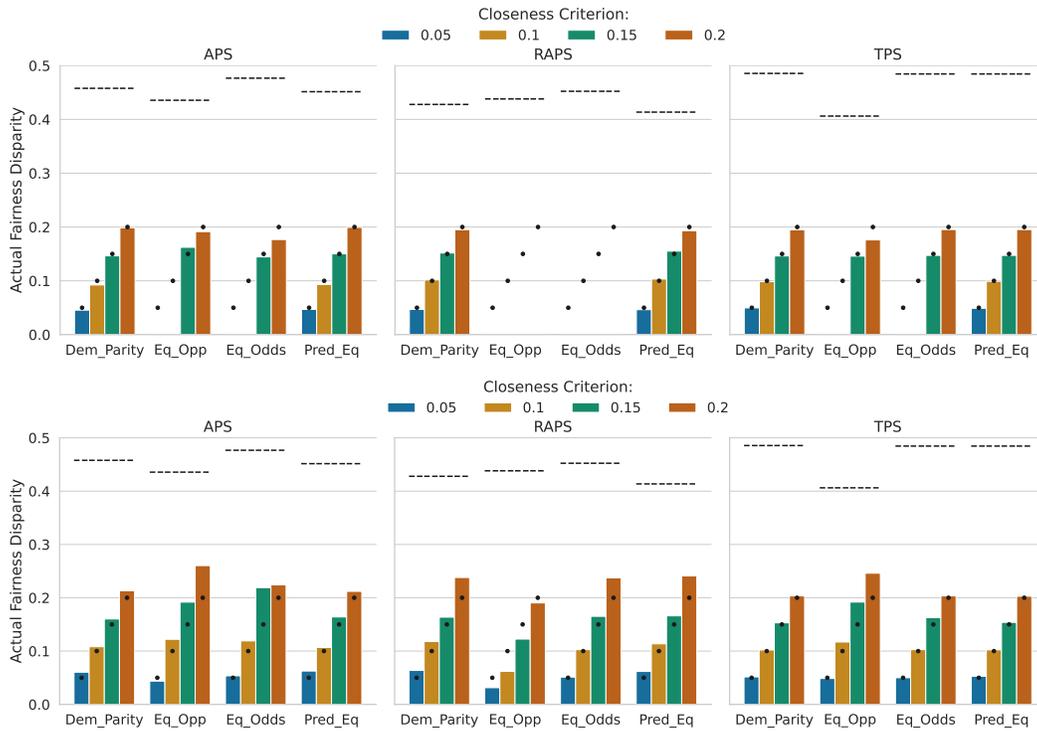


Figure E2: **ACSEducation**. Comparison of fairness disparities when using the CF Framework without (top) and with (bottom) classwise lambdas. We observe that the fairness disparities are better in the *left* plot. This is because by using a single λ , only the hardest-to-satisfy label will be at or around the coverage gap, c , unlike classwise which ensures all labels will be at or around the coverage gap, c . Since fewer labels have coverages around the coverage gap, for non-classwise in (left) the likelihood of being above the threshold is limited - as opposed to the classwise approach (right).

E.2 IMPACT OF CLASSWISE LAMBDA

E.2.1 ACSINCOME

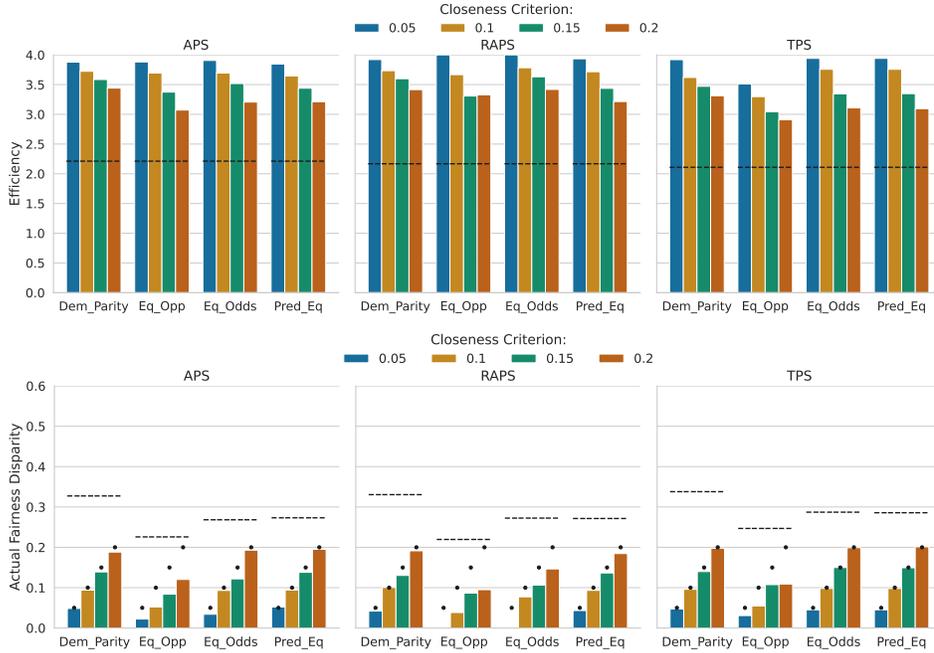


Figure E3: **ACSINCOME**. The efficiency (top) and fairness disparity (bottom) plots when **not** using classwise lambdas. We observe that the efficiency is worse, but the disparity control is much better than using classwise lambdas (see Figure 1).

E.2.2 CREDIT

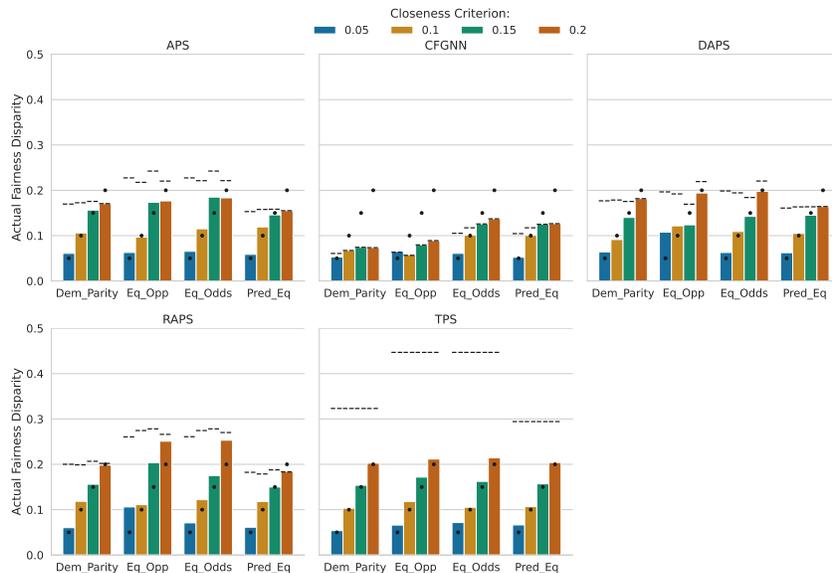


Figure E4: **Credit**. A granular version of Figure 2, which has a bar over each value of c rather than considering an average. This plot clarifies that the CF framework matches or exceeds the baseline conformal predictor's performance in fairness disparity.

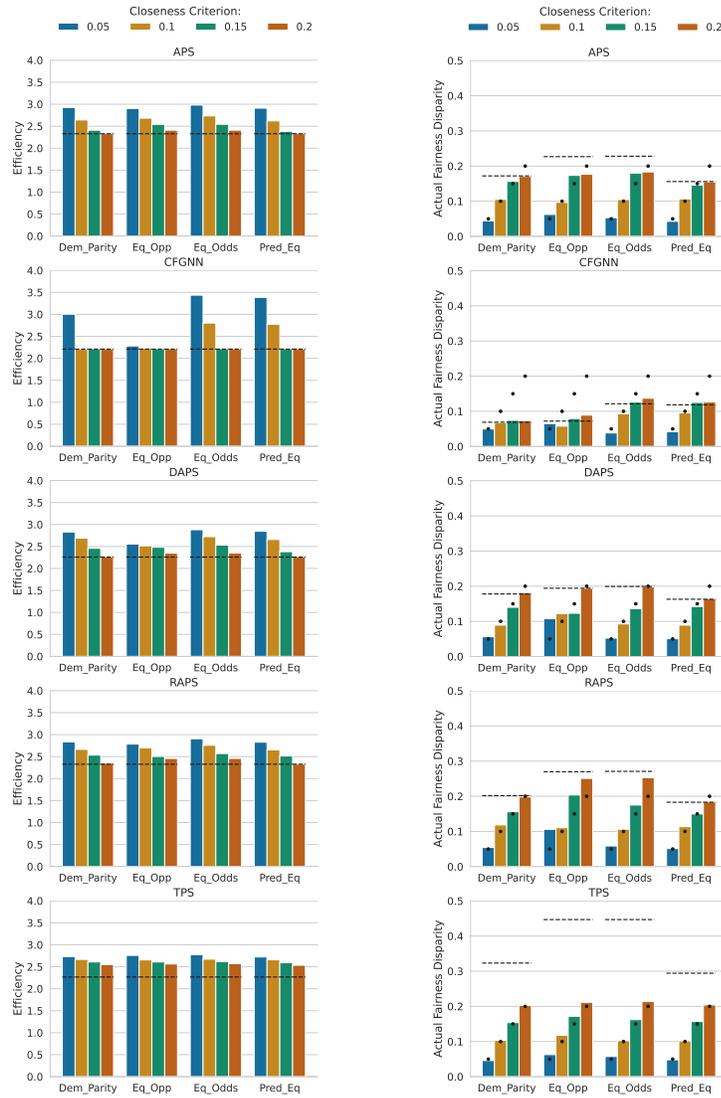


Figure E5: **Credit**. The efficiency (left) and fairness disparity (right) plots when **not** using classwise lambdas. We observe that the efficiency is worse, but the disparity control is much better than using classwise lambdas (see Figure 2).

E.3 IMPACT OF DIFFERENT POKEC SENSITIVE ATTRIBUTES

E.3.1 POKEC-N

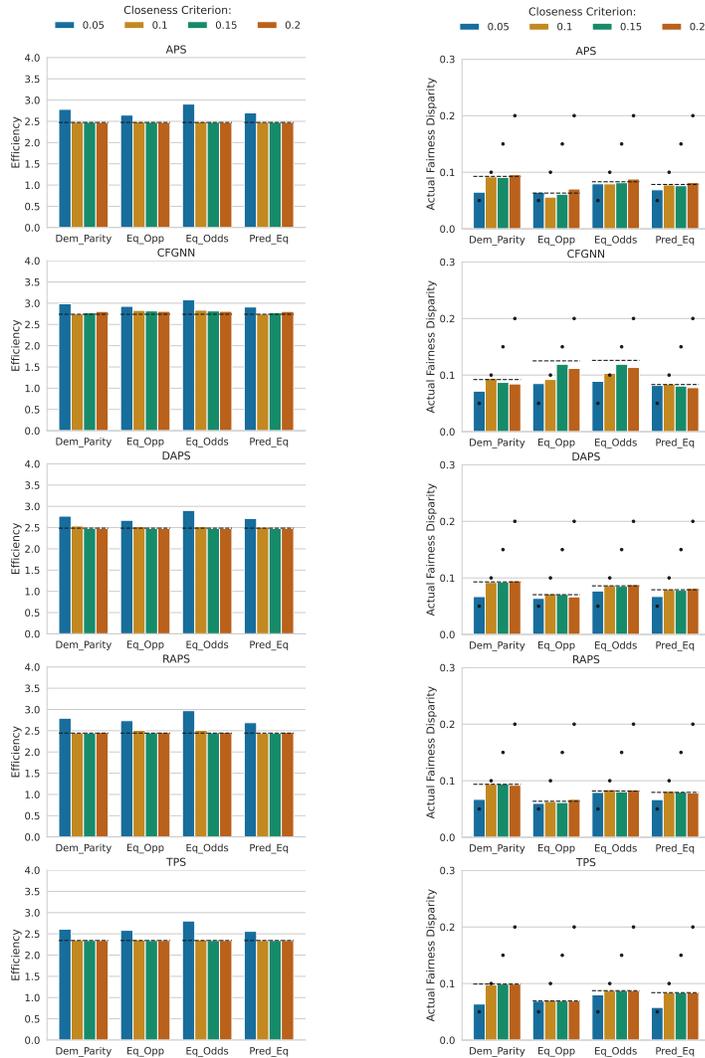


Figure E6: **Pokec-n**. The efficiency (top) and fairness violation (bottom) plots when considering only *gender* as the sensitive attribute. Observe that the baseline disparity here is smaller than the baseline in intersectional fairness (see Figure 3). Thus, when controlling for $c = 0.05$ coverage, there is a minimal change in efficiency across all the non-conformity scores.

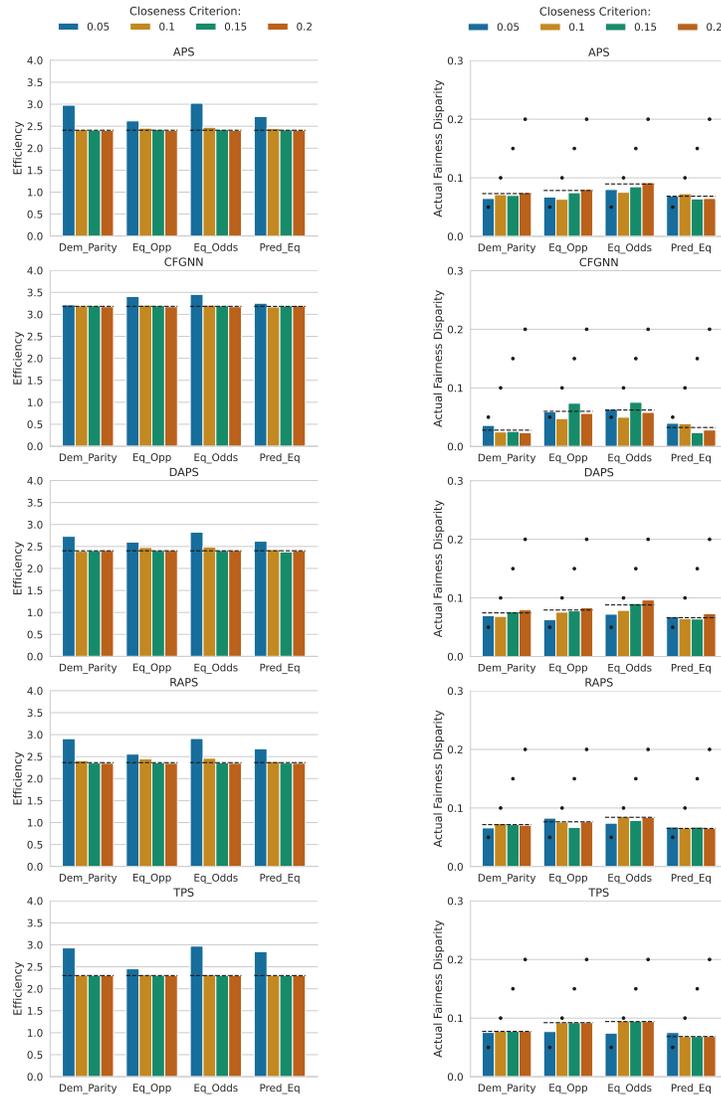


Figure E7: **Pokcc-n**. The efficiency (top) and fairness violation (bottom) plots when considering only *region* as the sensitive attribute. Compared to using *gender* in Figure E6, the prediction set sizes are larger for *region* when controlling for $c = 0.05$, thus illustrating that exact performance will vary based on the sensitive attribute and group distribution.

E.3.2 POKEC-Z

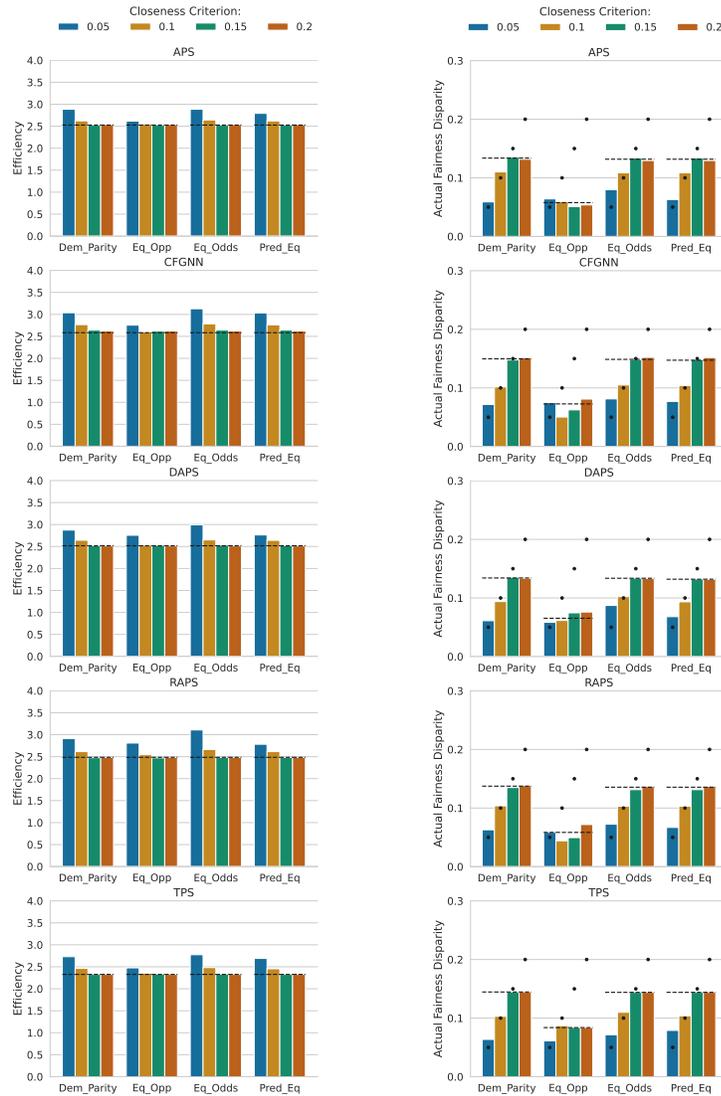


Figure E8: **Pokec-z**. The efficiency (top) and fairness violation (bottom) plots when considering only *gender* as the sensitive attribute. Unlike *Pokec-n* - using *gender* (see Figure E6 - we find that the baseline is unfair for several values of c - exemplifying the auditing capabilities of the CF framework. This result also demonstrates how fairness can vary at different localities since *Pokec-n* and *Pokec-z* are disjoint subgraphs of the same graph.

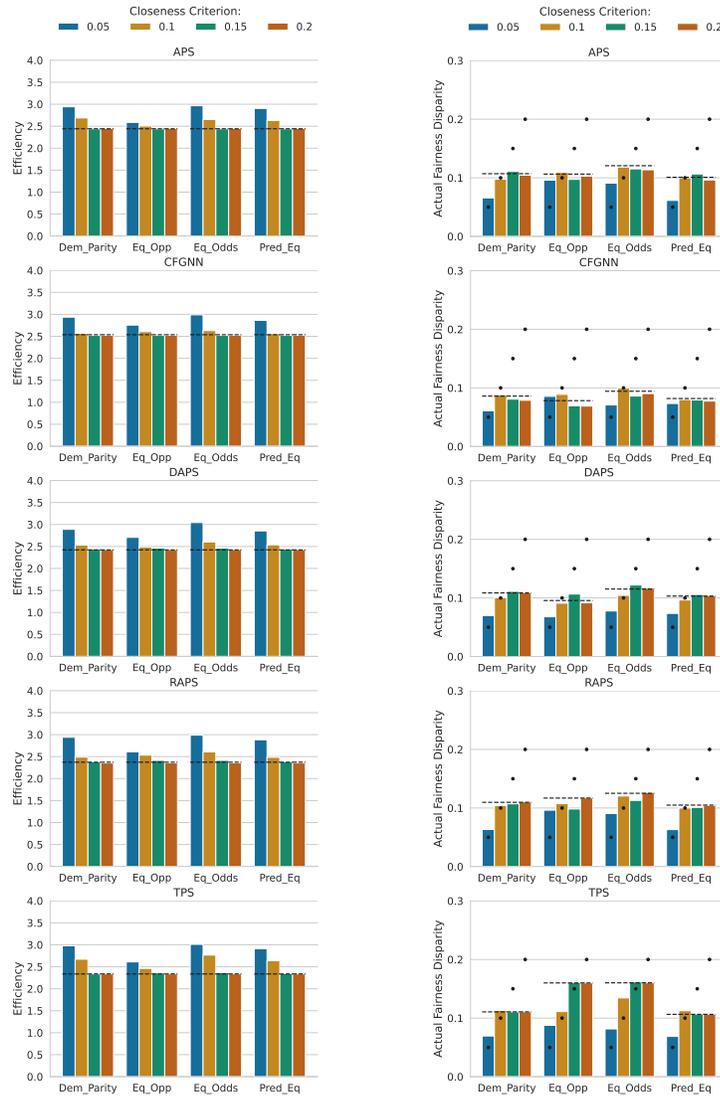


Figure E9: **Pokce-z**. The efficiency (top) and fairness violation (bottom) plots when considering only *region* as the sensitive attribute. In these plots for $c \geq 0.1$, the baselines are fair for all score functions except for TPS. The baseline for TPS being 'unfair' at $c = 0.1$ suggests TPS sacrifices fairness for efficiency.

E.4 PREDICTIVE PARITY PROXY RESULTS

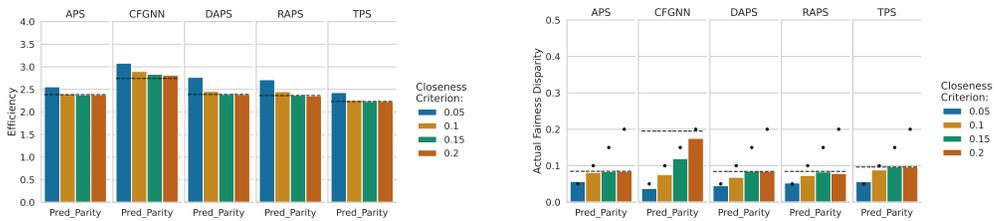


Figure E11: **Pokce-z**. Results using the Predictive Parity Proxy on a *graph* dataset. We can see that the CF Framework controls for this metric as well as or better than the base conformal predictor at all values of c and non-conformity scores.

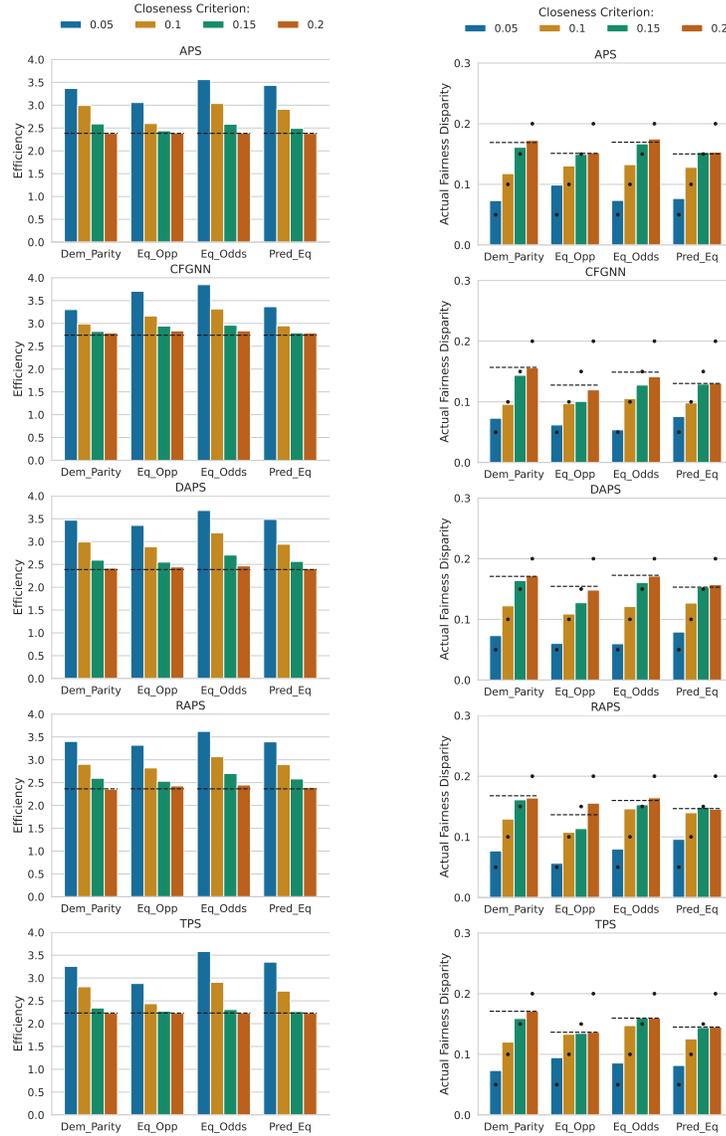


Figure E10: **Pokec-z**. Another example of intersectional fairness using both sensitive attributes of the Pokec-z dataset. Like Pokec-n, we find that the worst-case fairness disparity is better than the baseline for the CF Framework. We reiterate that the challenge of intersection fairness is the multiplicative increase in the number of groups that must be balanced. Thus, intersection fairness would benefit from more calibration points.

E.5 DISPARATE IMPACT RESULTS

Table E1: **ACSEducation**. The CF framework significantly enhances fairness (increases the ratio close to the desired 80% disparate impact rule), with a cost to efficiency. The classwise approach improves efficiency further while retaining a ratio near 80%.

ACSEducation		APS		RAPS		TPS	
Classwise		Base	CF	Base	CF	Base	CF
False	Disp. Impact	0.011	0.799	0.014	0.791	0.019	0.804
	Efficiency	2.982	5.662	3.031	5.658	2.828	5.705
True	Disp. Impact	0.011	0.781	0.014	0.761	0.019	0.797
	Efficiency	2.982	4.727	3.031	4.724	2.828	4.744

Table E2: **ACSIncome**. The CF framework significantly enhances fairness (increases the ratio close to the desired 80% disparate impact rule), with a cost to efficiency. The classwise approach improves efficiency further while retaining a ratio near 80%.

ACSIncome		APS		RAPS		TPS	
Classwise		Base	CF	Base	CF	Base	CF
False	Disp. Impact	0.397	0.815	0.387	0.847	0.356	0.804
	Efficiency	2.212	3.557	2.169	3.610	2.109	3.416
True	Disp. Impact	0.397	0.797	0.387	0.790	0.356	0.797
	Efficiency	2.212	2.674	2.169	2.752	2.109	2.679

Table E3: **Credit**. For the Credit dataset, we notice that the CF Framework improves over the baseline for APS, RAPS, TPS, and DAPS. We find that the baseline for TPS performs the worst of the 4 methods regarding disparate impact. Interestingly, the CFGNN baseline, on the other hand, maximizes the disparate impact while having the best efficiency. CFGNN demonstrates a case where CF does not perform worse than the baseline since the baseline was already ‘fair’. It also provides an example of where an audit via CF would find that the CFGNN conformal predictor is *a priori* fair.

Credit		APS		RAPS		TPS		CFGNN		DAPS	
Classwise		Base	CF	Base	CF	Base	CF	Base	CF	Base	CF
False	Disp. Impact	0.646	0.821	0.586	0.768	0.252	0.793	0.922	0.922	0.539	0.809
	Efficiency	2.325	2.561	2.326	2.559	2.268	2.620	2.202	2.202	2.254	2.573
True	Disp. Impact	0.646	0.821	0.586	0.768	0.252	0.793	0.922	0.922	0.539	0.809
	Efficiency	2.325	2.513	2.326	2.509	2.268	2.558	2.202	2.202	2.254	2.526

Table E4: **Pokec-n**. For Pokec-n, the CF Framework improves the disparate impact of the baseline conformal predictor. Similar to the other fairness metrics, the disparate impact worsens when considering intersectional fairness and is near or exceeds 80% when only one sensitive attribute is considered (i.e., *region* or *gender*).

		Pokec-n	APS		RAPS		TPS		CFGNN		DAPS	
Classwise	Group		Base	CF	Base	CF	Base	CF	Base	CF	Base	CF
False	Gender	Disp. Impact	0.798	0.802	0.793	0.811	0.777	0.796	0.784	0.797	0.800	0.806
		Efficiency	2.465	2.565	2.434	2.648	2.343	2.494	2.636	2.838	2.474	2.639
True	Gender	Disp. Impact	0.798	0.802	0.793	0.810	0.777	0.796	0.784	0.797	0.800	0.806
		Efficiency	2.465	2.473	2.434	2.457	2.343	2.358	2.636	2.666	2.474	2.494
False	Region	Disp. Impact	0.814	0.823	0.820	0.829	0.803	0.814	0.916	0.918	0.822	0.831
		Efficiency	2.401	2.498	2.373	2.552	2.300	2.426	3.168	3.185	2.413	2.584
True	Region	Disp. Impact	0.814	0.814	0.820	0.822	0.803	0.806	0.916	0.917	0.822	0.824
		Efficiency	2.401	2.430	2.374	2.404	2.300	2.332	3.168	3.180	2.413	2.437
False	Region & Gender	Disp. Impact	0.718	0.792	0.723	0.774	0.716	0.785	0.602	0.812	0.720	0.787
		Efficiency	2.485	3.037	2.435	3.012	2.341	2.970	2.537	3.563	2.509	2.991
True	Region & Gender	Disp. Impact	0.718	0.767	0.723	0.760	0.716	0.754	0.602	0.779	0.720	0.764
		Efficiency	2.485	2.739	2.435	2.656	2.341	2.566	2.537	2.964	2.509	2.709

Table E5: **Pokec-z**. For Pokec-z, the CF Framework improves the disparate impact of the baseline conformal predictor. Similar to the other fairness metrics, the disparate impact worsens when considering intersectional fairness and is near or exceeds 80% when only one sensitive attribute is considered (i.e., *region* or *gender*).

		Pokec-z	APS		RAPS		TPS		CFGNN		DAPS	
Classwise	Group		Base	CF								
False	Gender	Disp. Impact	0.798	0.802	0.737	0.800	0.737	0.800	0.784	0.797	0.800	0.806
		Efficiency	2.523	2.995	2.489	3.088	2.327	3.134	2.512	3.388	2.522	2.996
True	Gender	Disp. Impact	0.798	0.802	0.737	0.800	0.661	0.742	0.784	0.797	0.800	0.806
		Efficiency	2.523	2.572	2.489	2.571	2.327	2.415	2.512	2.642	2.522	2.579
False	Region	Disp. Impact	0.807	0.823	0.796	0.826	0.812	0.816	0.781	0.796	0.811	0.829
		Efficiency	2.429	2.569	2.369	2.592	2.337	2.497	2.546	2.644	2.416	2.586
True	Region	Disp. Impact	0.807	0.815	0.796	0.800	0.812	0.816	0.781	0.796	0.811	0.813
		Efficiency	2.429	2.475	2.369	2.403	2.337	2.391	2.546	2.603	2.416	2.453
False	Region & Gender	Disp. Impact	0.658	0.787	0.661	0.770	0.640	0.785	0.590	0.806	0.658	0.768
		Efficiency	2.408	3.173	2.359	3.214	2.265	3.175	2.512	3.471	2.409	3.133
True	Region & Gender	Disp. Impact	0.658	0.773	0.661	0.742	0.640	0.745	0.590	0.800	0.658	0.749
		Efficiency	2.408	2.799	2.359	2.741	2.265	2.627	2.512	2.788	2.409	2.743

E.6 COMPARISON WITH BATCHGCP

To produce conformal prediction sets with group-wise coverage, BatchGCP (Jung et al., 2023) learns a group-dependent threshold function to provide $1 - \alpha$ coverage for the correct label - i.e. $\Pr(y_{n+1} \in \mathcal{C}_{\hat{f}(\mathbf{x}_{n+1}; \lambda)}(\mathbf{x}_{n+1}) \mid \mathbf{x}_{n+1} \in g') = 1 - \alpha \forall g' \in \mathcal{G}$. To achieve this, a group-dependent threshold function, $\hat{f}(\mathbf{x}_{n+1}; \lambda)$, in Equation 9 is used to construct prediction sets by adding a correction to the base threshold function $f(x)$. For the scoring functions considered (i.e., APS), we define $f(x) \equiv \hat{q}(\alpha)$ as the $1 - \alpha$ quantile of the calibration scores. $\lambda \in \mathbb{R}^{|\mathcal{G}|}$ is a vector with $\lambda_{g'}$ corresponding to the entry for the g' group. The groups, g' , may intersect, allowing for \mathbf{x}_{n+1} to be a part of multiple groups.

$$\hat{f}(x; \lambda) = f(x) + \sum_{g' \in \mathcal{G}} \lambda_{g'} \mathbf{1}[x \in g'] \quad (9)$$

To determine the value of λ , Jung et al. (2023) solve the convex optimization problem in 10, where L_q is the pinball loss (Equation 11) - a function used to determine how close a threshold, $\hat{f}(x; \lambda)$, is to a specific quantile, $1 - \alpha$, for a given score, s . Using the pinball loss, λ^* - the optimal λ - is computed and used to construct prediction sets.

$$\lambda^* = \underset{\lambda}{\operatorname{argmin}} \mathbb{E}[L_q(\hat{f}(x; \lambda), s(x, y))] \quad (10)$$

$$L_{1-\alpha}(\tau, s) = (s - \tau)(1 - \alpha)\mathbf{1}[s > \tau] + (\tau - s)\alpha\mathbf{1}[s \leq \tau] \quad (11)$$

To compare against BatchGCP, we adapted the codebase³ provided by Jung et al. (2023) to accommodate APS as a non-conformity score and use the classification variant of the Folktables datasets since Jung et al. (2023) conducts experiments with those datasets. With the BatchGCP implementation, we get the threshold function, \hat{f} , and use it when constructing prediction sets. The fairness disparity is evaluated using Demographic Parity and Disparate Impact. Since BatchGCP aims to optimize *group*-conditional coverage, we only compare against those two fairness metrics, for a fair comparison.

While BatchGCP provides PAC guarantees on group-wise coverage, it does not *necessarily* provide the fairness guarantees our framework does, as empirically seen with the ACSIncome and ACSEducation datasets in Table E6 and Figures E12 and E13. We can dissect the poor performance on these datasets by looking at the per-group conditional coverages in Figure E14, where several groups are undercovered. We note that BatchGCP *requires group information at inference time*, which restricts it from settings where group information may be unavailable at inference time - this is not a limitation of the CF Framework.

Table E6: Comparing Base APS, BatchGCP, and CF framework under Disparate Impact. For the CF framework, we set $c = 0.8$ and use classwise lambdas. We observe that the CF Framework can achieve a disparate impact value much closer to $c = 0.8$ with little effect on efficiency.

		Base APS	BatchGCP	CF
ACSEducation	Disp. Impact	0.011	0.576	0.781
	Efficiency	2.982	2.893	4.727
ACSIncome	Disp. Impact	0.397	0.349	0.797
	Efficiency	2.212	2.200	2.674

³<https://github.com/ProgBelarus/BatchMultivaldConformal>

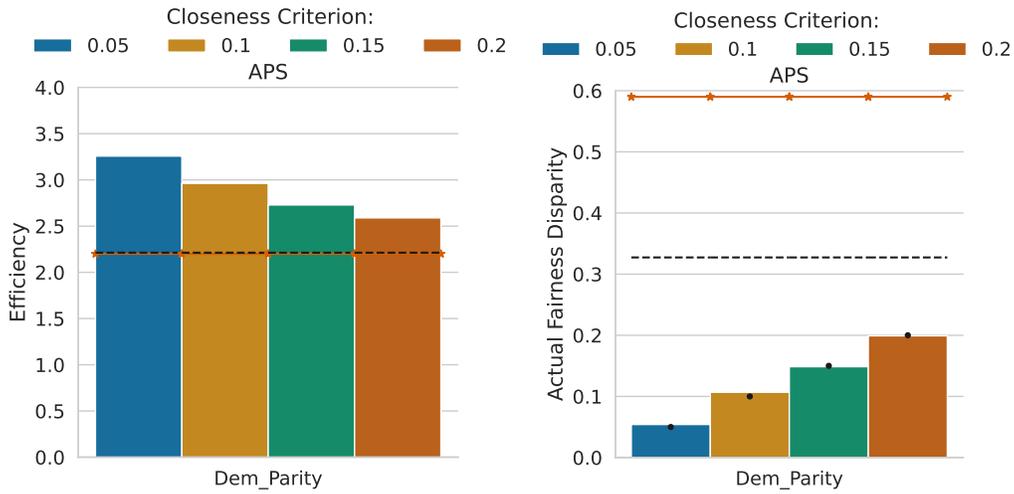


Figure E12: **ACSIncome**. Comparing Base APS, BatchGCP, and CF framework under Demographic Parity. For the CF framework, we vary $c \in \{0.05, 0.1, 0.15, 0.2\}$. We observe that the CF Framework achieves the smallest disparity for every value of c (seen on the right figure) with a small cost to the efficiency (as seen on the left figure). In the figure, Base APS = Black lines, BatchGCP = Red lines, and the CF framework is the bar charts. The black dots are the *desired* disparity level for the CF framework, which the CF framework achieves, while neither Base APS nor BatchGCP meets these thresholds.

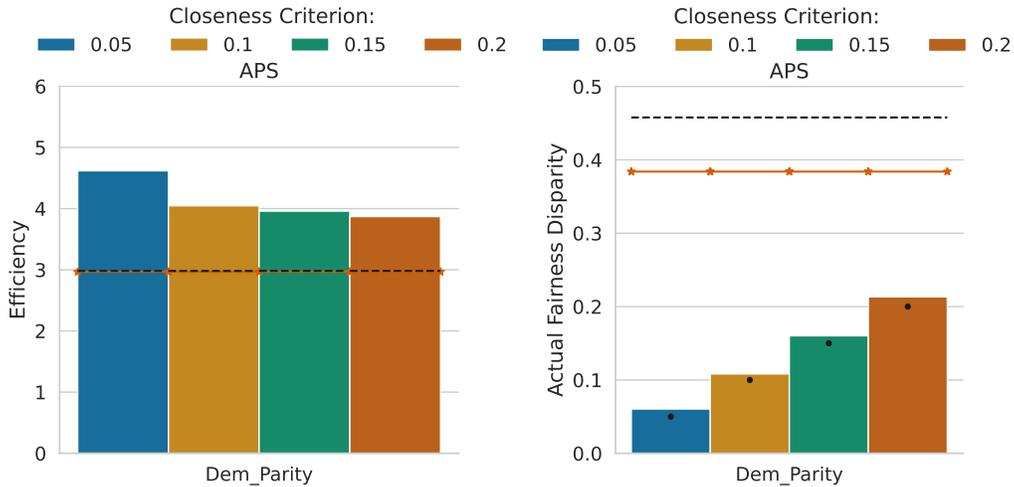


Figure E13: **ACSEducation**. Comparing Base APS, BatchGCP, and CF framework under Demographic Parity. For the CF framework, we vary $c \in \{0.05, 0.1, 0.15, 0.2\}$. We observe that the CF Framework achieves the smallest disparity for every value of c (seen on the right figure) with a small cost to the efficiency (as seen on the left figure). In the figure, Base APS = Black lines, BatchGCP = Red lines, and the CF framework is the bar charts. The black dots are the *desired* disparity level for the CF framework, which the CF framework achieves, while neither Base APS nor BatchGCP meets these thresholds.

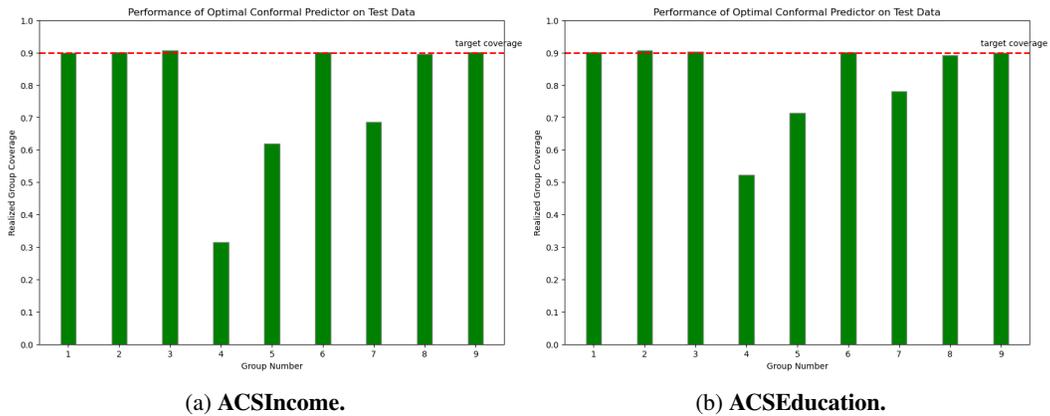


Figure E14: To dissect the poor performance of BatchGCP seen in Figures E12 and E13, we present the per-group conditional coverages for both datasets and see that certain groups are significantly undercovered (i.e., groups 4, 5, and 7 in both figures).