

# Fairness and Explainability: Bridging the Gap towards Fair Model Explanations

Yuying Zhao, Yu Wang, Tyler Derr

Vanderbilt University  
{yuying.zhao, yu.wang.1, tyler.derr}@vanderbilt.edu

## Abstract

While machine learning models have achieved unprecedented success in real-world applications, they might make biased/unfair decisions for specific demographic groups and hence result in discriminative outcomes. Although research efforts have been devoted to measuring and mitigating bias, they mainly study bias from the result-oriented perspective while neglecting the bias encoded in the decision-making procedure. This results in their inability to capture procedure-oriented bias, which therefore limits the ability to have a fully debiasing method. Fortunately, with the rapid development of explainable machine learning, explanations for predictions are now available to gain insights into the procedure. In this work, we bridge the gap between fairness and explainability by presenting a novel perspective of procedure-oriented fairness based on explanations. We identify the procedure-based bias by measuring the gap of explanation quality between different groups with Ratio-based and Value-based Explanation Fairness. The new metrics further motivate us to design an optimization objective to mitigate the procedure-based bias where we observe that it will also mitigate bias from the prediction. Based on our designed optimization objective, we propose a Comprehensive Fairness Algorithm (CFA), which simultaneously fulfills multiple objectives - improving traditional fairness, satisfying explanation fairness, and maintaining the utility performance. Extensive experiments on real-world datasets demonstrate the effectiveness of our proposed CFA and highlight the importance of considering fairness from the explainability perspective. Our code: <https://github.com/YuyingZhao/FairExplanations-CFA>.

## Introduction

Recent years have witnessed the unprecedented success of applying machine learning (ML) models in real-world domains, such as improving the efficiency of information retrieval (Fu et al. 2021; Wang et al. 2022b) and providing convenience with intelligent language translation (Dong et al. 2015). However, recent studies have revealed that historical data may include patterns of previous discriminatory decisions dominated by sensitive features such as gender, age, and race (Barocas, Hardt, and Narayanan 2017; Mehrabi et al. 2021). Training ML models with data including such

historical discrimination would explicitly inherit existing societal bias and further lead to cascading unfair decision-making in real-world applications (Fay and Williams 1993; Hanna and Linden 2012; Obermeyer et al. 2019).

In order to measure bias of decisions so that further optimization can be provided for mitigation, a wide range of metrics have been proposed to quantify unfairness (Dwork et al. 2012; Hardt, Price, and Srebro 2016; Barocas, Hardt, and Narayanan 2017). Measurements existing in the literature can be generally divided into group fairness and individual fairness where group fairness (Dwork et al. 2012; Hardt, Price, and Srebro 2016) measures the similarity of model predictions among different sensitive groups (e.g., gender, race, and income) and individual fairness (Kusner et al. 2017) measures the similarity of model predictions among similar individuals. However, all these measurements are computed based on the outcome of model predictions/decisions while ignoring the procedure leading to such predictions/decisions. In this way, the potential bias hidden in the decision-making process will be ignored and hence cannot be further mitigated by result-oriented debiasing methods, leading to a restricted mitigation.

Despite the fundamental importance of measuring and mitigating the procedure-oriented bias from explainability perspective, the related studies are still in their infancy (Dimanov et al. 2020; Fu et al. 2020). To fill this crucial gap, we widen the focus from solely measuring and mitigating the result-oriented bias to further identifying and alleviating the procedure-oriented bias. However, this is challenging due to the inherent complexity of ML models which aggravates the difficulty in understanding the procedure-oriented bias. Fortunately, with the development of explainability which aims to demystify the underlying mechanism of why a model makes prediction, we thus have a tool of generating human-understandable explanations for the prediction of a model, and therefore have insights into decision-making process. Given the explanations, a direct comparison would suffer from inability of automation due to the requirement of expert knowledge to understand the relationship between features and bias. The explanation quality, which is domain-agnostic, can serve as an indicator for fairness. If the quality of explanations are different for sensitive groups, the model treats them unfairly presenting higher-quality explanations for one group than the other, which indicates the model is biased.

To capture such bias encoded in the decision-making procedure, we bridge the gap between fairness and explainability and provide a novel procedure-oriented perspective based on explanation quality. We propose two group explanation fairness metrics measuring the gap between explanation quality among sensitive groups. First, **Ratio-based Explanation Fairness**  $\Delta_{\text{REF}}$  extends from result-oriented metric and quantifies the difference between ratios of instances with high-quality explanations for sensitive groups. However, due to the simplification from the continual scores to a discrete binary quality (e.g., high/low),  $\Delta_{\text{REF}}$  will ignore certain bias, motivating **Value-based Explanation Fairness**  $\Delta_{\text{VEF}}$  based on the detailed quality values. These two metrics cooperate to provide the novel explanation fairness perspective.

Traditional fairness measurements quantify bias in the result-oriented manner and our proposed explanation fairness measures quantify bias in the procedure-oriented way. These metrics together present a comprehensive view of fairness. Attempting to improve fairness from multiple perspectives while maintaining utility performance, we further design a **Comprehensive Fairness Algorithm (CFA)** based on minimizing representation distances between instances from distinct groups. This distance-based optimization is inspired by traditional result-oriented fairness methods (Dong et al. 2022a; Agarwal, Lakkaraju, and Zitnik 2021). Different from vanilla models that only aim for higher utility performance where the instance representations solely serve for the utility task, models with fairness consideration add fair constraints on the instance representations so that the instances from different groups are close to each other in the embedding space and thus avoid bias (i.e., the learned representations are task-related and also fair for sensitive groups). Naturally, we extend this idea to seek fair explanation quality in the embedding space which is composed of two parts: (a) the embeddings based on the original input features, and (b) the embeddings based on the partially masked input features according to feature attributions generated from explanations (inspired by explanation fidelity (Robnik-Šikonja and Bohanec 2018)). The objective function for explanation fairness thus becomes minimizing the representation distance between subgroups based on the input features and the masked features. Therefore, our learned representations achieve multiple goals: encoding task-related information and being fair for sensitive groups in both the prediction and the procedure. Our main contributions are three-fold:

- **Novel fairness perspectives/metrics.** Merging fairness and explainability domains, we propose explanation-based fairness perspectives and novel metrics ( $\Delta_{\text{REF}}$ ,  $\Delta_{\text{VEF}}$ ) that identify bias as gaps in explanation quality.
- **Comprehensive fairness algorithm.** We design Comprehensive Fairness Algorithm (CFA) to optimize multiple objectives simultaneously. We observe additionally optimizing for explanation fairness is beneficial for improving traditional fairness.
- **Experimental evaluation.** We conduct extensive experiments which verify the effectiveness of CFA and demonstrate the importance of considering procedure-oriented bias with explanations.

## Related Work

### Fairness in Machine Learning

Fairness in ML has raised increasingly interest in recent years (Pessach and Shmueli 2022; Dong et al. 2022b), as bias is commonly observed in ML models (Ntoutsis et al. 2020) imposing risks to the social development and impacting individual’s lives. Therefore, researchers have devoted much effort into this field. To measure the bias, various metrics have been proposed, including the difference of statistical parity ( $\Delta_{\text{SP}}$ ) (Dwork et al. 2012), equal opportunity ( $\Delta_{\text{EO}}$ ) (Hardt, Price, and Srebro 2016), equalized odds (Hardt, Price, and Srebro 2016), counterfactual fairness (Kusner et al. 2017), etc. To mitigate bias, researchers design debiasing solutions through pre-processing (Kamiran and Calders 2012), processing (Agarwal, Lakkaraju, and Zitnik 2021; Agarwal et al. 2018; Jiang and Nachum 2020), and post-processing (Hardt, Price, and Srebro 2016). Existing fairness metrics capture the bias from the predictions while ignoring the potential bias in the procedure. The debiasing methods to improve them therefore mitigate the result-oriented bias rather than procedure-oriented bias.

### Model Explainability

Explainability techniques (Guidotti et al. 2018) seek to demystify black box ML models by providing human-understandable explanations from varying viewpoints on which parts (of the input) are essential for prediction (Arrieta et al. 2020; Ras et al. 2022; Yuan et al. 2020). Based on how the important scores are calculated, they are categorized into gradient-based methods (Selvaraju et al. 2017; Lundberg and Lee 2017), perturbation-based methods (Dabkowski and Gal 2017), surrogate methods (Huang et al. 2022; Ribeiro, Singh, and Guestrin 2016), and decomposition methods (Montavon et al. 2017). These techniques help identify most salient features for the model decision, which provide insights into the procedure and reasons behind the model decisions.

### Intersection of Fairness and Explainability

Only few work study the intersection between fairness and explainability and most focus on explaining existing fairness metrics. They decompose model disparity into feature (Lundberg 2020; Begley et al. 2020) or path-specific contributions (Chiappa 2019; Pan et al. 2021) aiming to understand the origin of bias. Recently, in (Dimanov et al. 2020) fairness was studied from explainability where they observed a limitation of using sensitive feature importance to define fairness. Each of these prior works have not explicitly considered explanation fairness. To the best of our knowledge, only one work (Fu et al. 2020) designed such novel metric based on the diversity of explanations for fairness-aware explainable recommendation. However, it is limited to the specific domain. In comparison, we seek to find a novel perspective on fairness regarding explanation and quantify explanation fairness in a more general form that can be utilized in different domains.

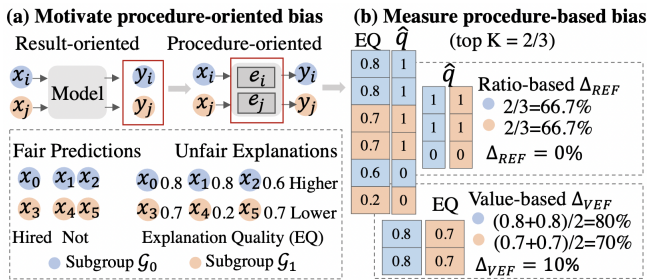


Figure 1: A motivating example of job hiring. From the result-oriented perspective, the prediction is fair since both two subgroups have the same statistical parity. From the procedure-oriented perspective, the explanation is unfair since explanation quality (EQ) of subgroup  $\mathcal{G}_0$  is generally higher than the one of  $\mathcal{G}_1$ .

## Novel Fairness Perspectives

In this section, we highlight the significance of procedure-oriented explanation fairness. Based on explanation quality, we propose Ratio-based and Value-based Explanation Fairness ( $\Delta_{REF}$  and  $\Delta_{VEF}$ ). For ease of understanding, we begin with the binary case where sensitive features  $s \in \{0, 1\}$  and the whole group  $\mathcal{G}$  is split into subgroups  $\mathcal{G}_0$  and  $\mathcal{G}_1$  based on the sensitive attribute (e.g., white and non-white subgroups divided by race) and investigate multiple-sensitive feature scenario. We include a discussion about result-oriented and procedure-oriented fairness in the end of this section.

### Motivating Explanation Quality

Although various fairness metrics have been proposed (Section ), almost all of them are directly computed based on the outputs of the predictions. Accordingly, the quantified unfairness could only reflect the result-oriented bias while ignoring the potential bias caused by the decision-making procedure. As shown in Figure 1 of a hiring example, model makes hiring decisions based on users’ features. From the result-oriented perspective, the model is fair under statistical parity constraint. However, from the procedure-oriented perspective, the model contains implicit bias during the decision-making process. The neglect of such procedure-oriented bias would restrict the debiasing methods to provide a comprehensive solution that is fair in both prediction and procedure. Therefore, the procedure leading to the decisions is critical for determining whether the decision is fair and should be explicitly considered. To get a better understanding of this process, we naturally leverage model explainability, which aims to identify the underlying mechanism of why a model makes the prediction and has raised much attention (Zhang and Zhu 2018). The demands on explainability are even more urgent in fairness domain where fairness and explainability are intertwined considering that bias identification largely hinges on the reason of prediction. Additionally, explainability will bring in extra benefit regarding model selection. When two models have comparable performance in utility and traditional fairness, explainability can provide another perspective to select the best model.

To define fairness from an explainability perspective, we first briefly introduce the explanations obtained from explainability methods, denoted as  $e_i$  for the  $i$ -th instance, which can be the important score per feature (e.g., which features are more salient for the hiring decision) or other forms. Based on the explanations, an intuitive way for fairness quantification is measuring the distribution difference of  $e_i$  and  $e_j$  where  $(i, j) \in (\mathcal{G}_0, \mathcal{G}_1)$ . However, the connection between bias and the distribution distance is controversial due to the inherent nature of diverse explanation in real-world and the lack of common understanding on fairness evaluation from attributes where expert knowledge is required and thus limit the automation across domains. Therefore, we turn the focus from explanation itself to domain-agnostic quality metric. Explanation quality (EQ) is a measurement of how good the explanation is at interpreting the model and its predictions. Various measurements (Robnik-Šikonja and Bohanec 2018) quantify EQ and provide a broad range of perspectives, which can be applied for calculating fairness. As shown in Figure 1, the unfairness in the procedure is revealed in EQ. When EQ differs among subgroups, the model treats them discriminatively, providing better-quality explanations for one group than the other. Therefore, EQ gap between subgroups indicates bias and should be mitigated. Inheriting the benefits from the variety of quality metrics, the proposed fairness metric can thus provide different perspectives regarding explainability.

### Ratio-based Explanation Fairness ( $\Delta_{REF}$ )

In order to quantify EQ gap, we borrow the idea from traditional fairness metrics  $\Delta_{SP}^1$ , which quantify bias based on the outputs  $\hat{y}_i$ . If the proportions of positive predictions are the same for subgroups, the model is fair. Regarding explanation, high quality is a positive label. Similar to the hiring example where each group deserves the same right of being hired, as an opportunity, the positive label should be obtained by subgroups fairly. Therefore, the proportions of instances with high EQ should be same for fairness consideration. We have the Ratio-based explanation fairness as

$$\Delta_{REF} = |P(\hat{q} = 1|s = 0) - P(\hat{q} = 1|s = 1)|, \quad (1)$$

where  $\hat{q}$  is the quality label for the explanation. If EQ is high,  $\hat{q} = 1$ . Otherwise,  $\hat{q} = 0$ . However, it would be challenging to define a threshold for high-quality (i.e., when EQ is larger than the threshold, the explanation is of high-quality) and the range of EQ might differ across explainers, limiting the generalizability. We transform the threshold to a top  $K$  criterion where top  $K$  percent of instances with highest EQ are regarded as of high quality (i.e.,  $\hat{q}_i$  is 1 when  $EQ_i$  belongs to the top  $K$  and 0 otherwise). The metric  $\Delta_{REF}$  measures the unfairness of high-quality proportion for subgroups and a smaller value relates to a fairer model.

### Value-based Explanation Fairness ( $\Delta_{VEF}$ )

Under certain cases, bias still exist with a low  $\Delta_{REF}$  due to the ignorance of detailed quality values in the top  $K$  (e.g.,  $\Delta_{REF} = 0$  in Figure 1 where instances in subgroups with

<sup>1</sup> $\Delta_{SP} = |P(\hat{y} = 1|s = 0) - P(\hat{y} = 1|s = 1)|$

high-quality explanations take up the same proportion but instances in  $\mathcal{G}_0$  in the top  $K$  always have larger EQ than  $\mathcal{G}_1$ . Therefore, we further compute the difference of subgroup’s average EQ as Value-based explanation fairness:

$$\Delta_{\text{VEF}} = \left| \frac{1}{|\mathcal{G}_0^K|} \sum_{i \in \mathcal{G}_0^K} \text{EQ}_i - \frac{1}{|\mathcal{G}_1^K|} \sum_{i \in \mathcal{G}_1^K} \text{EQ}_i \right|, \quad (2)$$

where  $\mathcal{G}_s^K$  denotes the top  $K$  instances with highest EQ whose sensitive feature equals  $s$ . The metric  $\Delta_{\text{VEF}}$  measures the unfairness of the average EQ of high-quality instances. Similarly, a smaller value indicates a higher level of fairness.

### Extending to Multi-class Sensitive Feature

The previous definitions are formulated for binary scenario, but they can be easily extended into multi-class sensitive feature scenarios following previous work in result-oriented fairness (Rahman et al. 2019; Spinelli et al. 2021). The main idea is to guarantee the fairness of EQ across all pairs of sensitive subgroups. Quantification strategies can be applied by leveraging either the variance of EQ across different sensitive subgroups (Rahman et al. 2019) or the maximum difference of EQ among all subgroup pairs (Spinelli et al. 2021).

### Comprehensive Fairness Algorithm

In light of the novel procedure-oriented fairness metrics defined in Section 3, we seek to develop a Comprehensive Fairness Algorithm (CFA) with multiple objectives, aiming to achieve better tradeoff performance across not only utility and traditional (result-oriented) fairness, but incorporating our proposed explanation (procedure-oriented) fairness. We note that the explanation fairness is complimentary to the existing traditional fairness metrics as compared to substitutional, and furthermore they all need to be considered under the context of utility performance (as it would be meaningless if having good fairness metrics but unusable utility performance). More specifically, we design the loss for CFA as

$$\mathcal{L}(\mathcal{G}) = \mathcal{L}_u(\mathcal{G}) + \alpha \mathcal{L}_f(\mathcal{G}_0, \mathcal{G}_1) + \beta \mathcal{L}_{exp}(\mathcal{G}_0, \mathcal{G}_1),$$

where  $\alpha, \beta$  are coefficients to balance these goals,  $\mathcal{L}_u$  is the utility loss for better utility performance,  $\mathcal{L}_f$  and  $\mathcal{L}_{exp}$  are the traditional and explanation fairness loss designed to mitigate the result/procedure-oriented bias. The framework is shown in Figure 2. In the following part, we first introduce  $\mathcal{L}_f$  and  $\mathcal{L}_{exp}$ , and provide CFA in the end where we find that mitigating procedure-oriented bias improves the result-oriented fairness and thus two loss terms are reduced to one.

### Traditional Fairness (Result-oriented)

Traditional fairness measures whether the output/prediction of the model is biased towards specific subgroups. The prediction is usually classified from the learned representations. Therefore, bias highly likely originates from the representation space. If the instance representations in subgroups in each class are close, the classifier will make fair decisions. We design an objective function based on the distance:

$$\mathcal{L}_f(\mathcal{G}_0, \mathcal{G}_1) = \mathbb{E}_{y \in Y} \mathbb{E}_{(i,j) \in (\mathcal{G}_0, \mathcal{G}_1)} \mathcal{D}(\mathbf{h}_i^y, \mathbf{h}_j^y), \quad (3)$$

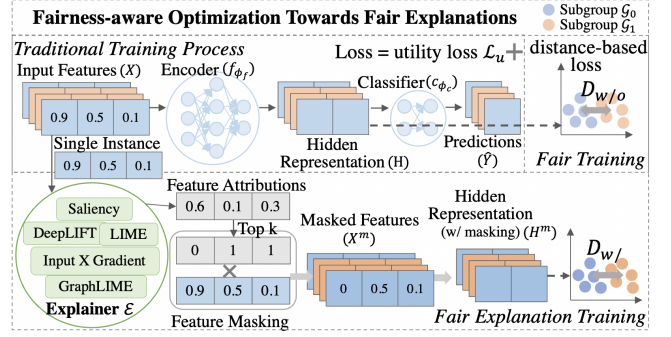


Figure 2: Overall framework of CFA.

where  $y \in Y$  is the utility label,  $\mathcal{D}(\cdot, \cdot)$  is a distance metric,  $\mathbf{h}_i^y$  and  $\mathbf{h}_j^y$  are the hidden representations of instances from subgroups whose utility label equals  $y$ . By minimizing  $\mathcal{L}_f$ , instances with same utility label are encouraged to stay close regardless of sensitive attributes. In Section 13, we conduct the ablation study to compare the effectiveness of different loss functions, including Sliced-Wasserstein distance (SW), Negative cosine similarity (Cosine), Kullback–Leibler divergence (KL) and Mean square error (MSE).

### Explanation Fairness (Procedure-oriented)

Explanation fairness measures explanation quality gap between subgroups. To reduce the gap, we first introduce the quality metric used in this paper called fidelity and then introduce the corresponding objective function which aims to minimize the distances between hidden representations of subgroups that are based on features with/without masking.

**Fidelity as Explanation Quality** Fidelity (Robnik-Šikonja and Bohanec 2018) studies the change of prediction performance (e.g., accuracy) after removing the top important features. Intuitively, if a feature is important for the prediction, when removing it, model performance will change significantly. A larger prediction change thus indicates higher importance. Fidelity has two versions - probability/accuracy-based. The former is defined as

$$\Delta_{F_i} = P_\theta(\mathbf{x}_i) - P_\theta(\mathbf{x}_i^{\mathbf{m}_i}), \quad (4)$$

where  $P_\theta$  is a general model parameterized by  $\theta$  aiming to predict the utility label  $y_i$ ,  $P_\theta(\mathbf{x}_i)$  is the predicted probability of  $y_i$ ,  $\mathbf{x}_i^{\mathbf{m}_i} = \mathbf{x}_i \odot \mathbf{m}_i$  is the masked feature where mask  $\mathbf{m}_i = [m_1, m_2, \dots, m_d] \in \{0, 1\}^d$  is generated as follows:  $m_i$  is 0 when the  $i$ -th feature has high (i.e., top  $k$ ) important score calculated from explainer  $\mathcal{E}$  otherwise 1. The latter substitutes the probability with binary score indicating whether the predicted label is true. In this paper, we mask the most important feature (i.e.,  $k = 1$ ) and adopt accuracy-based fidelity for stability consideration. The utility predictor  $P_\theta$  is composed of an encoder  $f_\phi$  and a classifier  $c_\phi$ .

**Distance-based Optimization** When substituting explanation quality to fidelity,  $\Delta_{\text{VEF}}$  (similar for  $\Delta_{\text{REF}}$ ) becomes

$$\Delta_{\text{VEF}} = \left| \frac{1}{|\mathcal{G}_0^K|} \sum_{i \in \mathcal{G}_0^K} \Delta_{F_i} - \frac{1}{|\mathcal{G}_1^K|} \sum_{i \in \mathcal{G}_1^K} \Delta_{F_i} \right|.$$

---

**Algorithm 1: Comprehensive Fairness Algorithm (CFA)**

---

**Input:** Features  $\mathbf{X}$  with utility labels  $\mathbf{Y}$ , top feature proportion  $k$ , encoder  $f_{\theta_f}$ , classifier  $c_{\theta_c}$ , explainer  $\mathcal{E}$ , coefficient  $\lambda$ , learning rate  $\epsilon$ , distance metric  $\mathcal{D}$ .

```

1 while not converged do
2   Hidden representation  $\mathbf{H} = f_{\theta_f}(\mathbf{X})$ 
3   Prediction  $\hat{\mathbf{Y}} = c_{\theta_c}(\mathbf{H})$ 
4   Mask  $\mathbf{M} = \mathcal{E}(\mathbf{X}, \hat{\mathbf{Y}}, f_{\theta_f}, c_{\theta_c}, k)$ 
5   Masked feature  $\mathbf{X}^m = \mathbf{X} \odot \mathbf{M}$ 
6   Hidden representation w/ masking  $\mathbf{H}^m = f_{\theta_f}(\mathbf{X}^m)$ 
7    $\mathcal{L}_u = -\sum_{i=1}^{|\mathbf{Y}|} (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$ 
8    $D_{w/o} = \mathcal{D}(\mathbf{H}_{\mathcal{G}_0}, \mathbf{H}_{\mathcal{G}_1})$ 
9    $D_{w/} = \mathcal{D}(\mathbf{H}^m_{\mathcal{G}_0}, \mathbf{H}^m_{\mathcal{G}_1})$ 
10   $\mathcal{L}_{exp} = D_{w/o} + D_{w/}$ 
11   $\mathcal{L} = \mathcal{L}_u + \lambda \mathcal{L}_{exp}$ 
12   $\theta = \theta - \epsilon \nabla_{\theta} \mathcal{L}$ 
13 return A fair utility predictor composed of  $f_{\theta_f}$  and  $c_{\theta_c}$ 

```

---

As shown in Eq. (4), each  $\Delta_{F_i}$  consists of two parts which are based on the original feature and the masked feature. We apply the same idea for optimizing traditional fairness in Section . If the hidden representations for instances with/without masking in subgroups are close to each other when utility label is the same,  $\Delta_{VEF}$  will be low. Therefore, to minimize the distance with and without masking, we have the objective function as

$$\mathcal{L}_{exp}(\mathcal{G}_0, \mathcal{G}_1) = \mathbb{E}_{y \in \mathcal{Y}} \mathbb{E}_{(i,j) \in (\mathcal{G}_0, \mathcal{G}_1)} (\mathcal{D}(\mathbf{h}_i^y, \mathbf{h}_j^y) + \mathcal{D}(\mathbf{h}_i^{m_y}, \mathbf{h}_j^{m_y})),$$

where  $\mathbf{h}_i^{m_y}$  and  $\mathbf{h}_j^{m_y}$  are the hidden representations based on masked features with other notations sharing the same meaning in Eq. (3). By minimizing  $\mathcal{L}_{exp}$ , instances with same utility label are encouraged to stay close both before and after masking regardless of their sensitive attributes.

### Overall Objective of CFA

The above objectives show that  $\mathcal{L}_{exp}$  contains  $\mathcal{L}_f$ , which indicates that alleviating procedure-oriented bias helps mitigating result-oriented bias under certain objective design. The final objective thus becomes  $\mathcal{L}(\mathcal{G}) = \mathcal{L}_u(\mathcal{G}) + \lambda \mathcal{L}_{exp}(\mathcal{G}_0, \mathcal{G}_1)$ , where  $\mathcal{L}_{exp}$  aims to mitigate the bias in both predictions and the procedure, and coefficient  $\lambda$  flexibly adjusts the optimization towards the desired goal. Our CFA framework is holistically presented in Algorithm 1.

## Experiments

In this section, we conduct extensive experiments to validate the proposed explanation fairness measurements and the effectiveness of the designed algorithm CFA<sup>2</sup>. In particular, we answer the following research questions:

- **RQ1:** How well can CFA mitigate the bias related to traditional result-oriented and procedure-oriented explanation fairness measurements (Section 13, Section 13)?
- **RQ2:** How well can CFA balance different categories of objectives compared with other baselines (Section 13)?

<sup>2</sup>Source code available at:

<https://github.com/YuyingZhao/FairExplanations-CFA>

- **RQ3:** Are the fairness metrics and algorithm general to other complex data domains (Section 13)?

Additionally, we further probe CFA with an ablation study of different loss functions (Section 13).

### Experimental Settings

**Dataset** We validate the proposed approach on four real-world benchmark datasets: German<sup>3</sup>, Recidivism (Jordan and Freiburger 2015), Math and Por (Cortez and Silva 2008), which are commonly adopted for fair ML (Le Quy et al. 2022).

**Baselines** Given no prior work has explicitly considered explanation fairness, we compare our methods with MLP model and two fair ML models, namely Reduction (Agarwal et al. 2018) and Reweight (Jiang and Nachum 2020). To validate the generalizability of our novel defined fairness measures on complex data, we include two fair models in the graph domain called NIFTY (Agarwal, Lakkaraju, and Zitnik 2021) and FairVGNN (Wang et al. 2022a).

**Evaluation Criteria** We evaluate model performance from four perspectives: model utility, result-oriented fairness, procedure-oriented explanation fairness, and the overall score. More specifically, the metrics are (1) *Model utility*: area under receiver operating characteristic curve (AUC), F1 score, and accuracy; (2) *Result-oriented fairness*: two widely-adopted fairness measurements  $\Delta_{SP}$  and  $\Delta_{EO}$ <sup>4</sup>; (3) *Procedure-oriented explanation fairness*: two proposed group fairness metrics  $\Delta_{REF}$  in Eq. (1) and  $\Delta_{VEF}$  in Eq. (2); (4) *Overall score*: the overall score regarding the previous three measurements via their average (per category):

$$\text{Score} = \frac{\text{AUC} + \text{F1} + \text{Acc}}{3.0} - \frac{\Delta_{SP} + \Delta_{EO}}{2.0} - \frac{\Delta_{REF} + \Delta_{VEF}}{2.0}.$$

**Setup** For a fair comparison, we record the best model hyperparameters based on the overall score in the validation set. After obtaining the optimal hyperparameters, we then run the model five times on the test dataset to get the average result. The explainer used for all methods is GraphLime (Huang et al. 2022) where for general datasets the closest instances are treated as neighbors.

### Performance Comparison

To answer **RQ1**, we evaluate CFA and general fairness methods on real-world datasets and report utility and fairness with standard deviations. The results are in Table 1 where the distance  $\mathcal{D}$  is Sliced-Wasserstein distance. Reduction and Reweight learn a logistic regression classifier using solvers to solve the optimization problem, which is different from the MLP-based implementation. They are more stable and have a smaller standard deviation. From the table, we observe:

- From the perspective of utility performance, CFA and other fairness methods have comparable performance with the vanilla MLP model, which indicates that little

<sup>3</sup>From UCI machine learning repository: <http://archive.ics.uci.edu/ml/index.php>

<sup>4</sup> $\Delta_{EO} = |P(\hat{y} = 1|y = 1, s = 0) - P(\hat{y} = 1|y = 1, s = 1)|$

	Metric	MLP	Reduction	Reweight	CFA
Recidivism	AUC $\uparrow$	86.12 $\pm$ 1.9	81.17 $\pm$ 0.0	<b>89.24 <math>\pm</math> 0.0</b>	89.02 $\pm$ 0.9
	F1 $\uparrow$	76.54 $\pm$ 2.5	76.69 $\pm$ 0.0	72.99 $\pm$ 0.0	<b>81.28 <math>\pm</math> 1.4</b>
	Acc $\uparrow$	83.48 $\pm$ 1.5	84.66 $\pm$ 0.0	83.70 $\pm$ 0.0	<b>87.17 <math>\pm</math> 0.8</b>
	$\Delta_{SP}$ $\downarrow$	6.07 $\pm$ 2.2	2.04 $\pm$ 0.0	4.27 $\pm$ 0.0	<b>1.16 <math>\pm</math> 0.5</b>
	$\Delta_{EO}$ $\downarrow$	3.19 $\pm$ 0.7	4.66 $\pm$ 0.0	3.37 $\pm$ 0.0	<b>1.14 <math>\pm</math> 0.4</b>
	$\Delta_{REF}$ $\downarrow$	4.45 $\pm$ 3.0	<b>0.53 <math>\pm</math> 0.0</b>	1.34 $\pm$ 0.9	1.98 $\pm$ 1.2
	$\Delta_{VEF}$ $\downarrow$	2.1 $\pm$ 1.4	<b>2.06 <math>\pm</math> 0.0</b>	3.22 $\pm$ 0.0	2.70 $\pm$ 0.8
Score $\uparrow$	74.15 $\pm$ 2.0	76.19 $\pm$ 0.0	75.88 $\pm$ 0.0	<b>82.33 <math>\pm</math> 0.6</b>	
German	AUC $\uparrow$	66.77 $\pm$ 2.1	63.95 $\pm$ 0.1	<b>68.03 <math>\pm</math> 0.0</b>	60.92 $\pm$ 5.2
	F1 $\uparrow$	71.11 $\pm$ 3.5	72.48 $\pm$ 1.2	74.71 $\pm$ 0.0	<b>81.14 <math>\pm</math> 2.3</b>
	Acc $\uparrow$	63.28 $\pm$ 3.2	64.40 $\pm$ 1.4	65.60 $\pm$ 0.0	<b>70.00 <math>\pm</math> 3.0</b>
	$\Delta_{SP}$ $\downarrow$	39.80 $\pm$ 9.3	24.55 $\pm$ 1.8	20.53 $\pm$ 0.0	<b>7.21 <math>\pm</math> 6.4</b>
	$\Delta_{EO}$ $\downarrow$	31.39 $\pm$ 11.5	16.03 $\pm$ 3.9	12.18 $\pm$ 0.0	<b>4.60 <math>\pm</math> 4.1</b>
	$\Delta_{REF}$ $\downarrow$	<b>4.50 <math>\pm</math> 2.8</b>	16.09 $\pm$ 8.6	7.16 $\pm$ 0.0	10.02 $\pm$ 4.2
	$\Delta_{VEF}$ $\downarrow$	<b>8.85 <math>\pm</math> 9.0</b>	23.55 $\pm$ 7.7	16.67 $\pm$ 0.0	12.87 $\pm$ 9.0
Score $\uparrow$	24.78 $\pm$ 12.3	26.83 $\pm$ 4.1	41.18 $\pm$ 0.0	<b>53.34 <math>\pm</math> 7.3</b>	
Por	AUC $\uparrow$	90.86 $\pm$ 0.4	67.64 $\pm$ 0.0	89.07 $\pm$ 0.0	<b>91.30 <math>\pm</math> 0.6</b>
	F1 $\uparrow$	58.41 $\pm$ 4.1	51.43 $\pm$ 0.0	51.43 $\pm$ 0.0	<b>60.55 <math>\pm</math> 4.7</b>
	Acc $\uparrow$	89.57 $\pm$ 0.8	89.57 $\pm$ 0.0	89.57 $\pm$ 0.0	<b>89.82 <math>\pm</math> 1.0</b>
	$\Delta_{SP}$ $\downarrow$	2.08 $\pm$ 0.8	1.93 $\pm$ 0.0	1.93 $\pm$ 0.0	<b>1.00 <math>\pm</math> 0.7</b>
	$\Delta_{EO}$ $\downarrow$	32.35 $\pm$ 7.1	<b>20.59 <math>\pm</math> 0.0</b>	<b>20.59 <math>\pm</math> 0.0</b>	27.65 $\pm$ 5.4
	$\Delta_{REF}$ $\downarrow$	8.68 $\pm$ 3.2	<b>1.37 <math>\pm</math> 0.0</b>	8.68 $\pm$ 0.0	4.66 $\pm$ 3.8
	$\Delta_{VEF}$ $\downarrow$	4.44 $\pm$ 2.2	<b>0.0 <math>\pm</math> 0.0</b>	7.69 $\pm$ 0.0	4.70 $\pm$ 3.7
Score $\uparrow$	55.83 $\pm$ 4.0	57.60 $\pm$ 0.0	57.25 $\pm$ 0.0	<b>61.55 <math>\pm</math> 3.3</b>	
Math	AUC $\uparrow$	95.73 $\pm$ 1.8	86.44 $\pm$ 0.1	95.06 $\pm$ 1.1	<b>96.97 <math>\pm</math> 0.6</b>
	F1 $\uparrow$	83.32 $\pm$ 3.3	82.65 $\pm$ 0.2	82.17 $\pm$ 2.7	<b>86.74 <math>\pm</math> 1.7</b>
	Acc $\uparrow$	88.60 $\pm$ 2.8	89.00 $\pm$ 0.0	88.00 $\pm$ 2.5	<b>91.00 <math>\pm</math> 1.3</b>
	$\Delta_{SP}$ $\downarrow$	<b>2.14 <math>\pm</math> 1.6</b>	3.59 $\pm$ 1.6	4.54 $\pm$ 1.6	4.59 $\pm$ 2.4
	$\Delta_{EO}$ $\downarrow$	16.89 $\pm$ 4.7	28.00 $\pm$ 2.7	23.78 $\pm$ 3.3	<b>6.22 <math>\pm</math> 3.2</b>
	$\Delta_{REF}$ $\downarrow$	9.60 $\pm$ 5.3	<b>2.08 <math>\pm</math> 2.6</b>	4.8 $\pm$ 0.0	6.08 $\pm$ 5.6
	$\Delta_{VEF}$ $\downarrow$	4.22 $\pm$ 5.2	8.22 $\pm$ 4.1	<b>3.56 <math>\pm</math> 4.4</b>	4.00 $\pm$ 4.9
Score $\uparrow$	72.79 $\pm$ 3.8	65.08 $\pm$ 1.5	70.08 $\pm$ 6.7	<b>81.1 <math>\pm</math> 5.6</b>	

Table 1: Model utility and fairness measurements of binary classification. The best and second best results are marked bold or underline, respectively. Additionally, the  $\uparrow$  represents the larger the better and  $\downarrow$  represents the opposite.

compensation has been made when pursuing higher fairness performance. In some dataset such as Recidivism, the utility performance of CFA is even higher indicating that the distance-based loss might provide a better-quality representation for the downstream tasks.

- From the perspective of traditional result-oriented fairness measurements (i.e.,  $\Delta_{SP}$  and  $\Delta_{EO}$ ), CFA has comparable or better performance when compared with existing fairness methods. This shows our model is able to provide fair predictions for subgroups.
- From the perspective of procedure-oriented fairness measurements (i.e.,  $\Delta_{REF}$  and  $\Delta_{VEF}$ ), we observe that CFA does not obtain the best performance in all datasets. This is partially due to the complex impact of multiple objectives and the criterion for model selection. When the difference in utility impacts more on the validation score for model selection, then models with higher utility might be selected at a cost of worse explanation fairness.
- From the global view, CFA obtains the highest overall score, which measures the multi-task performance, showing its strong ability of balancing multiple goals.

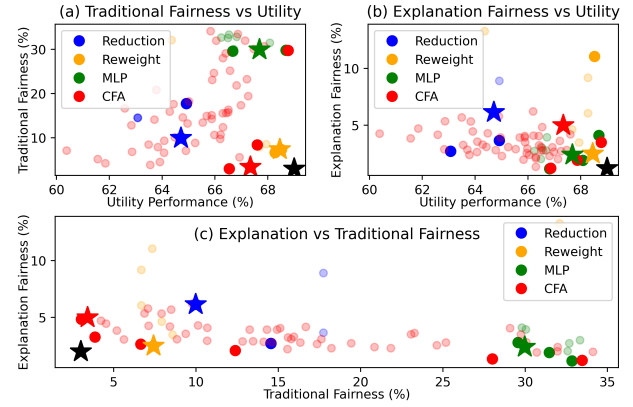


Figure 3: Performance comparison in German dataset.

### Tradeoff between Different Objectives

To answer **RQ2**, we explore the tradeoff between objectives in the validation records. Figure 3 shows the performance comparison in German dataset. Points relate to specific hyperparameters. Larger and darker circles are non-dominated solutions which form the Pareto frontier in multi-task and the others are dominated which either have poorer accuracy or fairness. Black star indicates the optimal solution direction and other stars correspond to the best hyperparameter setting based on the overall score. For Reduction, only three different results are obtained, showing its inability to provide diverse solutions. MLP model has poor performance in traditional fairness and achieve varying levels of explanation fairness. But if not selected based on the overall score, the model will gain slight improvement in utility performance at the cost of fairness, which validates the benefit of using overall score for model selection. While Reweight provides good performance in traditional fairness, most of its settings are unable to provide predictions with good explanation fairness, indicating that alleviating result-oriented bias does not necessarily mitigate procedure-oriented bias and this type of bias will be ignored by traditional fairness. This further highlights the necessity of considering explainability to provide a procedure-oriented perspective. CFA provides diverse solutions and thus satisfies different needs for utility and fairness. Although not always the one closest to the optimal solution, it is the only one method that maintains close to the optimal in all three conditions, indicating CFA has the best performance regarding balancing difference objectives which is also revealed in the overall score in Table 1.

### Effect of Coefficient $\lambda$

We explore the impact of the regularization coefficient  $\lambda$  which is used for balancing utility and fairness performance. Results based on the validation records are in Figure 4 where (a)-(b) record the average performance of utility, traditional/explanation fairness and (c)-(d) present the detailed explanation fairness  $\Delta_{REF}$  and  $\Delta_{VEF}$ . Figure 4 (b) shows that fairness performance increases while utility performance decreases along with the increase of  $\lambda$  in Por dataset. While in Figure 4 (a), fairness performance improves but utility

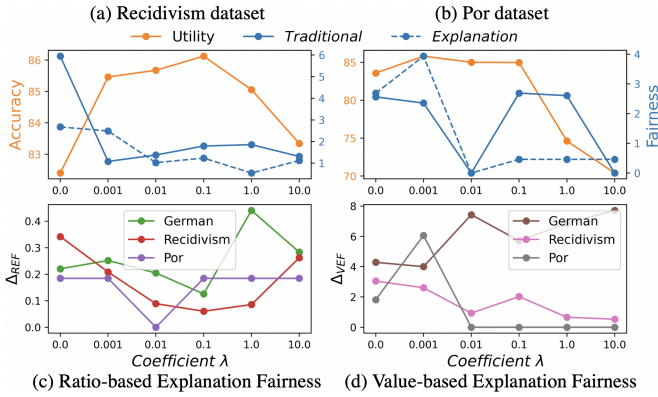


Figure 4: The effect of the coefficient  $\lambda$ .

performance also increases at first in Recidivism dataset, which indicates that in some scenario, good utility and fairness performance can exist together. From Figure 4 (c)-(d), we observe that explanation fairness decreases for Recidivism dataset. The trend is not clear for other datasets but we can find at least one point with better performance than that without fairness optimization which verifies that the fairness objective is beneficial for improving explanation fairness.

### Graph Domain Extension

To answer **RQ3**, we take graph as an example to show that our metrics are universal and can be applied to various domains. To handle the issue of feature propagation which results in sensitive information leakage, masking strategy is slightly different compared with i.i.d. data where mask is applied per instance. In graph domain, the corresponding channels in the neighborhood are also masked. We obtained results for binary classification for NIFTY (Agarwal, Lakkaraju, and Zitnik 2021) and FairVGNN (Wang et al. 2022a). Figure 5 (b) shows these methods with fair consideration have explanation bias encoded in the procedure.

Additionally, to test the effectiveness of proposed optimization, fairness loss term in Eq. (3) is added to the objective function (denoted as NIFTY<sup>+</sup> and FairVGNN<sup>+</sup>). Figure 5 shows performance in traditional/explanation fairness both improve, indicating that optimization is applicable to graph and successfully alleviates bias from both predictions and procedure. The impacts on utility performance are different - NIFTY experiences a larger decrease than FairVGNN, which is due to a lack of coefficient tuning of the newly added term, which leads to the difference. The result shows the potential of improving explanation fairness and tradeoff can be adjusted later by tuning the coefficient.

### Ablation Study

We conduct ablation study on the distance function. For ease of experiment, we only tune  $\lambda$  with other hyperparameters fixed to reduce the search space. Results in Table 2 show that all of them obtain best final score when compared with baseline results in Table 1, validating that CFA improves fairness despite of the choices of the distance measurements. Regarding different distance metrics, SW achieves better perfor-

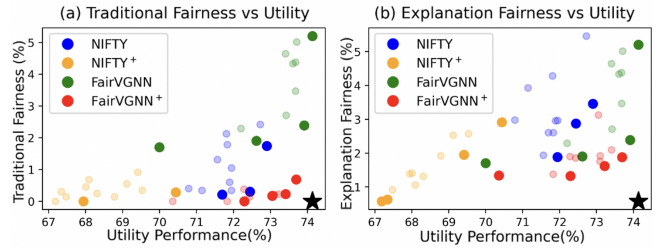


Figure 5: Performance comparison between graph-based methods with and without fair explanation optimization.

	Metric	SW	Cosine	KL	MSE
Recidivism	AUC $\uparrow$	88.35 $\pm$ 0.5	<b>91.73 <math>\pm</math> 2.8</b>	90.00 $\pm$ 2.0	89.43 $\pm$ 1.0
	F1 $\uparrow$	79.82 $\pm$ 0.6	82.03 $\pm$ 3.1	<b>82.26 <math>\pm</math> 3.9</b>	81.08 $\pm$ 1.8
	Acc $\uparrow$	86.13 $\pm$ 0.2	87.52 $\pm$ 2.4	<b>87.80 <math>\pm</math> 2.6</b>	87.55 $\pm$ 1.0
	$\Delta_{SP}$ $\downarrow$	<b>0.6 <math>\pm</math> 0.4</b>	3.54 $\pm$ 2.1	3.92 $\pm$ 2.4	1.98 $\pm$ 1.8
	$\Delta_{EO}$ $\downarrow$	<b>1.2 <math>\pm</math> 0.6</b>	1.71 $\pm$ 1.4	4.65 $\pm$ 2.7	4.02 $\pm$ 1.1
	$\Delta_{REF}$ $\downarrow$	1.50 $\pm$ 1.0	<b>1.35 <math>\pm</math> 1.2</b>	5.46 $\pm$ 3.5	2.31 $\pm$ 2.2
	$\Delta_{VEF}$ $\downarrow$	1.92 $\pm$ 2.0	<b>1.71 <math>\pm</math> 1.4</b>	4.65 $\pm$ 2.7	3.01 $\pm$ 2.5
	Score $\uparrow$	82.17 $\pm$ 0.9	<b>82.71 <math>\pm</math> 4.0</b>	78.44 $\pm$ 4.8	80.36 $\pm$ 1.6

Table 2: Performance of different distance functions.

mance in fairness aspects than KL owing to the knowledge of underlying space geometry; MSE has poor utility performance since enforcing groups to have same embeddings is not necessary; while Cosine is hypothesized to have the best performance due to ignoring scale which inherently relaxes the restriction enforced by MSE.

### Conclusion

In this paper, we investigate the potential bias during the decision-making procedure, which is ignored by traditional metrics. We provide a novel fairness perspective to raise the concern of such procedure-oriented bias. We utilize explainability to provide insights into the procedure and identify bias with two novel fairness metrics based on explanation quality. To simultaneously fulfill multiple goals - improving traditional fairness, satisfying explanation fairness, and maintaining utility performance, we design a Comprehensive Fairness Algorithm (CFA) for optimization. During the optimization, we uncover that optimizing procedure-oriented fairness is beneficial for result-oriented fairness. Experimental results demonstrate that our proposed explanation fairness captures bias ignored by previous result-oriented metrics, and the designed CFA effectively mitigates bias from multiple perspectives while maintaining good model utility. Additionally, our proposed metrics and the optimization strategy can be easily applied to other domains, showing a good generalizability. In the future, we plan to explore explanation fairness in inherently explainable models and further design fair and explainable GNNs. We also plan to extend and apply CFA towards fair model explanations in other data types (e.g., images and text) and in the setting of explanation supervision (Gao et al. 2022).

## References

- Agarwal, A.; Beygelzimer, A.; Dudík, M.; Langford, J.; and Wallach, H. 2018. A reductions approach to fair classification. In *International Conference on Machine Learning*, 60–69. PMLR.
- Agarwal, C.; Lakkaraju, H.; and Zitnik, M. 2021. Towards a unified framework for fair and stable graph representation learning. In *Uncertainty in Artificial Intelligence*, 2114–2124. PMLR.
- Arrieta, A. B.; Díaz-Rodríguez, N.; Del Ser, J.; Bennetot, A.; Tabik, S.; Barbado, A.; García, S.; Gil-López, S.; Molina, D.; Benjamins, R.; et al. 2020. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information fusion*, 58: 82–115.
- Barocas, S.; Hardt, M.; and Narayanan, A. 2017. Fairness in machine learning. *Nips tutorial*, 1: 2.
- Begley, T.; Schwedes, T.; Frye, C.; and Feige, I. 2020. Explainability for fair machine learning. *arXiv preprint arXiv:2010.07389*.
- Chiappa, S. 2019. Path-specific counterfactual fairness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 7801–7808.
- Cortez, P.; and Silva, A. 2008. Using data mining to predict secondary school student performance. *EUROSIS*.
- Dabkowski, P.; and Gal, Y. 2017. Real time image saliency for black box classifiers. In *Advances in Neural Information Processing Systems*, volume 30.
- Dimanov, B.; Bhatt, U.; Jamnik, M.; and Weller, A. 2020. You Shouldn't Trust Me: Learning Models Which Conceal Unfairness From Multiple Explanation Methods. In *Proceedings of the Workshop on Artificial Intelligence Safety, SafeAI@AAAI*, volume 2560 of *CEUR Workshop Proceedings*, 63–73.
- Dong, D.; Wu, H.; He, W.; Yu, D.; and Wang, H. 2015. Multi-task learning for multiple language translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 1723–1732.
- Dong, Y.; Liu, N.; Jalaian, B.; and Li, J. 2022a. Edits: Modeling and mitigating data bias for graph neural networks. In *Proceedings of the ACM Web Conference 2022*, 1259–1269.
- Dong, Y.; Ma, J.; Chen, C.; and Li, J. 2022b. Fairness in Graph Mining: A Survey. *arXiv preprint arXiv:2204.09888*.
- Dwork, C.; Hardt, M.; Pitassi, T.; Reingold, O.; and Zemel, R. 2012. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, 214–226.
- Fay, M.; and Williams, L. 1993. Gender bias and the availability of business loans. *Journal of Business Venturing*, 8(4): 363–376.
- Fu, Z.; Xian, Y.; Gao, R.; Zhao, J.; Huang, Q.; Ge, Y.; Xu, S.; Geng, S.; Shah, C.; Zhang, Y.; et al. 2020. Fairness-aware explainable recommendation over knowledge graphs. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 69–78.
- Fu, Z.; Xian, Y.; Zhu, Y.; Xu, S.; Li, Z.; De Melo, G.; and Zhang, Y. 2021. HOOPS: Human-in-the-Loop Graph Reasoning for Conversational Recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2415–2421.
- Gao, Y.; Sun, T. S.; Bai, G.; Gu, S.; Hong, S. R.; and Liang, Z. 2022. RES: A Robust Framework for Guiding Visual Explanation. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 432–442.
- Guidotti, R.; Monreale, A.; Ruggieri, S.; Turini, F.; Gian-notti, F.; and Pedreschi, D. 2018. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)*, 51(5): 1–42.
- Hanna, R. N.; and Linden, L. L. 2012. Discrimination in grading. *American Economic Journal: Economic Policy*, 4(4): 146–68.
- Hardt, M.; Price, E.; and Srebro, N. 2016. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems*, volume 29.
- Huang, Q.; Yamada, M.; Tian, Y.; Singh, D.; and Chang, Y. 2022. Graphlime: Local interpretable model explanations for graph neural networks. *IEEE Transactions on Knowledge and Data Engineering*.
- Jiang, H.; and Nachum, O. 2020. Identifying and correcting label bias in machine learning. In *International Conference on Artificial Intelligence and Statistics*, 702–712. PMLR.
- Jordan, K. L.; and Freiburger, T. L. 2015. The effect of race/ethnicity on sentencing: Examining sentence type, jail length, and prison length. *Journal of Ethnicity in Criminal Justice*, 13(3): 179–196.
- Kamiran, F.; and Calders, T. 2012. Data preprocessing techniques for classification without discrimination. *Knowledge and information systems*, 33(1): 1–33.
- Kusner, M. J.; Loftus, J.; Russell, C.; and Silva, R. 2017. Counterfactual fairness. *Advances in neural information processing systems*, 30.
- Le Quy, T.; Roy, A.; Iosifidis, V.; Zhang, W.; and Ntoutsis, E. 2022. A survey on datasets for fairness-aware machine learning. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, e1452.
- Lundberg, S. M. 2020. Explaining Quantitative Measures of Fairness. In *Fair & Responsible AI Workshop@ CHI2020*.
- Lundberg, S. M.; and Lee, S.-I. 2017. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, volume 30.
- Mehrabi, N.; Morstatter, F.; Saxena, N.; Lerman, K.; and Galstyan, A. 2021. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6): 1–35.
- Montavon, G.; Lapuschkin, S.; Binder, A.; Samek, W.; and Müller, K.-R. 2017. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern recognition*, 65: 211–222.

Ntoutsis, E.; Fafalios, P.; Gadiraju, U.; Iosifidis, V.; Nejd, W.; Vidal, M.-E.; Ruggieri, S.; Turini, F.; Papadopoulos, S.; Krasanakis, E.; et al. 2020. Bias in data-driven artificial intelligence systems—An introductory survey. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 10(3): e1356.

Obermeyer, Z.; Powers, B.; Vogeli, C.; and Mullainathan, S. 2019. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464): 447–453.

Pan, W.; Cui, S.; Bian, J.; Zhang, C.; and Wang, F. 2021. Explaining algorithmic fairness through fairness-aware causal path decomposition. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 1287–1297.

Pessach, D.; and Shmueli, E. 2022. A Review on Fairness in Machine Learning. *ACM Computing Surveys (CSUR)*, 55(3): 1–44.

Rahman, T.; Surma, B.; Backes, M.; and Zhang, Y. 2019. Fairwalk: Towards Fair Graph Embedding. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, 3289–3295. International Joint Conferences on Artificial Intelligence Organization.

Ras, G.; Xie, N.; van Gerven, M.; and Doran, D. 2022. Explainable deep learning: A field guide for the uninitiated. *Journal of Artificial Intelligence Research*, 73: 329–397.

Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2016. ” Why should i trust you?” Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 1135–1144.

Robnik-Šikonja, M.; and Bohanec, M. 2018. Perturbation-based explanations of prediction models. In *Human and machine learning*, 159–175. Springer.

Selvaraju, R. R.; Cogswell, M.; Das, A.; Vedantam, R.; Parikh, D.; and Batra, D. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, 618–626.

Spinelli, I.; Scardapane, S.; Hussain, A.; and Uncini, A. 2021. Biased edge dropout for enhancing fairness in graph representation learning. *arXiv preprint arXiv:2104.14210*.

Wang, Y.; Zhao, Y.; Dong, Y.; Chen, H.; Li, J.; and Derr, T. 2022a. Improving Fairness in Graph Neural Networks via Mitigating Sensitive Attribute Leakage. *arXiv preprint arXiv:2206.03426*.

Wang, Y.; Zhao, Y.; Zhang, Y.; and Derr, T. 2022b. Collaboration-Aware Graph Convolutional Networks for Recommendation Systems. *arXiv preprint arXiv:2207.06221*.

Yuan, H.; Yu, H.; Gui, S.; and Ji, S. 2020. Explainability in graph neural networks: A taxonomic survey. *arXiv preprint arXiv:2012.15445*.

Zhang, Q.-s.; and Zhu, S.-C. 2018. Visual interpretability for deep learning: a survey. *Frontiers of Information Technology & Electronic Engineering*, 19(1): 27–39.