# Disentangled Representation with Causal Constraints for Counterfactual Fairness

### Ziqi Xu
University of South Australia
Adelaide, Australia
ziqi.xu@mymail.unisa.edu.au

### Jixue Liu
University of South Australia
Adelaide, Australia
jixue.liu@unisa.edu.au

### Debo Cheng
University of South Australia
Adelaide, Australia
dobo.cheng@unisa.edu.au

### Jiuyong Li
University of South Australia
Adelaide, Australia
jiuyong.li@unisa.edu.au

### Lin Liu
University of South Australia
Adelaide, Australia
lin.liu@unisa.edu.au

### Ke Wang
Simon Fraser University
Burnaby, Canada
wangk@cs.sfu.ca

## ABSTRACT

Much research has been devoted to the problem of learning fair representations; however, they do not explicitly the relationship between latent representations. In many real-world applications, there may be causal relationships between latent representations. Furthermore, most fair representation learning methods focus on group-level fairness and are based on correlations, ignoring the causal relationships underlying the data. In this work, we theoretically demonstrate that using the structured representations enable downstream predictive models to achieve counterfactual fairness, and then we propose the Counterfactual Fairness Variational AutoEncoder (CF-VAE) to obtain structured representations with respect to domain knowledge. The experimental results show that the proposed method achieves better fairness and accuracy performance than the benchmark fairness methods.

## 1 INTRODUCTION

Machine learning algorithms have gradually penetrated into our life [34] and have been applied to decision-making for credit scoring [25], crime prediction [23] and loan assessment [10]. The fairness of these decisions and their impact on individuals or society have become an increasing concern. Some extreme unfair incidents have appeared in recent years. For example, COMPAS, a decision support model that estimates the risk of a defendant becoming a recidivist was found to predict higher risk for black people and lower risk for white people [4]; Google Photos are classifying black people as primates [56]; Facebook users receive a recommendation prompt when watching a video featuring blacks, asking them if they'd like to continue to watch videos about primates [31]. These incidents indicate that the machine learning models may become a source of unfairness, which may lead to serious social problems. Since most models are trained with data, which will lead to unfair decisions due to discrimination in the training data. Therefore, the key issue for solving unfair decisions becomes whether we can eliminate these biases embedded in the data through algorithms [27].

To obtain unbiased decisions, many methods [11, 15, 30, 32, 33, 36, 42, 44, 51] are proposed to learn fair representations through two competing goals: encoding data as much as possible, while eliminating any information that transfers through the sensitive attributes. To separate the information from sensitive attributes, various extensions of Variational Autoencoder (VAE) consider minimising the mutual information among latent representations [11, 30, 36, 42].

For example, Creager et al. [11] introduced disentanglement loss into the VAE objective function to decompose observed attributes into sensitive latents and non-sensitive latents to achieve subgroup level fairness; Park et al. [36] improved the above methods and proposed the mutual attribute latent (MAL) to retain only beneficial information for fair predictions.

All existing works of learning fair representations make the assumption that all observed attributes in the real-world can be represented by a number of latent representations. Nevertheless, the latent representations may have causal relationships among them. Let us consider an example where we aim to predict a person's salary using some observed attributes. Following the domain knowledge, we know that a person's salary is determined by two semantic concepts, intelligence and career respectively. We also note that a person's intelligence determines their career with high probability, which can be expressed as a conceptual level causal graph $\mathcal{G}_c$, that is, $Intelligence \rightarrow Career$. Figure 1a shows the causal graph that is learnt from the collected data, while the data itself is biased since the set of sensitive attributes $\mathbf{A}$ can affect the target attribute $Y$.

The existing methods [11, 30] follow Figure 1b to achieve fair predictions. Specifically, this method uses $\mathbf{Z_x}$ to represent the "concept" as mentioned while ensuring $\mathbf{Z_x}$ do not contain sensitive information that transfer through the path $\mathbf{A} \rightarrow \mathbf{X}$. However, this method may not satisfy the domain knowledge since there are causal relationships within these "concepts". Therefore, we need a method as shown in Figure 1c that not only ensures the representation of observed attributes with no sensitive information but also retains causal relationships with respect to domain knowledge. We note that our method builds on the premise that $\mathcal{G}_c$ is available, and we believe this assumption is valid. Fairness issues require humans to guide algorithms, and the causal graph should be given by humans rather than given by machine learning [6]. Compared with the complete version of the causal graph (i.e., a causal graph containing causal relationships between all observed attributes), $\mathcal{G}_c$ only covers the relationship between these "concepts" and is easier to obtain expert consensus.

On the measurement of fairness, all fair representation learning methods use fairness metrics based on correlation, including the VAE-based methods [11, 30, 36, 42]. It is well known that correlation or more generally association does not imply causation. Recent studies [38, 39] have shown that quantifying fairness based on correlation may produce higher deviations. Counterfactual fairness
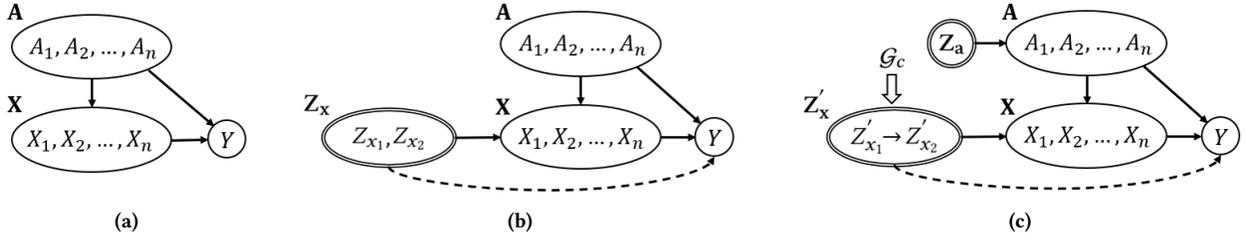
**Figure 1: (a) Causal graph for an example. (b) The process of existing work on learning fair representations to make predictions. (c) The process of our work. A is the set of sensitive attributes; X is the set of other observed attributes; $Z_a$ is the representation of A; $Y$ is the target attribute. The dotted line represent the prediction process that use the learnt representations obtained by the methods. $Z_x$ is the representation of X; $Z'_x$ is the structured representation of X with respect to the conceptual level causal graph $\mathcal{G}_c$. Both $Z_x$ and $Z'_x$ do not contain sensitive information that transform through the path A $\rightarrow$ X since these are learned based on the constraints mentioned in the different methods.**

is a fundamental framework based on causation. With counterfactual fairness, a decision is fair towards an individual if it is the same in the actual world and in the counterfactual world when the individual belonged to a different demographic group.

In this paper, we follow the counterfactual fairness and propose a VAE-based unsupervised fair representation learning method, namely Counterfactual Fairness Variational AutoEncoder (CF-VAE). We take all the observed attributes (except target attribute $Y$) as input, and disentangle the latent representations into $Z_a$ and $Z'_x$. With the causal constraints, $Z'_x$ retains the causal relationships with respect to domain knowledge while containing no sensitive information. We prove that $Z'_x$ is suitable to train the counterfactually fair predictive models. To the best of our knowledge, this work is the first unsupervised method that uses VAE-based techniques to learn the fair representations that enable counterfactual fairness for downstream predictive models. We make the following contributions in this paper:

- We propose CF-VAE, a novel VAE-based unsupervised counterfactual fairness method. CF-VAE can learn structured representations with no sensitive information and retain causal relationships with respect to the conceptual level causal graph determined by domain knowledge.
- We theoretically demonstrate that the structured representations obtained by CF-VAE are suitable for training counterfactually fair predictive models.
- We evaluate the effectiveness of the CF-VAE method on real-world datasets. The experiments show that CF-VAE outperforms existing benchmark fairness methods in both accuracy and fairness.

The rest of this paper is organised as follows. In Section 2, we discuss background knowledge, including our notations. The details of CF-VAE are shown in Section 3. In Section 4, we discuss the experiment results. In Section 5, we discuss related works. Finally, we conclude this paper in Section 6.

## 2 BACKGROUND

We use upper case letters to represent attributes and boldfaced upper case letters to denote the set of attributes. We use boldfaced lower case letters to represent the values of the set of attributes. The values of attributes are represented using lower case letters.

Let $\mathbf{A}$ be the set of sensitive attributes, which should not be used for predictive models; $\mathbf{X}$ be the set of other observed attributes, which may have causal relationships with $\mathbf{A}$; $\mathbf{V}$ be the set of all observed attributes, i.e., $\mathbf{V} = \{\mathbf{A}, \mathbf{X}\}$; $Y$ be the target attribute that may have causal relationships with attributes in $\mathbf{A}$ and $\mathbf{X}$. We use $\widehat{Y}(\cdot)$ to represent the predictor.

$\mathcal{G}_c$ is the conceptual level causal graph and represents domain knowledge. The nodes shown in $\mathcal{G}_c$ are "concepts", each of which represents a set of observed attributes that have similar meanings. Each "concept" has causal relationships with the other "concepts". For example, *Intelligence* is a "concept" in $\mathcal{G}_c$ and it may represent several observed attributes that have similar meanings, including *GPA*, *Education level* and *Major*.

We define that $\mathbf{Z_a}$ is the representation of $\mathbf{A}$; $\mathbf{Z_x}$ is the representation of $\mathbf{X}$ without embedding causal relationships; $\mathbf{Z'_x}$ is a structured version of $\mathbf{Z_x}$ under the causal constraints of domain knowledge and does not contain sensitive information.

### 2.1 Counterfactual Fairness

In this paper, a causal graph is used to represent a causal mechanism. In a causal graph, a directed edge, such as $V_j \rightarrow V_i$ denotes that $V_j$ is a parent (i.e., direct cause) and we use $pa_i$ to denote the set of parents of $V_i$. We follow Pearl's [39] notation and define a causal model as a triple $(\mathbf{U}, \mathbf{V}, \mathbf{F})$: $\mathbf{U}$ is a set of the latent background attributes, which are the factors not caused by any attributes in the set $\mathbf{V} = \{\mathbf{A}, \mathbf{X}\}$; $\mathbf{F}$ is a set of deterministic functions, $V_i = f_i(pa_i, U_{pa_i})$, such that $pa_i \subseteq \mathbf{V} \setminus \{V_i\}$ and $U_{pa_i} \subseteq \mathbf{U}$. Such equations are also known as structural equations [3]. Besides, some commonly used definitions in graphical causal modelling, such as faithfulness, $d$-separation and causal path can be found in [39, 40, 43].

With the causal model $(\mathbf{U}, \mathbf{V}, \mathbf{F})$, we have the following definition of counterfactual fairness:

DEFINITION 1 (COUNTERFACTUAL FAIRNESS [26]). *Predictor $\widehat{Y}(\cdot)$ is counterfactually fair if under any context $\mathbf{X} = \mathbf{x}$ and $\mathbf{A} = \mathbf{a}$,*

$$P(\widehat{Y}_{\mathbf{A}\leftarrow\mathbf{a}}(\mathbf{U}) = y \mid \mathbf{X} = \mathbf{x}, \mathbf{A} = \mathbf{a}) =$$
$$P(\widehat{Y}_{\mathbf{A}\leftarrow\bar{\mathbf{a}}}(\mathbf{U}) = y \mid \mathbf{X} = \mathbf{x}, \mathbf{A} = \mathbf{a}), \quad (1)$$

*for all y and any value ā attainable by* A.

Counterfactual fairness is considered to be related to individual fairness [26]. Individual fairness means that similar individuals should receive similar predicted outcomes. The concept of individual fairness when measuring the similarity of the individual is unknowable, which is similar to the unknowable distance between the real-world and the counterfactual world in counterfactual fairness [28].

## 2.2 Variational Autoencoder

Variational Autoencoder (VAE) was proposed by Kingma and Welling [24], which was originally applied to image dimensionality reduction. The objective of VAE is to maximise the Evidence Lower Bound (ELBO) $\mathcal{M}$, and derived as follows:

$$\log p(\mathbf{V}) \geq \mathbb{E}_{q(\mathbf{Z}|\mathbf{V})}[\log p(\mathbf{V}|\mathbf{Z}) + \log p(\mathbf{Z}) - \log q(\mathbf{V}|\mathbf{Z})]$$
$$=: \mathcal{M}_{\text{VAE}}, \quad (2)$$

which can also be rewritten as follows:

$$\mathcal{M}_{\text{VAE}} = \mathbb{E}_{q(\mathbf{Z}|\mathbf{V})}[\log p(\mathbf{V}|\mathbf{Z})] - D_{KL}[q(\mathbf{Z}|\mathbf{V})||p(\mathbf{Z})], \quad (3)$$

where $\mathbf{V}$ denotes the set of all observed attributes and $\mathbf{Z}$ denotes the set of learnt representations. The encoder $q$ encodes $\mathbf{V}$ into $\mathbf{Z}$, and the decoder $p$ reconstructs $\mathbf{V}$ from $\mathbf{Z}$.

The first part in Equation 3 can be considered as reconstruction error, i.e., the loss between the reconstructed and the original $\mathbf{V}$. The second part is the distribution distance between the Gaussian prior $p(\mathbf{Z}) = \mathcal{N}(0, \mathbf{I})$ and the $\mathbf{Z}$ encoded with $\mathbf{V}$. The training process of a VAE is to learn the parameters in $q$ and $p$ through the neural networks.

Higgins et al. [18] modified the above VAE objective function by adding a hyperparameter $\beta$ that balances latent channel capacity and independence constraints with reconstruction accuracy. Then, they devised a protocol to quantitatively compare the degree of disentanglement learnt by different models and argued that each dimension of a correctly disentangled representation should capture no more than one semantically meaningful concept. The ELBO of $\beta$-VAE is defined as:

$$\mathcal{M}_{\beta\text{-VAE}} = \mathbb{E}_{q(\mathbf{Z}|\mathbf{V})}[\log p(\mathbf{V}|\mathbf{Z})] - \beta D_{KL}[q(\mathbf{Z}|\mathbf{V})||p(\mathbf{Z})]. \quad (4)$$

Kim and Mnih [21] showed that $\mathbf{Z}$ can be considered as disentangled if each attribute in $\mathbf{Z}$, denoted as $Z_i$ is independent of each other. They minimised the total collection [45] of the latent representations as follows, and guaranteed disentanglement.

$$\text{Total Correction} = D_{KL}[q(\mathbf{Z})||\prod_{i=1}^{D_{\mathbf{Z}}} q(Z_i)], \quad (5)$$

where $D_{\mathbf{Z}}$ is dimension of $\mathbf{Z}$. They proposed Factor-VAE by using total correction and the ELBO of Factor-VAE is defined as:

$$\mathcal{M}_{\text{Factor-VAE}} = \mathcal{M}_{\text{VAE}} - \gamma D_{KL}[q(\mathbf{Z})||\prod_{i=1}^{D_{\mathbf{Z}}} q(Z_i)]. \quad (6)$$

## 3 PROPOSED METHOD

In this section, we first theoretically demonstrate that learning counterfactually fair representations are feasible. Then, we propose the Counterfactual Fairness Variational AutoEncoder (CF-VAE) to obtain the structured representations for predictors to achieve counterfactual fairness.

## 3.1 The Theory of Learning Counterfactually Fair Representations

We discuss what types of representations enable downstream predictive models to achieve counterfactual fairness. Following the work in [16], we define the three steps for counterfactual inference.

DEFINITION 2 (COUNTERFACTUAL INFERENCE [16]). *Given a causal model* $(\mathbf{U}, \mathbf{V}, \mathbf{F})$ *and evidence* $\mathbf{W}$, *where* $\mathbf{W} \subset \mathbf{V}$, *the counterfactual inference is the computation of probabilities* $P(Y_{\mathbf{A}\leftarrow\mathbf{a}}(\mathbf{U}|\mathbf{W} = \mathbf{w}))$.

- **Abduction**: *for a given prior on* $\mathbf{U}$, *compute the posterior distribution of* $\mathbf{U}$ *given the evidence* $\mathbf{W} = \mathbf{w}$;
- **Action**: *substitute the equations for* $\mathbf{A}$ *with the interventional values* $\mathbf{a}$, *resulting in the modified set of equations* $\mathbf{F}_{\mathbf{a}}$;
- **Prediction**: *compute the implied distribution on the remaining elements of* $\mathbf{V}$ *using* $\mathbf{F}_{\mathbf{a}}$ *and the posterior* $P(\mathbf{U}|\mathbf{W} = \mathbf{w})$.

Following the work in [26], the implication of counterfactual fairness is described as follows:

PROPOSITION 3.1 (IMPLICATION OF COUNTERFACTUAL FAIRNESS [26]). *Let* $\mathcal{G}$ *be the causal graph of the given model* $(\mathbf{U}, \mathbf{V}, \mathbf{F})$. *If there exists* $\mathbf{W}$ *be any non-descendant of* $\mathbf{A}$, *then downstream predictor* $\widehat{Y}(\mathbf{W})$ *will be counterfactually fair.*

We extend Proposition 3.1 to the fair representation learning and present the following theorem. We follow the similar proof process in work [26] to prove this theorem.

THEOREM 3.2. *Given the causal graph* $\mathcal{G}$, $\mathbf{Z}_{\mathbf{a}}$ *is the representation of sensitive attributes* $\mathbf{A}$, $\mathbf{Z}'_{\mathbf{x}}$ *is the structured representation of the other observed attributes* $\mathbf{X}$ *with respect to the conceptual level causal graph* $\mathcal{G}_c$. *We have* $\widehat{Y}(\mathbf{Z}'_{\mathbf{x}})$ *satisfy counterfactual fairness.*

PROOF. Given the causal graph $\mathcal{G}$ as shown in Figure 2, there is not a parent node of $\mathbf{A}$ in $\mathbf{X}$, and there is not a child node of $Y$ in $\mathbf{X}$. $\mathbf{X}$ contains four subsets: $\mathbf{X}_Y^{\mathbf{A}}$ is the subset of other observed attributes that are descendants of $\mathbf{A}$ and parents of $Y$; $\mathbf{X}_Y^{\mathbf{N}}$ is the subset of other observed attributes that are only parents of $Y$; $\mathbf{X}_{\mathbf{N}}^{\mathbf{N}}$ is the subset of other observed attributes that are no relationships with $\mathbf{A}$ and $Y$; $\mathbf{X}_{\mathbf{N}}^{\mathbf{A}}$ is the subset of other observed attributes that are only descendants of $\mathbf{A}$. After perfect representation learning, we obtain $\mathbf{Z}_{\mathbf{a}}$ and $\mathbf{Z}'_{\mathbf{x}}$.

We proof that $\mathbf{Z}'_{\mathbf{x}}$ is not the descendant of $\mathbf{A}$ with the following two subsets. For the first subsets $\{\mathbf{X}_Y^{\mathbf{A}}, \mathbf{X}_Y^{\mathbf{N}}, \mathbf{X}_{\mathbf{N}}^{\mathbf{A}}\}$, there are seven paths between $\mathbf{A}$ and $\mathbf{Z}'_{\mathbf{x}}$, including $\mathbf{A} \rightarrow \mathbf{X}_Y^{\mathbf{A}} \leftarrow \mathbf{Z}'_{\mathbf{x}}$, $\mathbf{A} \rightarrow \mathbf{X}_Y^{\mathbf{A}} \rightarrow Y \leftarrow \mathbf{Z}'_{\mathbf{x}}$, $\mathbf{A} \rightarrow \mathbf{X}_Y^{\mathbf{N}} \rightarrow Y \leftarrow \mathbf{X}_Y^{\mathbf{N}} \leftarrow \mathbf{Z}'_{\mathbf{x}}$, $\mathbf{A} \rightarrow Y \leftarrow \mathbf{X}_Y^{\mathbf{A}} \leftarrow \mathbf{Z}'_{\mathbf{x}}$, $\mathbf{A} \rightarrow Y \leftarrow \mathbf{Z}'_{\mathbf{x}}$, $\mathbf{A} \rightarrow Y \leftarrow \mathbf{X}_Y^{\mathbf{N}} \leftarrow \mathbf{Z}'_{\mathbf{x}}$ and $\mathbf{A} \rightarrow \mathbf{X}_{\mathbf{N}}^{\mathbf{A}} \leftarrow Y$. These seven paths are blocked by $\emptyset$ (i.e., $\mathbf{A}$ and $\mathbf{Z}'_{\mathbf{x}}$ are $d$-separated by $\emptyset$), since each path contains a collider either $\mathbf{X}_Y^{\mathbf{A}}$ or $Y$ or $\mathbf{X}_{\mathbf{N}}^{\mathbf{A}}$. For second subset $\mathbf{X}_{\mathbf{N}}^{\mathbf{N}}$, there is no path connecting $\mathbf{X}_{\mathbf{N}}^{\mathbf{N}}$ and $Y$. Hence, $\mathbf{Z}'_{\mathbf{x}}$ is not the descendant of $\mathbf{A}$. Therefore, $\widehat{Y}(\mathbf{Z}'_{\mathbf{x}})$ is counterfactually fair based on Proposition 3.1. □

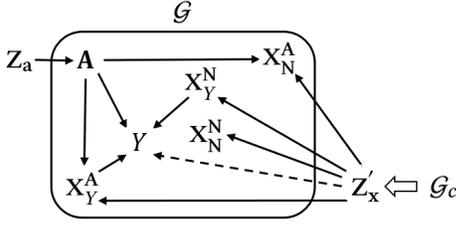We use Figure 2 to show whether the following predictors satisfy counterfactual fairness.

**Figure 2:** $\mathcal{G}$ **is the causal graph that represents the causal relationship between A, X** $= \{X_Y^A, X_Y^N, X_N^A, X_N^N\}$ **and** $Y$. **The dotted line represents the prediction process that uses** $Z_x'$.

- $\widehat{Y}(A, X)$: This model is unfair since it uses sensitive attributes to make prediction.
- $\widehat{Y}(X)$: This model satisfies fairness through awareness [13] but fails to achieve counterfactual fairness. The predictor $\widehat{Y}(X)$ does not use sensitive attributes explicitly, but it uses $X_Y^A$ and $X_N^A$ which are the descendants of A.
- $\widehat{Y}(Z_a, Z_x')$: This model is unfair because it uses sensitive attributes for prediction. The reason is that $Z_a$ is the representation of A, which should be consider as sensitive attributes either.
- $\widehat{Y}(X_Y^N, X_N^N)$: This model satisfies counterfactual fairness since both $X_Y^N$ and $X_N^N$ are non-descendants of A. However, this predictor losses a lot of useful information that embeds in other observed attributes, which means it may not achieve an acceptable prediction accuracy.
- $\widehat{Y}(Z_x')$: This model is counterfactually fair based on Theorem 3.2 and achieves higher accuracy than $\widehat{Y}(X_Y^N, X_N^N)$ as shown in our experiments.

## 3.2 CF-VAE

We first discuss the causal constraints and then explain the loss function of CF-VAE in detail. The architecture of CF-VAE is shown in Figure 3.

*3.2.1 Learning Representations with Causal Constraints.* We aim to retain causal relationships between "concepts" through a more easily accessible conceptual level causal graph $\mathcal{G}_c$ and embed these relationships in representations. Following the works in [48, 57], we transform these causal relationships in the form of an adjacency matrix C as causal constraints to construct $Z_x'$ and feed them to predictive models.

Many researchers study fair decision-making under the method of causality, in which an accessible causal graph is an important assumption. Zhang et al. [54] used the PC algorithm [19, 43] to learn the causal relationships of the dataset itself and used the causal graph to restrict the transfer of sensitive information along the specific paths. Learning causal graphs from observational data has been shown to be feasible in unbiased data, but in fairness computing, the dataset may be biased and the PC algorithm may not be suitable. Other researches [8, 35] assume that the complete version of causal graph is accessible and evolved from domain knowledge. In our paper, we adopt the second approach and further weaken it. We assume that $\mathcal{G}_c$ only covers the relationships between "concepts",

not all observed attributes, which is easy to obtain consensus from experts.

To formalise causal relationships, we consider $n$ "concepts" in the dataset, which means $Z_x'$ should have the same dimension as "concepts". The "concepts" in observations are causally structured by $\mathcal{G}_c$ with an adjacency matrix C. For simplicity, in this paper, the causal constraints are exactly implemented by a linear structural equation model as follows:

$$Z_x' = (I - C^T)^{-1} Z_x, \tag{7}$$

where I is the identity matrix, $Z_x$ is obtained from the encoder, $Z_x'$ is constructed from $Z_x$ and C. C is obtained from $\mathcal{G}_c$ with respect to domain knowledge. The parameters in C indicate that there are corresponding edges, and the values of the parameters indicate the weight of the causal relationships. It is worth noting that if the parameter value is zero, it means that such an edge does not exist, i.e., no causal relationship between these two "concepts".

As mentioned above, $Z_x$ is obtained from the encoder, we cannot guarantee that each attribute inside is independent. To ensure the independence of each attribute in $Z_x$, we follow the Factor-VAE in work [21] and employ the total correction regularisation (TCR) in our loss function. TCR also encourages the correctness of structured $Z_x'$ with respect to domain knowledge since there are no correlations in $Z_x$ before adding causal constraints. The TCR for our proposed CF-VAE is defined as:

$$\mathcal{L}_{TCR} = \gamma D_{KL}[q(Z_x) || \prod_{i=1}^{D_{Z_x}} q(Z_{x_i})], \tag{8}$$

where $\gamma$ is the weight value, $D_{Z_x}$ is dimension of $Z_x$.

*3.2.2 Learning Strategy.* We first explain the architecture of CF-VAE without using causal constraints. Then, we add causal constraints and orthogonality promoting regularisation (OPR) to obtain the loss function of CF-VAE.

In the inference model, the variational approximations of the posteriors are defined as:

$$q(Z_a|A) = \prod_{i=1}^{D_{Z_a}} \mathcal{N}(\mu = \hat{\mu}_{Z_{a_i}}, \sigma^2 = \hat{\sigma}_{Z_{a_i}}^2);$$

$$q(Z_x|X) = \prod_{i=1}^{D_{Z_x}} \mathcal{N}(\mu = \hat{\mu}_{Z_{x_i}}, \sigma^2 = \hat{\sigma}_{Z_{x_i}}^2), \tag{9}$$

where $\hat{\mu}_{a_i}, \hat{\mu}_{x_i}$ and $\hat{\sigma}_{a_i}^2, \hat{\sigma}_{x_i}^2$ are the means and variances of the Gaussian distributions parameterised by neural networks.

The generative model for A and X are defined as:

$$p(A|Z_a) = \prod_{i=1}^{D_A} p(A_i|Z_a); \; p(X|Z_x) = \prod_{i=1}^{D_X} p(X_i|Z_x), \tag{10}$$

where $D_A$ and $D_X$ are dimensions of A and X.

Following the setting in VAE [24], we choose Gaussian distribution as prior distributions, which are defined as:

$$p(Z_a) = \prod_{i=1}^{D_{Z_a}} \mathcal{N}(Z_{a_i}|0, 1); \; p(Z_x) = \prod_{i=1}^{D_{Z_x}} \mathcal{N}(Z_{x_i}|0, 1). \tag{11}$$

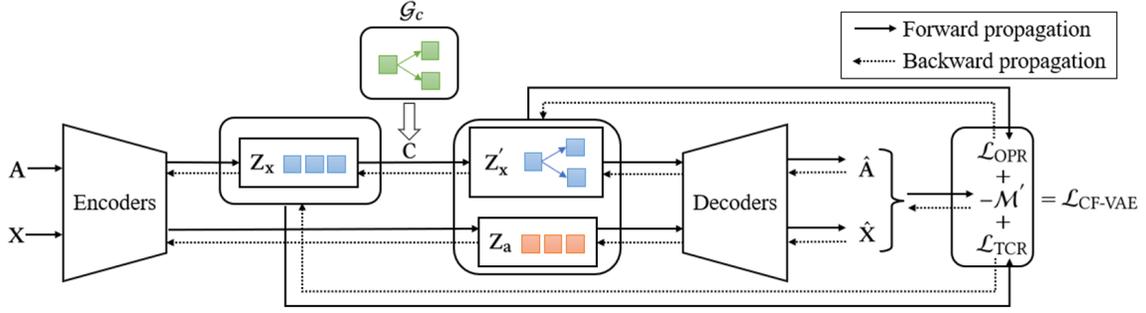Given the training samples, the parameters can be optimised by maximising the following ELBO:

**Figure 3: The architecture of CF-VAE. The adjacency Matrix C is used to construct $Z'_\mathbf{x}$ and is determined by the conceptual level causal graph $\mathcal{G}_c$ with respect to domain knowledge. $\mathcal{L}_{\text{TCR}}$ is used to ensure that each attribute in $Z_\mathbf{x}$ is independent of each other. $\mathcal{L}_{\text{OPR}}$ is used to encourage that $Z'_\mathbf{x}$ do not contain sensitive information. The loss function of CF-VAE is $\mathcal{L}_{\text{CF-VAE}}$.**

$$\mathcal{M} = \mathbb{E}_{q(\mathbf{Z_a}|\mathbf{A})}[\log p(\mathbf{A}|\mathbf{Z_a})] + \mathbb{E}_{q(\mathbf{Z_x}|\mathbf{X})}[\log p(\mathbf{X}|\mathbf{Z_x})] \\ - D_{KL}[q(\mathbf{Z_a}|\mathbf{A})||p(\mathbf{Z_a})] - D_{KL}[q(\mathbf{Z_x}|\mathbf{X})||p(\mathbf{Z_x})]. \tag{12}$$

We note that Equation 12 is not under causal constraints and still using $\mathbf{Z_x}$ to optimise. The $\mathcal{M}$ comprises four terms: the first and second term denote the reconstruction loss between the original $\{\mathbf{A}, \mathbf{X}\}$ and $\{\hat{\mathbf{A}}, \hat{\mathbf{X}}\}$; the third term and the fourth term are used for calculating the distribution distance between the prior knowledge and the latent representations that we obtained.

We follow Section 3.2.1 and add causal constraints in Equation 12. The updated ELBO is defined as:

$$\mathcal{M}' = \mathbb{E}_{q(\mathbf{Z_a}|\mathbf{A})}[\log p(\mathbf{A}|\mathbf{Z_a})] + \mathbb{E}_{q(\mathbf{Z'_x}|\mathbf{X})}[\log p(\mathbf{X}|\mathbf{Z'_x})] \\ - D_{KL}[q(\mathbf{Z_a}|\mathbf{A})||p(\mathbf{Z_a})] - D_{KL}[q(\mathbf{Z'_x}|\mathbf{X})||p(\mathbf{Z'_x})], \tag{13}$$

where

$$p(\mathbf{Z'_x}) = (\mathbf{I} - \mathbf{C}^T)^{-1} p(\mathbf{Z_x}); \ p(\mathbf{X}|\mathbf{Z'_x}) = \prod_{i=1}^{D_\mathbf{X}} p(X_i|\mathbf{Z'_x});$$

$$q(\mathbf{Z'_x}|\mathbf{X}) = \prod_{i=1}^{D_{Z'_\mathbf{x}}} \mathcal{N}(\mu = \hat{\mu}_{Z'_{\mathbf{x}_i}}, \sigma^2 = \hat{\sigma}^2_{Z'_{\mathbf{x}_i}}).$$

We introduce orthogonality to encourage disentanglement between $\mathbf{Z_a}$ and $\mathbf{Z'_x}$. Following the work in [49], we employ orthogonality promoting regularisation based on the pairwise cosine similarity among latent representations: if the cosine similarity is close to zero, then the latent representations are closer to being orthogonal and independent. The cosine similarity (CS) is defined as:

$$CS(\mathbf{E_1}, \mathbf{E_2}) = \frac{\mathbf{E_1}^T \mathbf{E_2}}{\|\mathbf{E_1}\|_2 \ \|\mathbf{E_2}\|_2}, \tag{14}$$

where $\|\cdot\|_2$ is the $l_2$ norm.

To encourage orthogonality between two vectors $\mathbf{E_1}$ and $\mathbf{E_2}$, we can make their inner product $\mathbf{E_1}^T \mathbf{E_2}$ close to zero and their $l_2$ norm $\|\mathbf{E_1}\|_2$, $\|\mathbf{E_2}\|_2$ close to one [47]. The orthogonality promoting regularisation (OPR) for our proposed CF-VAE is defined as:

$$\mathcal{L}_{\text{OPR}} = \frac{1}{B} \sum_{i=1}^{B} CS(\mathbf{Z_{a_i}}, \mathbf{Z'_{x_i}}), \tag{15}$$

where $B$ denotes the batch size for neural network.

In conclusion, the loss function of our proposed CF-VAE is defined as:

$$\mathcal{L}_{\text{CF-VAE}} = -\mathcal{M}' + \mathcal{L}_{\text{TCR}} + \mathcal{L}_{\text{OPR}}. \tag{16}$$

## 4 EXPERIMENTS

In this section, we conduct extensive experiments to evaluate CF-VAE on real-world datasets. Before showing the detailed results, we first present the details of selected methods and the evaluation metrics.

### 4.1 Framework Comparison

The proposed CF-VAE is considered as a pre-processing technique to address fairness issues since it obtains structured representations for downstream predictive models to achieve counterfactual fairness. Hence, we compare CF-VAE with traditional and VAE-based pre-processing methods. For traditional methods, we select baselines including ReWeighting (RW) [20], Disparate Impart Remover (DIR) [14] and Optimized Preprocessing (OP) [5]. Both of them are available in AIF360 [1]. For VAE-based methods, we compare with VFAE [30] and FFVAE [11]. Both of them are implemented in Pytorch [37]. We also obtain the Full model for comparison, which uses all attributes in the dataset to make predictions.

We do not choose the basic VAE (e.g., VAE [24], $\beta$-VAE [18], and Factor-VAE [21]) for comparison in this experiment, since they are not optimised for fairness problems. In addition, we do not use VAE-based inference models [8, 22, 41, 48] for comparison, because the purpose of these inference models is to generate counterfactual data or to estimate effects, which is different from our goals.

We select several well-known predictive models to simulate the downstream prediction process. Linear Regression ($\text{LR}_\text{R}$), Stochastic Gradient Descent Regression ($\text{SGD}_\text{R}$) and Multi-layer Perceptron Regression ($\text{MLP}_\text{R}$) are used for regression tasks; Logistic Regression ($\text{LR}_\text{C}$), Stochastic Gradient Descent Classification ($\text{SGD}_\text{C}$) and Multi-layer Perceptron Classification ($\text{MLP}_\text{C}$) are used for classification tasks. For each predictive model, we run 10 times and record the mean and error of the results for evaluation metrics, which are explained in detail in Section 4.2.
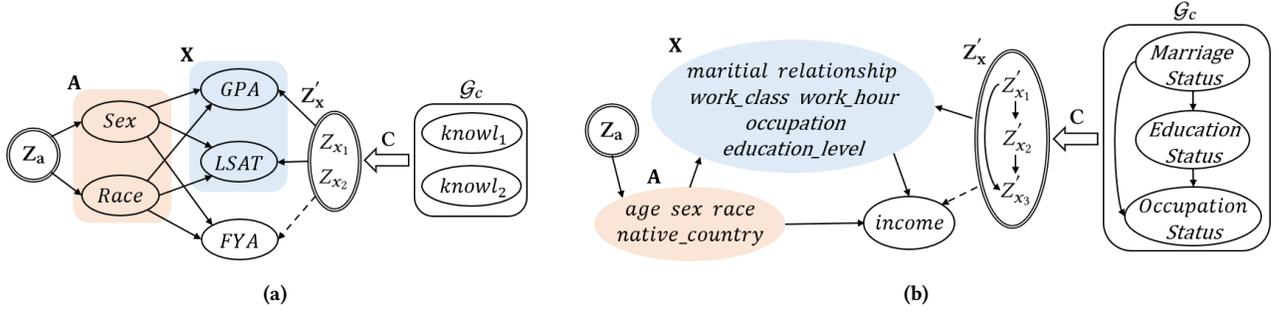
**Figure 4: (a) The process of CF-VAE for Law school dataset. (b) The process of CF-VAE for Adult dataset. $Z_a$ is the representation of A; $Z'_x$ is the structured representation of X. The adjacency matrix C is used to construct $Z'_x$ with respect to $\mathcal{G}_c$. The dotted line represents the prediction process of $Y$ that uses $Z'_x$.**

## 4.2 Evaluation Metrics

*4.2.1 Fairness.* There are no metrics to quantify counterfactual fairness since we can only obtain real-world data. We propose the situation test to measure fairness for different predictive models. The situation test has already been widely used in United States to detect individual discrimination [2]. In our experiment, we construct a matched pair for each individual by inverting the values of sensitive attributes. We take this matched pair as the input to the predictive model, and the predictive model is fair if the predictions of the matched pair are the same as the original pair.

We define unfairness score (UFS) to measure the result of the situation test. Specifically, the form of score differs for different predictive models. For regression tasks, we define $UFS_R$ that measure the bias between prediction results for the matched pair and the original pair; For classification tasks, $UFS_C$ is defined as how many individuals' prediction results are changed after intervening the values of sensitive attributes. $UFS_R$ and $UFS_C$ are described as follows:

$$UFS_R = \sqrt{\frac{1}{N}\sum_{i=1}^{N}\left(\widehat{Y}_{A\leftarrow a}(Z'_{x_i}) - \widehat{Y}_{A\leftarrow \bar{a}}(Z'_{x_i})\right)^2} \;;$$

$$UFS_C = \frac{1}{N}\sum_{i=1}^{N} xor\left(\widehat{Y}_{A\leftarrow a}(Z'_{x_i}), \widehat{Y}_{A\leftarrow \bar{a}}(Z'_{x_i})\right),$$

(17)

where $N$ is the number of samples for evaluation.

The lower UFS value means that the predictive models achieve higher individual fairness.

*4.2.2 Accuracy.* We evaluate the performance on prediction with the following metrics. For regression tasks, we use Root Mean Square Error (RMSE) to compare the error between prediction results and target attributes' values. For classification tasks, we use accuracy to evaluate various predictive models.

## 4.3 Law School

The law school dataset comes from a survey [46] of admissions information from 163 law schools in the United States. It contains information of 21,790 law students, including their entrance exam scores (LSAT), their grade point average (GPA) collected prior to law school, and their first-year average grade (FYA). The school

expects to predict if the applicants will have a high FYA. Gender and race are sensitive attributes in this dataset, and the school also wants to ensure that predictions are not affected by sensitive attributes. However, LSAT, GPA and FYA scores may be biased due to socio-environmental factors. The process of CF-VAE for the Law school dataset is shown in Figure 4a.

*4.3.1 Implementation Details.* We divide the Law school dataset into 70% training set for training the representation models, 30% testing set for evaluating the accuracy of the predictive models, and inverting the values of sensitive attributes in the testing set to generate the auditing set for evaluating the fairness of the predictive models.

We use the same $\mathcal{G}_c$ as shown in work [26] to model latent "concepts" of *GPA* and *LSAT*. Since $knowl_1$ and $knowl_2$ have no causal relationship, the parameters in adjacency matrix C are set to zero. As a results, we set the $D_{Z'_x} = 2$ and set the weight value in $\mathcal{L}_{TCR}$ as $\gamma = 10$.

*4.3.2 Fairness.* The purpose is to demonstrate our method can achieve better fairness performance than other VAE-based methods. As shown in Table 1, since the Full model uses sensitive attributes to make predictions, inverting sensitive attributes has the highest impact on the individual's prediction results, which means that the model is unfair. RW, DIR and OP achieves fair predictions by modifying the dataset compared to the Full model. Both VFAE and FFVAE disentangle the sensitive attributes with latent representations, so the influence of inverting the sensitive attributes on the prediction results is small. Our method achieves the lowest $UFS_R$, 0.013, 0.025, and 0.044 for $LR_R$, $SGD_R$, and $MLP_R$ respectively, which means CF-VAE disentangle $Z'_x$ and $Z_a$ more precisely.

*4.3.3 Accuracy.* The accuracy results are shown in Table 1. The Full model is unfair and it uses sensitive information to more accurately predict FYA and thus achieves the highest accuracy. The proposed CF-VAE achieves the best fairness aware accuracy in all predictive models than other methods. Our method not only achieves counterfactual fairness for downstream predictors but also flexible for choosing predictive models.

**Table 1: The results for Law School dataset. The best fairness aware RMSE and the best $UFS_R$ are shown in bold.**

| Model | Accuracy (RMSE) ↓ | | | Fairness ($UFS_R$) ↓ | | |
|---|---|---|---|---|---|---|
| | $LR_R$ | $SGD_R$ | $MLP_R$ | $LR_R$ | $SGD_R$ | $MLP_R$ |
| Full | 0.865 ± 0.007 | 0.867 ± 0.007 | 0.865 ± 0.007 | 0.660 ± 0.019 | 0.762 ± 0.019 | 0.760 ± 0.045 |
| RW | 0.955 ± 0.013 | 0.956 ± 0.012 | 0.953 ± 0.012 | 0.067 ± 0.002 | 0.067 ± 0.001 | 0.079 ± 0.003 |
| DIR | 0.943 ± 0.009 | 0.944 ± 0.009 | 0.941 ± 0.010 | 0.060 ± 0.001 | 0.060 ± 0.001 | 0.070 ± 0.002 |
| OP | 0.959 ± 0.011 | 0.960 ± 0.011 | 0.956 ± 0.010 | 0.047 ± 0.001 | 0.046 ± 0.001 | 0.055 ± 0.003 |
| VFAE | 0.932 ± 0.007 | 0.933 ± 0.007 | 0.934 ± 0.007 | 0.035 ± 0.010 | 0.074 ± 0.017 | 0.096 ± 0.010 |
| FFVAE | 0.933 ± 0.005 | 0.934 ± 0.004 | 0.935 ± 0.005 | 0.032 ± 0.007 | 0.060 ± 0.022 | 0.097 ± 0.008 |
| CF-VAE | **0.931 ± 0.006** | **0.932 ± 0.006** | **0.932 ± 0.006** | **0.013 ± 0.006** | **0.025 ± 0.011** | **0.044 ± 0.006** |

**Table 2: The results for Adult dataset. The best fairness aware accuracy and the best $UFS_C$ are shown in bold.**

| Model | Accuracy ↑ | | | Fairness ($UFS_C$) ↓ | | |
|---|---|---|---|---|---|---|
| | $LR_C$ | $SGD_C$ | $MLP_C$ | $LR_C$ | $SGD_C$ | $MLP_C$ |
| Full | 0.802 ± 0.002 | 0.803 ± 0.004 | 0.831 ± 0.004 | 0.068 ± 0.003 | 0.060 ± 0.018 | 0.034 ± 0.009 |
| RW | 0.797 ± 0.001 | 0.792 ± 0.002 | 0.819 ± 0.001 | 0.038 ± 0.001 | 0.029 ± 0.002 | 0.052 ± 0.001 |
| DIR | 0.800 ± 0.001 | 0.793 ± 0.003 | 0.817 ± 0.001 | 0.035 ± 0.001 | 0.027 ± 0.002 | 0.046 ± 0.001 |
| OP | 0.780 ± 0.002 | 0.779 ± 0.003 | 0.783 ± 0.002 | 0.032 ± 0.003 | 0.030 ± 0.004 | 0.033 ± 0.005 |
| VFAE | 0.785 ± 0.001 | 0.781 ± 0.003 | 0.819 ± 0.004 | 0.062 ± 0.002 | 0.041 ± 0.010 | 0.025 ± 0.003 |
| FFVAE | 0.785 ± 0.003 | 0.782 ± 0.001 | 0.814 ± 0.005 | 0.062 ± 0.001 | 0.044 ± 0.010 | 0.032 ± 0.010 |
| CF-VAE | **0.801 ± 0.002** | **0.794 ± 0.004** | **0.820 ± 0.002** | **0.031 ± 0.002** | **0.020 ± 0.006** | **0.024 ± 0.004** |

## 4.4 Adult

The Adult dataset comes from the UCI repository [12] contains 14 attributes including race, age, education information, marital information as well as capital gain and loss for 48,842 individuals. The process of CF-VAE is shown in Figure 4b.

*4.4.1 Implementation Details.* We pre-process the dataset by deleting missing information and encoding discrete attributes. After that, we get 45,222 individuals and the downstream tasks' goal is to predict whether the individual's income is above $50,000. We set *race*, *age*, *sex* and *native country* as **A**; *marital*, *relationship*, *work class*, *work hour*, *occupation* and *education level* as **X**. We divide the Adult dataset into 70% training set for training representation models, 30% testing set for evaluating the accuracy of the predictive models, and select 10,000 individuals with female that income below 50K as the auditing set.

We use the same $\mathcal{G}_c$ as shown in previous research [8, 35] to model the latent "concepts". We set the $D_{Z'_x} = 3$ and set the weight value in $\mathcal{L}_{TCR}$ as $\gamma = 10$. The adjacency matrix **C** is defined as:

$$\mathbf{C} = \begin{vmatrix} 0 & \lambda_{12} & \lambda_{13} \\ 0 & 0 & \lambda_{23} \\ 0 & 0 & 0 \end{vmatrix} \tag{18}$$

Then, we construct $\mathbf{Z'_x}$ from $\mathbf{Z_x}$ and **C** as follows:

$$Z'_{x_1} = Z_{x_1}; \; Z'_{x_2} = \lambda_{12} Z_{x_1} + Z_{x_2}; \\ Z'_{x_3} = \lambda_{13} Z_{x_1} + \lambda_{23} Z_{x_2} + Z_{x_3}. \tag{19}$$

We set parameter $\{\lambda_{12} = 1, \lambda_{13} = 1, \lambda_{23} = 1\}$ to denote that edges within latent representations, i.e., $Z'_{x_1} \rightarrow Z'_{x_2}, Z'_{x_1} \rightarrow Z'_{x_3}, Z'_{x_2} \rightarrow Z'_{x_3}$.

*4.4.2 Fairness.* The fairness results are shown in Table 2, the Full model achieves the worst $UFS_C$, since it use **A** to predict *income*. Both baseline fairness models and other VAE-based methods improve fairness to a certain extent. The proposed CF-VAE achieves the best $UFS_C$, only 3.1%, 2.0% and 2.4% of individuals' results are affected by sensitive attributes' values inversions in $LR_C$, $SGD_C$ and $MLP_C$, respectively. Our method achieves better fairness performance than other methods, since it remains causal relationships in latent representations with respect to $\mathcal{G}_c$ and disentangles structured representations with sensitive attributes.

*4.4.3 Accuracy.* The Full model uses all observed attributes for predictions. It is worth noting that the Full model does not achieve the 85% accuracy shown in [12], because we omit capital gain and loss, and achieve similar accuracy as shown in the work [35].

The accuracy results are shown in Table 2. In order to achieve fairness, VFAE and FFVAE lose about 2% of their accuracy performance. RW, DIR and OP modify the dataset resulting in a loss of predictive performance. The proposed CF-VAE not only guarantees the fairness performance but also retains the causal relationships to improve accuracy. CF-VAE loses less information than other VAE-base methods and achieves the best fairness aware accuracy performance in all predictive models, i.e., 80.1%, 79.4% and 82.0% in $LR_C$, $SGD_C$ and $MLP_C$, respectively.

## 4.5 Ablation Study

We follow the same procedure in [7] to generate synthetic datasets and conduct an ablation study to validate the contribution of each component in our method as shown in Table 3.

The Full model uses all the observed attributes to train the predictors. The predictors achieve the best accuracy but the worst fairness

**Table 3: The results of ablation study. The Full model and VFAE are shown in the first two rows. The third row is the method without causal constraints. The fourth row is the method without employing OPR. Our proposed CF-VAE is shown in the last row. The best fairness aware RMSE and the best UFS$_R$ are shown in bold, and the runner-up results are underlined.**

| Loss function | Accuracy (RMSE) ↓ | | | Fairness (UFS$_R$) ↓ | | |
|---|---|---|---|---|---|---|
| | LR$_R$ | SGD$_R$ | MLP$_R$ | LR$_R$ | SGD$_R$ | MLP$_R$ |
| - | 0.078 ± 0.001 | 0.081 ± 0.001 | 0.081 ± 0.001 | 0.102 ± 0.001 | 0.098 ± 0.001 | 0.106 ± 0.002 |
| $-\mathcal{M}_{\text{VFAE}}$ | 0.126 ± 0.002 | 0.126 ± 0.002 | 0.145 ± 0.002 | 0.006 ± 0.001 | 0.010 ± 0.002 | 0.104 ± 0.005 |
| $-\mathcal{M}$ | 0.125 ± 0.001 | 0.125 ± 0.001 | 0.145 ± 0.001 | 0.007 ± 0.001 | 0.011 ± 0.003 | 0.105 ± 0.003 |
| $-\mathcal{M}' + \mathcal{L}_{\text{TCR}}$ | **0.109 ± 0.001** | <u>0.111 ± 0.001</u> | <u>0.122 ± 0.002</u> | <u>0.003 ± 0.001</u> | **0.004 ± 0.002** | <u>0.071 ± 0.002</u> |
| $-\mathcal{M}' + \mathcal{L}_{\text{TCR}} + \mathcal{L}_{\text{OPR}}$ | **0.109 ± 0.001** | **0.110 ± 0.001** | **0.121 ± 0.001** | **0.002 ± 0.001** | <u>0.005 ± 0.002</u> | **0.070 ± 0.002** |

performance as shown in the first row in Table 3. VFAE is the basic VAE-based unsupervised fair representation learning method. We set it to be the baseline in the second row in Table 3. The third row is CF-VAE without adding causal constraints, which achieves similar results as VFAE since both methods remove sensitive information from the learnt representations.

Then, we employ causal constraints and add TCR ($\gamma = 10$) in the loss function. As shown in the fourth row in Table 3, this step retains causal relationships in latent representations and improves both accuracy and fairness performance than previous rows. The last step is to encourage $Z'_x$ and $Z_a$ are disentangled by adding OPR. Our proposed CF-VAE achieves the best accuracy performance and UFS$_R$ among most predictive models as shown in the last row in Table 3.

## 5 RELATED WORKS

The machine learning literature has increasingly focused on exploring how algorithms can protect marginalised populations from unfair treatment. An important research area is how to quantify fairness, which can be divided into two categories, the statistical framework and the causal framework.

In the statistical framework, Demographic parity was defined by Zemel et al. [51], which is used to measure group level fairness. Other similar metrics include equalised odds [17], predictive rate parity [50]. Dwork et al. [13] proposed a measurement to quantify individual level fairness, that is, similar individuals should have similar treatments, and they use distance functions to measure how similar between individuals. In the causal framework, the (conditional) average causal effect is used to quantify fairness between groups [29]; Natural direct and natural indirect effects are used to quantify specific fairness [35, 52, 55]; When unfair causal paths are identified by domain knowledge, Chiappa [8] used the path-specific causal effects to quantify fairness on approved paths. For more related works, please refer to the literature review [9, 34, 53].

Our work is related to learning fair representations, which aims to encode data information into a lower space while removing sensitive information, and remaining causal relationships with respect to domain knowledge for building counterfactually fair predictive models. VAE [24] and $\beta$-VAE [18], as introduced in Section 2.2, have inspired several studies in fair representation learning. Louizos et al. [30] first introduced VAE for learning fair representation to disentangle the sensitive information and non-sensitive information,

they proposed a semi-supervised method to encourage disentanglement by using "Maximum Mean Discrepanc" (MMD). However, Zemel et al. [51], Gitiaux and Rangwala [15] argued that in real-world applications, the organisations that collect the data cannot predict the downstream uses of the data and the models that might be used. Due to this, there are many following up works focusing on unsupervised learning fair representation. For example, Creager et al. [11] proposed an algorithm that can achieve group level fairness by adding demographic parity as a constraint in objection function; Song et al. [42] developed an information theory-based method for learning maximally expressive representations subject to fairness constraints that allows users to control the fairness of representations by specifying limits on unfairness.

Our approach combines counterfactual fairness and unsupervised representation learning to provide the proper representations to help predictive models achieve counterfactual fairness. We extend the definition of counterfactual fairness [26] to the representation learning. Based on the current literature review, our work is the first method to use VAE-based techniques for unsupervised representation and satisfy counterfactual fairness. Furthermore, we innovatively embed domain knowledge into representations by adding causal constraints with respect to domain knowledge.

## 6 CONCLUSION

In this paper, we investigate unsupervised counterfactually fair representation learning and propose a novel method named CF-VAE which considers causal relationships with respect to domain knowledge. We theoretically demonstrate that the structured representations obtained by CF-VAE enable predictive models to achieve counterfactual fairness. Experimental results on real-world datasets show that CF-VAE achieves better accuracy and fairness performance on downstream predictive models than the benchmark fairness methods. Ablation study on synthetic datasets shows that causal constraints with total correction regularisation achieve better accuracy performance and orthogonality promoting regularisation encourages disentanglement with sensitive attributes.

# REFERENCES

[1] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilović, et al. 2019. AI Fairness 360: An Extensible Toolkit for Detecting and Mitigating Algorithmic Bias. *IBM Journal of Research and Development* 63, 4/5 (2019), 4–1.

[2] Marc Bendick. 2007. Situation Testing for Employment Discrimination in the United States of America. *Horizons stratégiques* 3 (2007), 17–39.

[3] Kenneth A. Bollen. 1989. *Structural Equations with Latent Variables*. Wiley.

[4] Tim Brennan, William Dieterich, and Beate Ehret. 2009. Evaluating the Predictive Validity of the COMPAS Risk and Needs Assessment System. *Criminal Justice and Behavior* 36, 1 (2009), 21–40.

[5] Flávio P. Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and Kush R. Varshney. 2017. Optimized Pre-Processing for Discrimination Prevention. In *NeurIPS 2017*. 3992–4001.

[6] Alycia N. Carey and Xintao Wu. 2022. The Causal Fairness Field Guide: Perspectives From Social and Formal Sciences. *Frontiers Big Data* 5 (2022), 892837.

[7] Debo Cheng, Jiuyong Li, Lin Liu, Kui Yu, Thuc Duy Le, and Jixue Liu. 2022. Toward Unique and Unbiased Causal Effect Estimation From Data With Hidden Variables. *IEEE Transactions on Neural Networks and Learning Systems* (2022), 1–13.

[8] Silvia Chiappa. 2019. Path-Specific Counterfactual Fairness. In *AAAI 2019*. 7801–7808.

[9] Sam Corbett-Davies and Sharad Goel. 2018. The Measure and Mismeasure of Fairness: A Critical Review of Fair Machine Learning. *CoRR* abs/1808.00023 (2018). arXiv:1808.00023

[10] Alexandru Coșer, Monica Mihaela Maer-matei, and Crişan Albu. 2019. Predictive Models for Loan Default Risk Assessment. *Economic Computation & Economic Cybernetics Studies & Research* 53, 2 (2019).

[11] Elliot Creager, David Madras, Jörn-Henrik Jacobsen, Marissa A. Weis, Kevin Swersky, Toniann Pitassi, and Richard S. Zemel. 2019. Flexibly Fair Representation Learning by Disentanglement. In *ICML 2019*. 1436–1445.

[12] Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository. http://archive.ics.uci.edu/ml

[13] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard S. Zemel. 2012. Fairness Through Awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference, ITCS 2012*. 214–226.

[14] Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and Removing Disparate Impact. In *SIGKDD 2015*. 259–268.

[15] Xavier Gitiaux and Huzefa Rangwala. 2021. Learning Smooth and Fair Representations. In *AISTATS 2021*. 253–261.

[16] Madelyn Glymour, Judea Pearl, and Nicholas P Jewell. 2016. *Causal Inference in Statistics: A Primer*. John Wiley & Sons.

[17] Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of Opportunity in Supervised Learning. In *NeurIPS 2016*. 3315–3323.

[18] Irina Higgins, Loïc Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. 2017. beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. In *ICLR 2017*. 1–22.

[19] Markus Kalisch and Peter Bühlmann. 2007. Estimating High-Dimensional Directed Acyclic Graphs with the PC-Algorithm. *Journal of Machine Learning Research* 8 (2007), 613–636.

[20] Faisal Kamiran and Toon Calders. 2012. Data Preprocessing Techniques for Classification without Discrimination. *Knowledge and Information Systems* 33, 1 (2012), 1–33.

[21] Hyunjik Kim and Andriy Mnih. 2018. Disentangling by Factorising. In *ICML 2018*. 2654–2663.

[22] Hyemi Kim, Seungjae Shin, JoonHo Jang, Kyungwoo Song, Weonyoung Joo, Wanmo Kang, and Il-Chul Moon. 2021. Counterfactual Fairness with Disentangled Causal Effect Variational Autoencoder. In *AAAI 2021*. 8128–8136.

[23] Suhong Kim, Param Joshi, Parminder Singh Kalsi, and Pooya Taheri. 2018. Crime Analysis Through Machine Learning. In *Proceedings of the 9th Annual Information Technology, Electronics and Mobile Communication Conference, IEMCON 2018*. 415–420.

[24] Diederik P. Kingma and Max Welling. 2014. Auto-Encoding Variational Bayes. In *ICLR 2014*. 1–14.

[25] Jochen Kruppa, Alexandra Schwarz, Gerhard Arminger, and Andreas Ziegler. 2013. Consumer Credit Risk: Individual Probability Estimates Using Machine Learning. *Expert Systems with Applications* 40, 13 (2013), 5125–5131.

[26] Matt J. Kusner, Joshua R. Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual Fairness. In *NeurIPS 2017*. 4066–4076.

[27] Nicol Turner Lee, Paul Resnick, and Genie Barton. 2019. Algorithmic Bias Detection and Mitigation: Best Practices and Policies to Reduce Consumer Harms. *Brookings Institute: Washington, DC, USA* (2019).

[28] David Lewis. 2013. *Counterfactuals*. John Wiley & Sons.

[29] Jiuyong Li, Jixue Liu, Lin Liu, Thuc Duy Le, Saisai Ma, and Yizhao Han. 2017. Discrimination Detection by Causal Effect Estimation. In *IEEE BigData 2017*. 1087–1094.

[30] Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard S. Zemel. 2016. The Variational Fair Autoencoder. In *ICLR 2016*. 1–11.

[31] Ryan Mac. 2021. *Facebook apologizes after AI Puts 'primates' label on video of black men*. Retrieved March, 2022 from https://www.nytimes.com/2021/09/03/technology/facebook-ai-race-primates.html

[32] David Madras, Elliot Creager, Toniann Pitassi, and Richard S. Zemel. 2018. Learning Adversarially Fair and Transferable Representations. In *ICML 2018*. 3381–3390.

[33] David Madras, Elliot Creager, Toniann Pitassi, and Richard S. Zemel. 2019. Fairness through Causal Awareness: Learning Causal Latent-Variable Models for Biased Data. In *Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT* 2019*. 349–358.

[34] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A Survey on Bias and Fairness in Machine Learning. *Comput. Surveys* 54, 6 (2021), 115:1–115:35.

[35] Razieh Nabi and Ilya Shpitser. 2018. Fair Inference on Outcomes. In *AAAI 2018*. 1931–1940.

[36] Sungho Park, Sunhee Hwang, Dohyung Kim, and Hyeran Byun. 2021. Learning Disentangled Representation for Fair Facial Attribute Classification via Fairness-aware Information Alignment. In *AAAI 2021*. 2403–2411.

[37] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *NeurIPS 2019*. 8024–8035.

[38] Judea Pearl. 2009. Causal Inference in Statistics: An Overview. *Statistics Surveys* 3 (2009), 96–146.

[39] Judea Pearl. 2009. *Causality*. Cambridge University Press.

[40] Thomas Richardson and Peter Spirtes. 2002. Ancestral Graph Markov Models. *The Annals of Statistics* 30, 4 (2002), 962–1030.

[41] Mhd Hasan Sarhan, Nassir Navab, Abouzar Eslami, and Shadi Albarqouni. 2020. Fairness by Learning Orthogonal Disentangled Representations. In *ECCV 2020*. 746–761.

[42] Jiaming Song, Pratyusha Kalluri, Aditya Grover, Shengjia Zhao, and Stefano Ermon. 2019. Learning Controllable Fair Representations. In *AISTATS 2019*. 2164–2173.

[43] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. 2000. *Causation, Prediction, and Search*. MIT press.

[44] Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. 2019. Balanced Datasets Are Not Enough: Estimating and Mitigating Gender Bias in Deep Image Representations. In *ICCV 2019*. IEEE, 5309–5318.

[45] Michael Satosi Watanabe. 1960. Information Theoretical Analysis of Multivariate Correlation. *IBM Journal of Research and Development* 4, 1 (1960), 66–82.

[46] Linda F Wightman. 1998. LSAC National Longitudinal Bar Passage Study. LSAC Research Report Series. (1998).

[47] Pengtao Xie, Wei Wu, Yichen Zhu, and Eric P. Xing. 2018. Orthogonality-Promoting Distance Metric Learning: Convex Relaxation and Theoretical Analysis. In *ICML 2018*. 5399–5408.

[48] Mengyue Yang, Furui Liu, Zhitang Chen, Xinwei Shen, Jianye Hao, and Jun Wang. 2021. CausalVAE: Disentangled Representation Learning via Neural Structural Causal Models. In *CVPR 2021*. 9593–9602.

[49] Yang Yu, Yu-Feng Li, and Zhi-Hua Zhou. 2011. Diversity Regularized Machine. In *IJCAI 2011*. 1603–1608.

[50] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez-Rodriguez, and Krishna P. Gummadi. 2017. Fairness Beyond Disparate Treatment & Disparate Impact: Learning Classification without Disparate Mistreatment. In *WWW 2017*. 1171–1180.

[51] Richard S. Zemel, Yu Wu, Kevin Swersky, Toniann Pitassi, and Cynthia Dwork. 2013. Learning Fair Representations. In *ICML 2013*. 325–333.

[52] Junzhe Zhang and Elias Bareinboim. 2018. Fairness in Decision-Making - The Causal Explanation Formula. In *AAAI 2018*. 2037–2045.

[53] Lu Zhang and Xintao Wu. 2017. Anti-discrimination Learning: A Causal Modeling-based Framework. *International Journal of Data Science and Analytics* 4, 1 (2017), 1–16.

[54] Lu Zhang, Yongkai Wu, and Xintao Wu. 2017. Achieving Non-Discrimination in Data Release. In *SIGKDD 2017*. 1335–1344.

[55] Lu Zhang, Yongkai Wu, and Xintao Wu. 2017. A Causal Framework for Discovering and Removing Direct and Indirect Discrimination. In *IJCAI 2017*. 3929–3935.

[56] Maggie Zhang. 2015. *Google photos tags two african-americans as gorillas through facial recognition software*. Retrieved March, 2022 from https://www.forbes.com/sites/mzhang/2015/07/01/google-photos-tags-two-african-americans-as-gorillas-through-facial-recognition-software.html

[57] Xun Zheng, Bryon Aragam, Pradeep Ravikumar, and Eric P. Xing. 2018. DAGs with NO TEARS: Continuous Optimization for Structure Learning. In *NeurIPS 2018*. 9492–9503.