# Hierarchical VampPrior Variational Fair Auto-Encoder

## Philip Botros 1 Jakub M. Tomczak 1

## **Abstract**

Decision making is a process that is extremely prone to different biases. In this paper we consider learning fair representations that aim at removing nuisance (sensitive) information from the decision process. For this purpose, we propose to use deep generative modeling and adapt a hierarchical Variational Auto-Encoder to learn these fair representations. Moreover, we utilize the mutual information as a useful regularizer for enforcing fairness of a representation. In experiments on two benchmark datasets and two scenarios where the sensitive variables are fully and partially observable, we show that the proposed approach either outperforms or performs on par with the current best model.

## 1. Introduction

Reducing bias in machine learning algorithms has been an active area of discussion recently after the reliance on algorithmic decision making has been greatly increased. Consider the case of credit assignment, mortgage approvals or the provision of health care, where there are growing concerns that biases based on historical data prevent a fair process.

In these cases it is not sufficient to prevent the decision maker from having access to the sensitive variable since this information has already been leaked into other features (Dwork et al., 2012; McNamara et al., 2017; Menon & Williamson, 2017; Zafar et al., 2017). To correct for this and ensure a fair process, a new representation has to be created where all the sensitive information is removed.

More formally, we can consider this problem as the task of learning fair representations, where the goal is to learn a representation **z** that maximizes the information about the class label **y**, while removing the sensitive information **s**.

Presented at the ICML 2018 workshop on Theoretical Foundations and Applications of Deep Generative Models, Stockholm, Sweden, 2018. Copyright 2018 by the author(s).

Learning rich representations from vast amounts of data using deep generative models remains one of the major challenges of machine learning. In recent years, different approaches to achieving this goal were proposed by formulating alternative training objectives to the log-likelihood (Goodfellow et al., 2014) or by utilizing variational inference that leads to a highly scalable framework now known as the variational auto-encoders (VAE) (Kingma & Welling, 2013; Rezende et al., 2014).

The use of a deep generative model for fair classification has already been explored by (Louizos et al., 2015) who proposed the Variational Fair Auto-Encoder (VFAE). They, however, do not consider the partially-supervised case with partially observed s which is more applicable in real-world settings nor do they address the problem of inactive stochastic units inherent in deep latent variable models.

Furthermore, even though the formulation of the graphical model encourages separation between the sensitive variable and the latent representation, some sensitive information can remain if this information is correlated with the prediction task. Therefore, an additional regularization term is necessary to further enforce fairness of the representation. We follow this line of thinking and explore the mutual information as a fairness regularizer to ensure that the sensitive information is removed completely.

The contribution of the paper is twofold:

- We propose a new deep generative model for learning fair representations and empirically show that it outperforms the VFAE when s is observed partially and performs on par when the sensitive variable is fully observed.
- We introduce the mutual information as a new fairness regularizer that is better suited for the realistic case where the sensitive variables are partially observed.

## 2. Fair Deep Generative Models

Learning fair representations aims at becoming invariant to nuisance or sensitive factors while retaining as much of the remaining relevant information as possible (Louizos et al., 2015; Zemel et al., 2013). From the probabilistic modeling point of view, the problem could be formulated in terms of a

<sup>&</sup>lt;sup>1</sup>AMLAB, University of Amsterdam. Correspondence to: Philip Botros <philip.botros@student.uva.nl>, Jakub M. Tomczak <j.m.tomczak@uva.nl>.

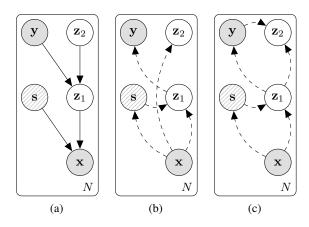


Figure 1. (a) Generative part of both models. (b) Variational part H-VFAE. (c) Variational part VFAE.

set of independent factors working on the input  $\mathbf{x} \in \mathcal{X}$ ,  $\mathcal{X}$  is a D-dimensional discrete or continuous space, namely, the (partially)-observed discrete sensitive (nuisance) variable  $\mathbf{s} \in \mathcal{S}$ , typically  $\mathcal{S} = \{0,1\}$ , and the continuous unobserved latent variable  $\mathbf{z}_1 \in \mathbb{R}^{M_1}$ . Additionally, since the goal is to learn features that are invariant to  $\mathbf{s}$  without losing information about the label  $\mathbf{y}$ , a hierarchy of latent variables could be introduced. In this paper, we assume a second layer of latent variables  $\mathbf{z}_2 \in \mathbb{R}^{M_2}$ . All the label independent noise inherent in  $\mathbf{x}$  is modeled in the hierarchical latent representation, while also allowing the model to correlate the discrete label information  $\mathbf{y}$  with the invariant features  $\mathbf{z}_1$ . As a result, the following generative process could be considered:

$$\mathbf{y} \sim \text{Cat}(\mathbf{y})$$
 (1)

$$\mathbf{z}_2 \sim p(\mathbf{z}_2) \tag{2}$$

$$\mathbf{z}_1 \sim p_{\theta}(\mathbf{z}_1|\mathbf{z}_2,\mathbf{y})$$
 (3)

$$\mathbf{x} \sim p_{\theta}(\mathbf{x}|\mathbf{z}_1, \mathbf{s}),$$
 (4)

where  $Cat(\cdot)$  denotes the categorical distribution, see Figure 1 for the probabilistic graphical model. This formulation can be recast as an inference problem where the objective is to learn the posterior  $p(\mathbf{z_1}, \mathbf{z_2}, \mathbf{y} | \mathbf{x}, \mathbf{s})$  in the case of observed  $\mathbf{s}$ , and  $p(\mathbf{z_1}, \mathbf{z_2}, \mathbf{y}, \mathbf{s} | \mathbf{x})$  in the case of unobserved  $\mathbf{s}$ . Since calculating the true posterior is infeasible, we will use variational inference and the methodology of variational auto-encoders (VAE) (Kingma & Welling, 2013).

## 3. Variational Fair Auto-Encoder

Depending on assumed dependencies among random variables in the variational posterior, the application of variational inference to the generative process presented in the previous section may result in different VAE architectures. Louizos et al. (2015) assumed s is always given and, thus, they proposed to factorize the variational posterior

 $q_{\phi}(\mathbf{z}_1, \mathbf{z}_2, \mathbf{y}|\mathbf{x}, \mathbf{s})$  as  $q_{\phi}(\mathbf{z}_1|\mathbf{x}, \mathbf{s})q_{\phi}(\mathbf{y}|\mathbf{z}_1)q_{\phi}(\mathbf{z}_2|\mathbf{z}_1, \mathbf{y})$ , see Figure 1(c). The final model is defined as follows:

$$q_{\phi}(\mathbf{z}_{1}|\mathbf{x},\mathbf{s}) = \mathcal{N}(\mathbf{z}_{1}|\boldsymbol{\mu}_{\phi}(\mathbf{x},\mathbf{s}),\boldsymbol{\sigma}_{\phi}(\mathbf{x},\mathbf{s}))$$
 (5)

$$q_{\phi}(\mathbf{y}|\mathbf{z}_1) = \operatorname{Cat}(\mathbf{y}|\boldsymbol{\pi}_{\phi}(\mathbf{z}_1)) \tag{6}$$

$$q_{\phi}(\mathbf{z}_{2}|\mathbf{z}_{1},\mathbf{y}) = \mathcal{N}(\mathbf{z}_{2}|\boldsymbol{\mu}_{\phi}(\mathbf{z}_{1},\mathbf{y}),\boldsymbol{\sigma}_{\phi}(\mathbf{z}_{1},\mathbf{y}))$$
(7)

$$p(\mathbf{z}_2) = \mathcal{N}(\mathbf{z}_2|\mathbf{0}, \mathbf{I}) \tag{8}$$

$$p_{\theta}(\mathbf{z}_1|\mathbf{z}_2, \mathbf{y}) = \mathcal{N}(\mathbf{z}_1|\boldsymbol{\mu}_{\phi}(\mathbf{z}_2, \mathbf{y}), \boldsymbol{\sigma}_{\phi}(\mathbf{z}_2, \mathbf{y}))$$
(9)

$$p_{\theta}(\mathbf{x}|\mathbf{z}_1, \mathbf{s}) = f_{\theta}(\mathbf{z}_1, \mathbf{s}), \tag{10}$$

where all distributions are parameterized by neural networks, and  $f_{\theta}(\mathbf{z}_1, \mathbf{s})$  is a distribution suited for the data that is modeled. Following this formulation, we aim at maximizing the variational (evidence) lower bound on  $\ln p(\mathbf{x}|\mathbf{s})$  (ELBO):

$$\mathcal{L}_s(\mathbf{x}, y, \mathbf{s}) = \mathbb{E}_{q_{\theta}(\mathbf{z}_1 | \mathbf{x}, \mathbf{s})} [\ln p_{\theta}(\mathbf{x} | \mathbf{z}_1, \mathbf{s}) -$$
(11)

$$KL(q_{\phi}(\mathbf{z}_{2}|\mathbf{z}_{1},\mathbf{y})||p(\mathbf{z}_{2}))]+\tag{12}$$

$$\mathbb{E}_{q_{\phi}(\mathbf{z}_{1},\mathbf{z}_{2}|\mathbf{x},\mathbf{s},\mathbf{y})}[\ln p_{\theta}(\mathbf{z}_{1}|\mathbf{z}_{2},\mathbf{y}) - \ln q_{\phi}(\mathbf{z}_{1}|\mathbf{x},\mathbf{s})] + (13)$$

$$\alpha \mathbb{E}_{q_{\phi}(\mathbf{z}_{1}|\mathbf{x},\mathbf{s})}[\ln q_{\phi}(\mathbf{y}|\mathbf{z}_{1})] \tag{14}$$

where  $\alpha>0$  is an additional parameter to control the influence of the classifier during training. The ELBO can be jointly optimized with respect to the parameters  $\phi,\theta$  of the inference and generative model, respectively, using the reparameterization trick (Kingma & Welling, 2013). We refer to this model as the Variational Fair Auto-Encoder (VFAE).

## 4. Hierarchical VampPrior VFAE

The VFAE is shown to be successful in learning fair representations (Louizos et al., 2015). However, it has been shown that in general deep VAEs suffer from the inactive latent variable problem (Sønderby et al., 2016), following from a top-down multi-layered generative process while the variational part is bottom-up. Therefore, we propose to change the variational part of the VFAE with inputs fed directly to the deepest layer such that the final encoder changes to  $q_{\phi}(\mathbf{z}_2|\mathbf{x})$ . This has the effect of enforcing a dependency between the data and the latent units at the deepest level during the generative process, preventing the latent units from regularization towards the prior, i.e. setting  $q_{\phi}(\mathbf{z}|\mathbf{x}) = p(\mathbf{z})$ .

Furthermore, this formulation allows for easy integration of a recently proposed powerful prior, the Variational Mixture of Posteriors Prior (VampPrior) with a hierarchical architecture (Tomczak & Welling, 2017). We hypothesize that the quality of the VFAE could be improved by utilizing a different family of variational posteriors coupled with a powerful prior over the latent representation utilizing the new hierarchical structure.

Eventually, we end up with a different structure of the variational posterior (see Figure 1(b) for the probabilistic graphical model):

$$q_{\phi}(\mathbf{z}_1, \mathbf{z}_2, \mathbf{y}|\mathbf{x}, \mathbf{s}) = q_{\phi}(\mathbf{z}_1|\mathbf{x}, \mathbf{s})q_{\phi}(\mathbf{y}|\mathbf{z}_1)q_{\phi}(\mathbf{z}_2|\mathbf{x}).$$
 (15)

Now, we can consider the problem of finding the prior that optimizes the lower bound given the data. The solution is simply the aggregated posterior (Hoffman & Johnson, 2016):

$$p_{\lambda}^{*}(\mathbf{z}) = \frac{1}{N} \sum_{n=1}^{N} q_{\phi}(\mathbf{z}|\mathbf{x}_{n}). \tag{16}$$

This could, however, lead to overfitting and would be very expensive to compute for every training iteration. A computationally efficient alternative, which also prevents from overfitting by restricting  $K \ll N$ , is an approximation using a mixture of variational posteriors with learnable pseudoinputs (Tomczak & Welling, 2017):

$$p_{\lambda}(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^{K} q_{\phi}(\mathbf{z}|\mathbf{u}_{k}), \tag{17}$$

where K is the number of pseudo-inputs and  $\mathbf{u}_k$  denotes the k-th pseudo-input that is of the same dimension as the input.

The final model is defined as follows:

$$q_{\phi}(\mathbf{z}_1|\mathbf{x},\mathbf{s}) = \mathcal{N}(\mathbf{z}_1|\boldsymbol{\mu}_{\phi}(\mathbf{x},\mathbf{s}),\boldsymbol{\sigma}_{\phi}(\mathbf{x},\mathbf{s}))$$
(18)

$$q_{\phi}(\mathbf{y}|\mathbf{z}_1) = \operatorname{Cat}(\mathbf{y}|\boldsymbol{\pi}_{\phi}(\mathbf{z}_1)) \tag{19}$$

$$q_{\phi}(\mathbf{z}_2|\mathbf{x}) = \mathcal{N}(\mathbf{z}_2|\boldsymbol{\mu}_{\phi}(\mathbf{x}), \boldsymbol{\sigma}_{\phi}(\mathbf{x}))$$
 (20)

$$p_{\lambda}(\mathbf{z}_2) = \frac{1}{K} \sum_{k=1}^{K} q_{\phi}(\mathbf{z}_2 | \mathbf{u}_k)$$
 (21)

$$p_{\theta}(\mathbf{z}_1|\mathbf{z}_2, \mathbf{y}) = \mathcal{N}(\mathbf{z}_1|\boldsymbol{\mu}_{\phi}(\mathbf{z}_2, \mathbf{y}), \boldsymbol{\sigma}_{\phi}(\mathbf{z}_2, \mathbf{y}))$$
(22)

$$p_{\theta}(\mathbf{x}|\mathbf{z}_1, \mathbf{s}) = f_{\theta}(\mathbf{z}_1, \mathbf{s}) \tag{23}$$

The objective function is the ELBO in the following form:

$$\mathcal{L}_s(\mathbf{x}, y, \mathbf{s}) = \mathbb{E}_{q_{\theta}(\mathbf{z}_1 | \mathbf{x}, \mathbf{s})} [\ln p_{\theta}(\mathbf{x} | \mathbf{z}_1, \mathbf{s})] - \tag{24}$$

$$KL(q_{\phi}(\mathbf{z}_{2}|\mathbf{x})||p_{\lambda}(\mathbf{z}_{2})) - \tag{25}$$

$$\mathbb{E}_{q_{\phi}(\mathbf{z}_{2}|\mathbf{x})}[\text{KL}(q_{\phi}(\mathbf{z}_{1}|\mathbf{x},\mathbf{s}) || p_{\theta}(\mathbf{z}_{1}|\mathbf{z}_{2},\mathbf{y}))] + (26)$$

$$\alpha \mathbb{E}_{q_{\phi}(\mathbf{z}_1|\mathbf{x},\mathbf{s})} [\ln q_{\phi}(\mathbf{y}|\mathbf{z}_1)]. \tag{27}$$

We refer to this model as Hierarchical VampPrior Variational Fair Auto-Encoder (H-VFAE + VP).

Alternatively, we can consider a simpler case where the VampPrior is replaced by the standard Gaussian prior,  $p_{\lambda}(\mathbf{z}_2) = \mathcal{N}(\mathbf{z}_2|\mathbf{0},\mathbf{I})$ . We will call this model the Hierarchical Variational Fair Auto-Encoder (H-VFAE).

## 5. Encouraging learning fair representations

If the sensitive variable is correlated with the prediction task, information about s can still remain in the latent representation  $z_1$ . To remove this information, two fairness penalties

are discussed, which can easily be added to the lower bound as a regularizer.

### 5.1. MMD regularizer

Louizos et al. (2015) originally proposed to use the Maximum Mean Discrepancy (MMD) measure to regularize the marginal  $q_{\phi}(\mathbf{z}_1|\mathbf{s})$ . The rationale behind applying the MMD is that it compares statistics of two samples, and if they are similar, the MMD indicates that they were drawn from the same distribution. The distance between the empirical statistics  $\varphi$  of two datasets can be computed in the following manner:

$$\|\frac{1}{N_0} \sum_{i=1}^{N_0} \varphi(\mathbf{z}_0) - \frac{1}{N_1} \sum_{i=1}^{N_1} \varphi(\mathbf{z}_1)\|^2.$$
 (28)

An unbiased MMD estimator (Gretton et al., 2007) is obtained by expanding the square and is solely composed of inner products on which the kernel trick can be applied:

$$\ell_{MMD} = \mathbb{E}_{p(\mathbf{z}_0), p(\mathbf{z}_0')}[k(\mathbf{z}_0, \mathbf{z}_0')] + \tag{29}$$

$$\mathbb{E}_{q(\mathbf{z}_1),q(\mathbf{z}_1')}[k(\mathbf{z}_1,\mathbf{z}_1')] - 2\mathbb{E}_{p(\mathbf{z}_0),q(\mathbf{z}_1)}[k(\mathbf{z}_0,\mathbf{z}_1)]. \quad (30)$$

Optimizing for the MMD regularizer has the effect of matching the moments of marginal distributions  $q_{\phi}(\mathbf{z}_1|\mathbf{s}=0)$  and  $q_{\phi}(\mathbf{z}_1|\mathbf{s}=1)$ , while still allowing individual elements to differ. In our case, the MMD regularizer is the following:

$$\ell_{MMD} = \mathbb{E}_{\tilde{p}(\mathbf{x})} \Big[ \| \mathbb{E}_{q_{\phi}(\mathbf{z}_1 | \mathbf{x}, \mathbf{s} = 0)} [\varphi(\mathbf{z}_1)] -$$
(31)

$$\mathbb{E}_{q_{\phi}(\mathbf{z}_1|\mathbf{x},\mathbf{s}=1)}[\varphi(\mathbf{z}_1)]\|^2, \qquad (32)$$

where  $\tilde{p}(\mathbf{x})$  denotes the empirical distribution.

The behavior of the MMD regularizer is schematically presented in Figure 2. Two marginal distributions are matched by matching their respective moments.

### 5.1.1. FAST MMD REGULARIZER

To prevent computing the expensive full MMD estimator, random kitchen sinks (Rahimi & Recht, 2009) can be used to compute the feature expansion  $\varphi(\mathbf{z})$  to serve as an approximation to the MMD regularizer. The idea is to draw a random matrix  $\mathbf{W} \in \mathbb{R}^{M \times K}$ , with M as the dimensionality of  $\mathbf{z}$  and K as the number of random features, where each entry is drawn from a standard isotropic Gaussian. Additionally, a M-dimensional uniform random vector  $\mathbf{b}$  is drawn with entries in  $[0, 2\pi]$ . The feature expansion can be then computed as follows (Louizos et al., 2015):

$$\varphi_{\mathbf{W}}(\mathbf{z}) = \sqrt{\frac{2}{D}\cos\left(\sqrt{\frac{2}{\gamma}}\mathbf{z}\mathbf{W} + \mathbf{b}\right)},$$
 (33)

where  $\gamma=2M$ . The inner product of these randomized feature expansions converges to a kernel function given an increasing number of features.

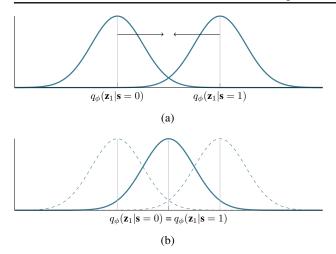


Figure 2. A schematic presentation of the behavior of the MMD regularizer.

## 5.2. Mutual information regularizer

Another manner of enhancing fairness is to force independence between the representation and the sensitive variable. A natural candidate for this purpose is the mutual information, which represents a measure of mutual dependence between two random variables. In our case we are interested in the conditional mutual information between  $\mathbf{z}_1$  and  $\mathbf{s}$  for given  $\mathbf{x}$ , that is:

$$MI(\mathbf{z}_{1}, \mathbf{s} | \mathbf{x}) = \mathbb{E}_{\tilde{p}(\mathbf{x})q_{\phi}(\mathbf{z}_{1}, \mathbf{s} | \mathbf{x})} \left[ \ln \frac{q_{\phi}(\mathbf{z}_{1}, \mathbf{s} | \mathbf{x})}{q_{\phi}(\mathbf{z}_{1} | \mathbf{x})q_{\phi}(\mathbf{s} | \mathbf{x})} \right].$$
(34)

Furthermore assuming that the posterior factorizes as  $q_{\phi}(\mathbf{s}|\mathbf{x})q_{\phi}(\mathbf{z}_1|\mathbf{x},\mathbf{s})$  we obtain an easily computable and differentiable estimator:

$$\ell_{MI} = \mathbb{E}_{\tilde{p}(\mathbf{x})q_{\phi}(\mathbf{z}_{1},\mathbf{s}|\mathbf{x})} \left[ \ln \frac{q_{\phi}(\mathbf{z}_{1}|\mathbf{x},\mathbf{s})}{q_{\phi}(\mathbf{z}_{1}|\mathbf{x})} \right]$$
(35)

$$= \mathbb{E}_{\tilde{p}(\mathbf{x})q_{\phi}(\mathbf{z}_{1},\mathbf{s}|\mathbf{x})} \left[ \ln \frac{q_{\phi}(\mathbf{z}_{1}|\mathbf{x},\mathbf{s})}{\sum_{s} q_{\phi}(\mathbf{s}|\mathbf{x})q_{\phi}(\mathbf{z}_{1}|\mathbf{x},\mathbf{s})} \right], \quad (36)$$

which can be approximated using Monte Carlo samples. Typically,  ${\bf s}$  is low-dimensional (e.g., it is binary), hence, calculating the mixture distribution in the denominator is easily tractable. Stochastic gradient ascent can now be performed on  $\nabla_{\phi}\ell_{MI}$  to regularize the encoder  $q_{\phi}({\bf z}_1|{\bf x},{\bf s})$ . The behavior of the Mutual Information regularizer is schematically presented in Figure 3. Notice that in contrast to the MMD regularizer, here we try to match each mode separately to the mixture and the optimal solution is attained when both modes overlap. In other words, the encoder does not use the information about the sensitive variable solely if it produces the same distribution for any value of  ${\bf s}$ .

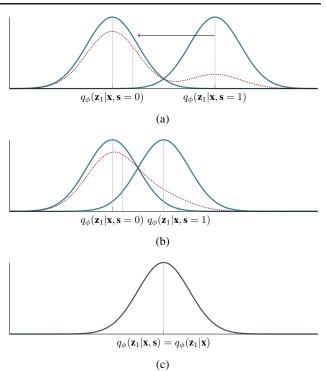


Figure 3. A schematic representation of the behavior of the MI regularizer.

## 6. Partial supervision of the sensitive variable

Typically, it is assumed that the sensitive variable is fully observable. However, in many real-life applications s is only partially-observable. In this case, the generative approach allows the model to infer the sensitive variable so that all data could be used during training.

The supervised model can be easily extended to cope with these examples where no sensitive variables are provided by adding a variational categorical distribution  $q_{\phi}(\mathbf{s}|\mathbf{x})$  to the model. Besides  $\mathcal{L}_s$  we now jointly optimize the unsupervised lower bound where we can sum over all possible values of  $\mathbf{s}$  or we can use differentiable samples from the random discrete node  $q_{\phi}(\mathbf{s}|\mathbf{x})$  obtained by the reparameterization trick using the concrete distribution (Jang et al., 2016; Maddison et al., 2016):

$$\mathcal{L}_{u}(\mathbf{x}, y) = \mathbb{E}_{q_{\phi}(\mathbf{s}|\mathbf{x})} \left[ \mathcal{L}_{s}(\mathbf{x}, y, \mathbf{s}) \right] - KL(q_{\phi}(\mathbf{s}|\mathbf{x}) || p(\mathbf{s})), \tag{37}$$

where the following distributions are introduced:

$$q_{\phi}(\mathbf{s}|\mathbf{x}) = \text{Cat}(\mathbf{s}|\boldsymbol{\pi}_{\phi}(\mathbf{x})) \tag{38}$$

$$p(\mathbf{s}) = \operatorname{Cat}(\mathbf{s}|\pi),\tag{39}$$

and  $\pi$  denote a priori probabilities of sensitive variables.

Eventually, we can combine both the supervised and unsupervised objectives together that gives our final objective function for given training data  $\mathcal{D}$ :

$$\mathcal{L}(\mathcal{D}) = \mathbb{E}_{\tilde{p}(\mathbf{x}, y, \mathbf{s})} [\mathcal{L}_{s}(\mathbf{x}, y, \mathbf{s})] + \mathbb{E}_{\tilde{p}(\mathbf{x}, \mathbf{s})} [\mathcal{L}_{u}(\mathbf{x}, y)] + (40)$$

$$\mathbb{E}_{\tilde{p}(\mathbf{x}, \mathbf{s})} [-\ln q_{\phi}(\mathbf{s}|\mathbf{x})] - \mathbb{E}_{\tilde{p}(\mathbf{x}, y)} [\lambda \ell(\mathbf{z}_{1})], \quad (41)$$

where  $\tilde{p}(\cdot)$  denotes the empirical distribution,  $\lambda > 0$ , and  $\ell(\mathbf{z}_1)$  is either the MMD regularizer or the MI regularizer.

## 7. Experiments

## 7.1. Datasets

Experiments were run on the German and Adult datasets with the same training, validation and test splits as used in (Zemel et al., 2013). The German dataset consists of credit data and the objective is to predict if a person has a good or bad credit rating. The sensitive variable here is the age of the individual. The Adult income dataset consists of census data and the prediction task is to determine whether a person makes over 50.000 dollars a year. The sensitive variable here is the gender. We binarized both datasets and used a Bernoulli distribution for the final decoder, i.e.  $p_{\theta}(\mathbf{x}|\mathbf{z}_1,\mathbf{s}) = \text{Bern}(\mathbf{x}|\pi_{\phi}(\mathbf{z}_1,\mathbf{s}))$  similarly to (Louizos et al., 2015).

### 7.2. Settings

The same neural network architectures as in (Louizos et al., 2015) were used for all experiments. For the small German dataset, a hidden layer of 60 units was used for all encoders and decoders, with a stochastic latent dimensionality of 30 units. For the Adult income dataset, 100 hidden units were used for all encoders and decoders, while the dimensionality of the latent space was increased to 50.

Like (Louizos et al., 2015) we took  $\alpha=1$  for the supervised setting, while  $\beta$  was cross-validated due to the varying nature of the strength of the regularizers. For the partially-supervised setting, we set  $\alpha=20$  as it was observed that the partially-supervised regularized models were well suited to handle the increased dependency on the classification error. Optimization was done with the Adam optimizer (Kingma & Ba, 2014), where the default settings were used.

#### 7.3. Evaluation

The primary goal of the paper is fair classification, therefore, all models need to be evaluated with respect to the classification accuracy and with respect to information being available about the sensitive variables.

To measure the information about the sensitive variables remaining in the predictive features, a logistic regression classifier was trained to predict the state of our sensitive variable given the features from the variational posterior  $q_{\phi}(\mathbf{z}_1|\mathbf{x},\mathbf{s})$ .

Furthermore, since another objective of fair classification is group fairness, ensuring equal treatment between different groups, the probabilistic discriminative metric from (Louizos et al., 2015) is used:

$$DS = \left| \frac{\sum_{n|s_n=0} p(\hat{y}_n)}{N_{s=0}} - \frac{\sum_{n|s_n=1} p(\hat{y}_n)}{N_{s=1}} \right|, \quad (42)$$

This metric has the simple interpretation of measuring the difference in classification predictions between different groups.

Both models were tested on the full supervision and partial supervision of s, where the fraction of observed sensitive variables was set to 0.05. To test the regularization penalties we evaluate the models with  $\ell_{MI}$ ,  $\ell_{MMD}$  and no regularizer

Additionally, the influence of the VampPrior was isolated by also providing a baseline of our proposed model with the standard Gaussian prior instead of a richer one, denoted by H-VFAE in our experiments.

### 7.4. Experiment with the full supervision of s

In the first experiment, our model and the newly proposed regularizer are evaluated on the case with fully observed sensitive variables.<sup>1</sup> The results are presented in Table 1.

First of all, we notice that without any regularization the proposed family of variational posteriors performs similarly to the VFAE in terms of the classification accuracy on y, however, it performs worst on the DS metric. It is also worth to note that our model greatly benefits from the VampPrior.

The benefits of the VampPrior on our new architecture are easily observed by noting that the prediction accuracy on average increases on both datasets while still having a regularized effect, resulting in a lower amount of information available about s.

Moreover, the effect of applying the VampPrior is presented in Figure 4 (the crosses represent means of the components  $q_{\phi}(\mathbf{z}_1|\mathbf{u}_k)$ ). Notice that the prior is highly multi-modal and it covers the latent space in places where the encoder places  $\mathbf{z}_1$  for given  $\mathbf{x}$  and  $\mathbf{s}$ . If the standard Gaussian prior is used, the encoder would be forced to put most of the points close to the origin. Also note that the latent representations are almost indistinguishable for the two groups.

Additionally, the new model with the VampPrior outperforms the original VFAE on predictive capabilities on both datasets while keeping the information about s similar. Furthermore, all models are as good as invariant against classification with respect to s on the features  $\mathbf{z}_1$ .

<sup>&</sup>lt;sup>1</sup>In order to have comparable results to the original paper on VFAE, we used the same experiment setting as in (Louizos et al., 2015).

<i>Table 1.</i> Results on the fully	supervised case.	Best results of methods w	ith fairness regularization in bold.

Model	GERMAN Y	ADULT Y	GERMAN S	ADULT S	GERMAN DS	ADULT DS
RANDOM	71.1	75.0	80.1	67.0	-	_
VFAE	72.4	82.0	80.1	66.1	2.7	5.4
H-VFAE	72.5	80.5	80.1	68.3	10.1	11.3
H-VFAE + VP	72.4	81.9	80.1	67.2	5.5	7.9
VFAE + MMD	72.7	81.3	80.1	67.4	0.6	2.5
H-VFAE + MMD	74.2	81.0	80.1	67.2	7.2	8.6
H-VFAE + VP + MMD	74.4	82.2	80.1	67.2	1.6	3.3
VFAE + MI	72.9	81.6	80.1	67.4	4.2	3.4
H-VFAE + MI	73.6	82.1	80.1	67.4	5.1	5.6
H-VFAE + VP + MI	73.4	82.1	80.1	67.3	3.7	2.5

Table 2. Results on the partially supervised case. Best results of methods with fairness regularization in bold.

MODEL	GERMAN Y	ADULT Y	GERMAN S	ADULT S	GERMAN DS	ADULT DS
RANDOM	71.1	75.0	80.1	67.0	-	-
VFAE	75.5	84.8	80.1	69.7	8.6	11.4
H-VFAE + VP	75.1	84.5	80.1	69.4	8.8	10.7
VFAE + MMD	73.4	81.5	80.1	67.4	3.4	8.1
H-VFAE + VP + MMD	73.4	81.7	80.1	67.4	3.2	6.1
VFAE + MI	72.8	82.0	80.1	67.4	3.3	5.6
H-VFAE + VP + MI	74.1	82.3	80.1	67.4	3.1	4.9

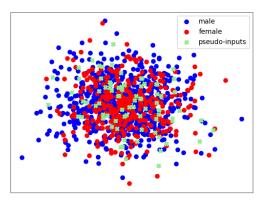


Figure 4. The 2D latent space visualization for  $\mathbf{z}_1$  in the Adult dataset. The colors correspond to gender values and the crosses represent means for components of the VampPrior.

If we look at the regularization penalties, it is clear that both the MMD and MI penalties have a significant regularizing effect on the information retained in s. The MMD regularizer is shown to be a marginally better fit for the supervised case with slightly lower scores on s while retaining the same classification accuracy. A possible explanation for that is that for sufficiently large training samples, the MMD

is better approximated than the MI regularizer. Another explanation is that we might need to perform more thorough hyperparameter searches. Even after our quite exhaustive search, it might be the case that there is a warm spot for which the MI regularizer might give better results in terms of the DS metric.

All in all, the H-VFAE+VP with the MMD regularizer seems to give the best trade-off between predictive accuracy and losing information about s in the supervised case, outperforming all the other models.

### 7.5. Experiment with the partial supervision of s

In the second experiment a more challenging task is considered, where the values of s are partially observed. First, we see that the classification accuracy of the unregularized models is high, this is however coupled with a significant increase in retained sensitive information. Again, the unregularized models seem ill suited for the task of fair classification. Note that the increase in classification accuracy compared to the supervised case is due to the increased weight on  $\alpha$ . Since there was no baseline from previous work on this task, we performed a more thorough hyperparameter search.

Interestingly, the MI regularizer outperforms the MMD regu-

larizer over the whole spectrum with higher accuracy scores and lower discriminative scores in the partially observed case. It seems the MI regularizer is better suited to handle the estimation uncertainty of  $q_{\phi}(\mathbf{s}|\mathbf{x})$ . Moreover, the MMD requires known s and since we have less fully supervised training examples, the estimation of the regularizer is worse. Possibly, the MI regularizer is more robust to this problem.

In conclusion, on the partially supervised task our proposed model outperforms the VFAE with respect to both classification accuracy and sensitive information retained.

### 8. Conclusion

In the paper we proposed the Hierarchical VampPrior Variational Auto-Encoder for learning fair representations. Additionally, we used the mutual information as a regularizer for obtaining fair representations, an alternative to the currently used MMD regularizer. In the experiments we considered two cases: (i) fully observable s scenario, and (ii) a case with partially supervised s. Especially the second task is interesting because in many real-life situations the information about s is missing, e.g., in domain adaptation the domain label could be unknown or hard to achieve.

The obtained results on two benchmark datasets show that our model together with the VampPrior obtains very promising results. The MMD regularizer seems to be preferred when the supervised training sample is large enough. Otherwise, using the MI regularizer provides the best results. Nevertheless, more thorough experiments are needed to reach a definite conclusion. Moreover, in this work we used a single Monte Carlo sample to approximate the MI regularizers. The obtained results might possibly be better if a larger sample would be used. We leave investigating these issues for future work.

## Acknowledgements

The research conducted by Jakub M. Tomczak was funded by the European Commission within the Marie Skodowska-Curie Individual Fellowship (Grant No. 702666, "Deep learning and Bayesian inference for medical imaging").

#### References

- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pp. 214–226. ACM, 2012.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *Advances in neural* information processing systems, pp. 2672–2680, 2014.

- Gretton, A., Borgwardt, K. M., Rasch, M., Schölkopf, B., and Smola, A. J. A kernel method for the two-sample-problem. In *Advances in neural information processing systems*, pp. 513–520, 2007.
- Hoffman, M. D. and Johnson, M. J. ELBO surgery: yet another way to carve up the variational evidence lower bound. In *Workshop in Advances in Approximate Bayesian Inference*, *NIPS*, 2016.
- Jang, E., Gu, S., and Poole, B. Categorical reparameterization with gumbel-softmax. arXiv preprint arXiv:1611.01144, 2016.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Louizos, C., Swersky, K., Li, Y., Welling, M., and Zemel, R. The variational fair autoencoder. *arXiv preprint arXiv:1511.00830*, 2015.
- Maddison, C. J., Mnih, A., and Teh, Y. W. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv* preprint arXiv:1611.00712, 2016.
- McNamara, D., Ong, C. S., and Williamson, R. C. Provably fair representations. *arXiv preprint arXiv:1710.04394*, 2017.
- Menon, A. K. and Williamson, R. C. The cost of fairness in classification. *arXiv preprint arXiv:1705.09055*, 2017.
- Rahimi, A. and Recht, B. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In *Advances in neural information processing systems*, pp. 1313–1320, 2009.
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- Sønderby, C. K., Raiko, T., Maaløe, L., Sønderby, S. K., and Winther, O. Ladder variational autoencoders. In *Advances in neural information processing systems*, pp. 3738–3746, 2016.
- Tomczak, J. M. and Welling, M. VAE with a VampPrior. *CoRR*, abs/1705.07120, 2017. URL http://arxiv.org/abs/1705.07120.
- Zafar, M. B., Valera, I., Gomez Rodriguez, M., and Gummadi, K. P. Fairness constraints: Mechanisms for fair classification. *arXiv preprint arXiv:1507.05259*, 2017.
- Zemel, R., Wu, Y., Swersky, K., Pitassi, T., and Dwork, C. Learning fair representations. In *International Conference on Machine Learning*, pp. 325–333, 2013.