

# Fairness-Aware Unsupervised Feature Selection

Xiaoying Xing<sup>1</sup>, Hongfu Liu<sup>2</sup>, Chen Chen<sup>3</sup> and Jundong Li<sup>3</sup>

<sup>1</sup>Tsinghua University, Beijing, China 100084

<sup>2</sup>Brandeis University, Waltham, MA, USA 02453

<sup>3</sup>University of Virginia, Charlottesville, VA, USA 22904

xingxy0505@gmail.com, hongfuliu@brandeis.edu, {zrh6du, jl6qk}@virginia.edu

## Abstract

Feature selection is a prevalent data preprocessing paradigm for various learning tasks. Due to the expensive cost of acquiring supervision information, unsupervised feature selection sparks great interests recently. However, existing unsupervised feature selection algorithms do not have fairness considerations and suffer from a high risk of amplifying discrimination by selecting features that are over associated with protected attributes such as gender, race, and ethnicity. In this paper, we make an initial investigation of the fairness-aware unsupervised feature selection problem and develop a principled framework, which leverages kernel alignment to find a subset of high-quality features that can best preserve the information in the original feature space while being minimally correlated with protected attributes. Specifically, different from the mainstream in-processing debiasing methods, our proposed framework can be regarded as a model-agnostic debiasing strategy that eliminates biases and discrimination before downstream learning algorithms are involved. Experimental results on multiple real-world datasets demonstrate that our framework achieves a good trade-off between utility maximization and fairness promotion.

## 1 Introduction

Feature selection is an effective data preprocessing strategy for a myriad of data mining and machine learning tasks [Guyon and Elisseeff, 2003; Li *et al.*, 2017a]. It aims to select a subset of relevant features from the original feature space while eliminating the adverse impact of irrelevant, redundant, and noisy features. In contrast to the prevalent deep learning models for representation learning, feature selection gives learning models better readability and interpretability by maintaining the physical meanings of original features, thus is often preferred in high-stake applications (e.g., clinical diagnosis [Inbarani *et al.*, 2014], employment [Sobnath *et al.*, 2020], and financial analytics [Liang *et al.*, 2015]). According to the label availability, traditional feature selection algorithms can be mainly categorized as supervised methods and unsupervised methods [Li *et al.*, 2017a]. As supervision

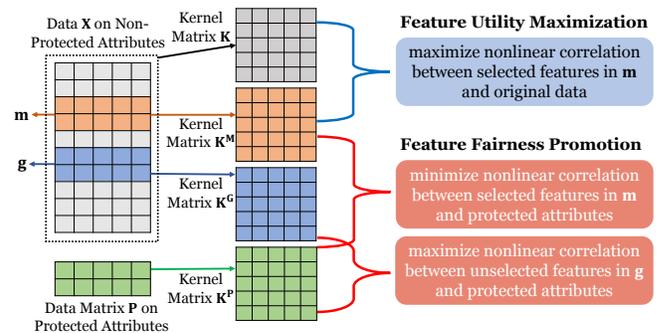


Figure 1: An illustration of the proposed fairness-aware unsupervised feature selection framework FUFs, which has two components: feature utility maximization and feature fairness promotion.

information is often costly to amass in many real-world scenarios, unsupervised feature selection methods has attracted an increasing amount of attention in recent years.

Despite the successful adoption of feature selection algorithms in various high-stake decision-making scenarios, the existing selection algorithms often do not have the fairness considerations and may suffer from a risk of amplifying stereotypes and exhibiting discriminatory actions toward specific groups or populations by over associating protected attributes<sup>1</sup> (e.g., gender, race, and ethnicity) [Chouldechova and Roth, 2018; Du *et al.*, 2020; Mehrabi *et al.*, 2019; Grgić-Hlača *et al.*, 2018]. Although it is intuitive to manually remove these protected attributes in the selected feature subset to avoid direct discrimination, a number of non-protected attributes that are highly correlated with the protected attributes may still be selected by the algorithms and result in unintentional discrimination problems (e.g., residential zip code of a person may indicate the race information because of the population of residential areas) [Zhang *et al.*, 2016; Kallus *et al.*, 2019].

In this paper, we make an initial investigation of the fairness issues of unsupervised feature selection and develop a general model-agnostic debiasing strategy at the input level. Our proposed efforts have the potential to alleviate unwanted biases before applying downstream learning algorithms and are complementary to the mainstream in-processing algorithm-

<sup>1</sup>We use protected and sensitive features interchangeably. Meanwhile, we also use attributes and features interchangeably.

mic fairness research [Mehrabi *et al.*, 2019]. However, developing a fairness-aware unsupervised feature selection framework remains a daunting task, mainly because of the following challenges. Firstly, feature selection should achieve a good trade-off between fairness and feature utility—finding a subset of features that do not exhibit discrimination while still benefiting downstream tasks. However, without label information as supervision signals, we are in short of effective evaluation criteria to quantify these two targets simultaneously. Secondly, due to the trade-off between utility and fairness, it is difficult to simultaneously achieve the maximums of both. Thus, it is necessary to explicitly exclude the fairness-related features from the selected set, which have strong correlations with protected attributes.

To tackle the aforementioned challenges, we propose a novel Fairness-aware Unsupervised Feature Selection (FUFS) framework (as shown in Fig. 1). In essence, to ensure that the selected features do not cause much utility loss for downstream learning algorithms, we select features that can maximally preserve the information in the original feature space. Additionally, we impose fairness constraints to enforce the protected attributes being minimally correlated with the selected features while over associating with a small number of unselected features. All the above considerations are modeled in a joint optimization framework. The major contributions of our work can be summarized as follows:

- We address a crucial and newly emerging problem, fairness-aware unsupervised feature selection, which is essential to debiasing the input data before downstream learning algorithms are involved.
- We propose a novel FUFS framework, which selects high-quality features by preserving information embedded in the original feature space and obeying the fairness considerations to eliminate sensitive information.
- We formulate two desiderata of fairness-aware unsupervised feature selection (i.e., feature utility maximization and feature fairness promotion) as an optimization problem with a principled solution.
- We validate the selected features by feature utility and fairness measurements, where empirical evaluations corroborate the superiority of our proposed framework.

## 2 Problem Statement

In this section, we first summarize the notations used in this paper, and then formally define our research problem of *fairness-aware unsupervised feature selection*.

We use bold uppercase letters for matrices (e.g.,  $\mathbf{A}$ ), bold lowercase letters for vectors (e.g.,  $\mathbf{a}$ ), normal lowercase letters for scalars (e.g.,  $a$ ). Also, we use  $a_i$  to denote the  $i$ -th element of the vector  $\mathbf{a}$ ,  $\mathbf{A}_{*j}$  to denote the  $j$ -th column of matrix  $\mathbf{A}$ ,  $\mathbf{A}_{ij}$  to denote the  $i$ -th row and  $j$ -th column of matrix  $\mathbf{A}$ , and  $\text{diag}(\mathbf{a})$  to represent the diagonalization of vector  $\mathbf{a}$ . Meanwhile,  $\mathbf{1}$  denotes a column vector whose elements are all 1,  $\mathbf{I}$  denotes an identity matrix.  $\|\mathbf{a}\|_1$  denotes the  $\ell_1$ -norm of the vector  $\mathbf{a}$ , respectively.  $\|\mathbf{A}\|_F$  denotes the Frobenius norm of the matrix  $\mathbf{A}$ .  $\text{Tr}(\mathbf{A})$  denotes the trace of the matrix  $\mathbf{A}$  when it is a square matrix.

Table 1: Notations and descriptions.

Notation	Description
$n$	Number of instances
$d$	Number of non-protected attributes
$p$	Number of protected attributes
$\mathbf{X} \in \mathbb{R}^{d \times n}$	Data matrix on non-protected attributes
$\mathbf{P} \in \mathbb{R}^{p \times n}$	Data matrix on protected attributes
$\mathbf{m} \in \{0, 1\}^d$	Indicator vector for selected features
$\mathbf{g} \in \{0, 1\}^d$	Indicator vector for features that are highly correlated with protected attributes
$\mathbf{M} \in \mathbb{R}^{d \times n}$	Data subset on the indicator vector $\mathbf{m}$
$\mathbf{G} \in \mathbb{R}^{d \times n}$	Data subset on the indicator vector $\mathbf{g}$
$\mathbf{K} \in \mathbb{R}^{n \times n}$	Kernel matrix on the input data $\mathbf{X}$
$\mathbf{K}^{\mathbf{P}} \in \mathbb{R}^{n \times n}$	Kernel matrix on matrix $\mathbf{P}$
$\mathbf{K}^{\mathbf{M}} \in \mathbb{R}^{n \times n}$	Kernel matrix on matrix $\mathbf{M}$
$\mathbf{K}^{\mathbf{G}} \in \mathbb{R}^{n \times n}$	Kernel matrix on matrix $\mathbf{G}$

In this work, we assume there are  $n$  data instances, the matrix  $\mathbf{P} \in \mathbb{R}^{p \times n}$  denotes the set of  $p$  protected attributes for instances (e.g., age, gender, and race), and the matrix  $\mathbf{X} \in \mathbb{R}^{d \times n}$  denotes the set of  $d$  non-protected attributes (in most cases we have  $p \ll d$ ). The main symbols are summarized in Table 1. With these notations, we formally define our research problem as follows.

**Problem Definition** (Fairness-Aware Unsupervised Feature Selection). *Given the input data  $\mathbf{X} \in \mathbb{R}^{d \times n}$  and  $\mathbf{P} \in \mathbb{R}^{p \times n}$  with  $d$  non-protected attributes and  $p$  protected attributes, respectively, the problem of fairness-aware unsupervised feature selection aims to select a subset of  $k$  features among  $d$  non-protected attributes ( $k \ll d$ ) which can maximally preserve the information in the original feature space while being minimally correlated with the protected attributes.*

## 3 The Proposed Framework - FUFS

In this section, we present our proposed Fairness-aware Unsupervised Feature Selection (FUFS) framework in detail. An overview of the proposed framework is illustrated in Fig. 1.

### 3.1 Maximizing Feature Utility

As label information is not available in an unsupervised scenario, we need to seek alternative evaluation criteria to assess the importance of features. One principled metric is capable of ensuring that the selected features can well capture the information embedded in the original feature space. In other words, we would like to maximize the correlation between the selected features and the original ones. However, since the original features could be high-dimensional, complex nonlinear correlations could exist between these two features spaces. Hence, we aim to measure their nonlinear correlation with kernel alignment [Cristianini *et al.*, 2006; Wei *et al.*, 2016] techniques.

Suppose the vector  $\mathbf{m} \in \{0, 1\}^d$  is the feature selection indicator vector such that  $\mathbf{1}^\top \mathbf{m} = k$ , where  $m_i = 1$  if the  $i$ -th feature is selected, otherwise  $m_i = 0$ . The data matrix on the selected features can be obtained as  $\mathbf{M} = \text{diag}(\mathbf{m})\mathbf{X}$ . Then we define a kernel  $\kappa$  which implicitly computes the similarity between instances in a high-dimensional reproducing kernel Hilbert space (RKHS) [Aronszajn, 1950], such that

$\mathbf{K}_{ij} = \kappa(\mathbf{X}_{*i}, \mathbf{X}_{*j})$  and  $\mathbf{K}_{ij}^M = \kappa(\mathbf{M}_{*i}, \mathbf{M}_{*j})$ . In practice, we can choose polynomial kernel or RBF kernel. Denoting the centering matrix as  $\mathbf{H} = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$ , these two kernel matrices after centering can be denoted as  $\mathbf{K}_c = \mathbf{H}\mathbf{K}\mathbf{H}$  and  $\mathbf{K}_c^M = \mathbf{H}\mathbf{K}^M\mathbf{H}$ , respectively.

With these two centered kernel matrices, we can characterize the inherent nonlinear correlation between these two feature spaces with the centered kernel alignment:

$$\rho(\mathbf{K}, \mathbf{K}^M) = \text{Tr}(\mathbf{K}_c\mathbf{K}_c^M) = \text{Tr}(\mathbf{H}\mathbf{K}\mathbf{H}\mathbf{H}\mathbf{K}^M\mathbf{H}). \quad (1)$$

With the observation that  $\mathbf{H}\mathbf{H} = \mathbf{H}$  and  $\text{Tr}(\mathbf{A}\mathbf{B}) = \text{Tr}(\mathbf{B}\mathbf{A})$  (where  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ ), we can further simplify  $\rho(\mathbf{K}, \mathbf{K}^M)$  as  $\text{Tr}(\mathbf{H}\mathbf{K}\mathbf{H}\mathbf{K}^M)$ . Our goal expects that the selected features in  $\mathbf{m}$  can maximally preserve the information embedded in the original feature space.

### 3.2 Promoting Feature Fairness

Although maximizing the centered kernel alignment function in Eq. (1) helps to select important features that preserve the information of original features, it does not have any fairness considerations such that the selected features may have a risk being associated with the protected attributes in the matrix  $\mathbf{P}$ . To tackle this issue, we further impose fairness constraints to make the selected features in  $\mathbf{M}$  not well aligned with the protected attributes  $\mathbf{P}$ . To achieve this goal, suppose  $\mathbf{K}^P \in \mathbb{R}^{n \times n}$  is the kernel matrix computed from  $\mathbf{P}$ , we can also leverage centered kernel alignment to minimize the nonlinear correlation between  $\mathbf{M}$  and  $\mathbf{P}$  in the kernel space:

$$\rho(\mathbf{K}^M, \mathbf{K}^P) = \text{Tr}(\mathbf{H}\mathbf{K}^M\mathbf{H}\mathbf{K}^P). \quad (2)$$

In this way, we can guarantee that the selected features in  $\mathbf{M}$  do not have high correlation with the sensitive information.

Additionally, to further enforcing that the sensitive information is eliminated in the selected features, a small number of unselected features need to exhibit a high degree of correlation with the protected attributes. Toward this goal, we further define a decomposition indicator  $\mathbf{g} \in \{0, 1\}^d$  to indicate the index of non-protected attributes that are highly correlated with  $\mathbf{P}$ , where  $\mathbf{1}^\top \mathbf{g} = l$ , and  $l$  denotes the number of sensitive features we need to eliminate. Ideally, the nonzero indices of  $\mathbf{g}$  should not overlap with those of  $\mathbf{m}$ . Hence, the data matrix  $\mathbf{G}$  corresponding to  $\mathbf{g}$  can be obtained as follows:

$$\mathbf{G} = \text{diag}(\mathbf{g})(\mathbf{I} - \text{diag}(\mathbf{m}))\mathbf{X}. \quad (3)$$

Assume the corresponding kernel matrix is  $\mathbf{K}^G \in \mathbb{R}^{n \times n}$ , then the centered kernel alignment can also be leveraged to maximize the nonlinear correlation between  $\mathbf{G}$  and  $\mathbf{P}$ :

$$\rho(\mathbf{K}^G, \mathbf{K}^P) = \text{Tr}(\mathbf{H}\mathbf{K}^G\mathbf{H}\mathbf{K}^P). \quad (4)$$

### 3.3 Objective Function of FUFs

The above two subsections discuss two desiderata of fairness-aware unsupervised feature selection—preserving the information of original features and imposing fairness constraints. Combining them together, we obtain a joint constrained optimization problem as follows:

$$\begin{aligned} \min_{\mathbf{m}, \mathbf{g}} & -\text{Tr}(\mathbf{H}\mathbf{K}\mathbf{H}\mathbf{K}^M) + \alpha \text{Tr}(\mathbf{H}\mathbf{K}^M\mathbf{H}\mathbf{K}^P) \\ & - \alpha \text{Tr}(\mathbf{H}\mathbf{K}^G\mathbf{H}\mathbf{K}^P) \\ \text{s.t. } & \mathbf{m}, \mathbf{g} \in \{0, 1\}^d, \mathbf{1}^\top \mathbf{m} = k, \mathbf{1}^\top \mathbf{g} = l. \end{aligned} \quad (5)$$

---

#### Algorithm 1 The Proposed Framework FUFs.

---

**Input:** protected attribute matrix  $\mathbf{P} \in \mathbb{R}^{p \times n}$ , non-protected attribute matrix  $\mathbf{X} \in \mathbb{R}^{d \times n}$ , number of features  $k$  to select, hyperparameters  $\alpha, \beta$ .

**Output:** the top- $k$  most important features that both preserve the information of original data and do not correlate with protected attributes.

- 1: Initialize indicator vectors  $\mathbf{m}, \mathbf{g}$
  - 2: Construct kernel matrices  $\mathbf{K}$  from  $\mathbf{X}$ ,  $\mathbf{K}^P$  from  $\mathbf{P}$ ,  $\mathbf{K}^M$  from  $\text{diag}(\mathbf{m})\mathbf{X}$ ,  $\mathbf{K}^G$  from  $\text{diag}(\mathbf{g})(\mathbf{I} - \text{diag}(\mathbf{m}))\mathbf{X}$
  - 3: **while** not convergence **do**
  - 4:  $g_i \leftarrow g_i - \eta \partial \mathcal{L} / \partial g_i, i = 1, \dots, d$
  - 5:  $g_i = \min(1, \max(0, g_i)), i = 1, \dots, d$
  - 6:  $m_i \leftarrow m_i - \eta \partial \mathcal{L} / \partial m_i, i = 1, \dots, d$
  - 7:  $m_i = \min(1, \max(0, m_i)), i = 1, \dots, d$
  - 8: Update matrix  $\mathbf{M}$  and kernel matrix  $\mathbf{K}^M$
  - 9: Update matrix  $\mathbf{G}$  and kernel matrix  $\mathbf{K}^G$
  - 10: **end while**
  - 11: Rank features according to the entries in  $\mathbf{m}$
- 

In the above formulation,  $\alpha$  is a hyperparameter that can control how strong we would like to enforce the fairness of unsupervised feature selection.

The optimization problem in Eq. (6) is not easy to solve because it is not joint convex w.r.t.  $\mathbf{m}$  and  $\mathbf{g}$  simultaneously. Although we can employ alternating optimization scheme for a local optimum, the whole optimization still remains difficult due to the discrete nature of variables  $\mathbf{m}$  and  $\mathbf{g}$ . To address this issue, we relax the discrete constraints by reformulating it as a real-valued vector in the range of  $[0, 1]$ . By rewriting the summation constraints  $\mathbf{1}^\top \mathbf{m} = k$  and  $\mathbf{1}^\top \mathbf{g} = l$  in the form of Lagrangian, we have the following optimization problem:

$$\begin{aligned} \min_{\mathbf{m}, \mathbf{g}} \mathcal{L} &= -\text{Tr}(\mathbf{H}\mathbf{K}\mathbf{H}\mathbf{K}^M) + \alpha \text{Tr}(\mathbf{H}\mathbf{K}^M\mathbf{H}\mathbf{K}^P) \\ & - \alpha \text{Tr}(\mathbf{H}\mathbf{K}^G\mathbf{H}\mathbf{K}^P) + \beta(\|\mathbf{m}\|_1 + \|\mathbf{g}\|_1) \\ \text{s.t. } & \mathbf{m}, \mathbf{g} \in [0, 1]^d, \end{aligned} \quad (6)$$

where the  $\ell_1$ -norm is introduced for the sparsity of model parameters  $\mathbf{m}$  and  $\mathbf{g}$ . The hyperparameter  $\beta$  is used to control the number of selected features that are relevant and do not correlate with protected attributes and the number of unselected features that are highly correlated with protected attributes, respectively.

**Updating  $\mathbf{m}$  and  $\mathbf{g}$ .** We update two model parameters  $\mathbf{m}$  and  $\mathbf{g}$  alternatively until the objective function converges to a local optimum. The update rules are as follows:

$$m_i \leftarrow P[m_i - \eta \partial \mathcal{L} / \partial m_i], \quad g_i \leftarrow P[g_i - \eta \partial \mathcal{L} / \partial g_i], \quad (7)$$

where  $P[x]$  is a box projection operator which projects  $x$  into a bounded range. Specifically, since we relax the constraints of  $m_i$  and  $g_i$  in the range of  $[0, 1]$ , we have  $P[x] = 0$  if  $x < 0$ ,  $P[x] = 1$  if  $x > 1$ , and otherwise  $P[x] = x$ . In the above update rules,  $\eta$  is the learning rate.

With these updating rules, the pseudo code of the proposed FUFs framework is illustrated in Algorithm 1.

## 4 Experimental Evaluations

In this section, we conduct experiments on real-world datasets to evaluate the performance of the proposed fairness-aware unsupervised feature selection framework FUFs in terms of both utility and fairness measurements. Before presenting the detailed experiments and findings, we first introduce the experimental settings.

### 4.1 Experimental Setup

**Datasets.** We perform experiments on four public available datasets. (1) CRIME<sup>2</sup>: This dataset combines census data, law enforcement data, and crime data of US communities. We define the percentage of population for African American as a protected attribute. We define two clusters by the number of violent crimes, and the cutoff threshold is 0.15 crimes per 100K population. In total, we have 2,215 communities described by 147 different attributes. (2) ADOLESCENT<sup>3</sup>: This dataset comes from a longitudinal study of adolescents in Grades 7-12. The attributes are obtained from personal information of the interviewees and their answers to an exhaustive questionnaire. Bio-sex of the interviewee is regarded as the protected attribute and we categorize the interviewees into two clusters by whether their Picture Vocabulary test score is more than 65. In total, this dataset contains 6,504 instances and 2,793 attributes. (3) GOOGLE+<sup>4</sup>: This dataset comes from Google+, which contains user features and social relations within multiple social circles. Each instance refers to a user and attributes are obtained from personal information of users. Gender is regarded as the protected attribute. We have two clusters defined by the social circles that the users belong to without overlapping. The dataset consists of 2,437 users and 1,695 features. (4) TOXICITY<sup>5</sup>: This dataset is obtained from a Toxic Comment Classification Challenge, where each comment is considered as an instance. We apply a tokenizer to transform text data to numerical values. The identity label ‘female’ is regarded as the protected attribute. The features are from identity labels and comment texts. There are two clusters defined by whether the comment is regarded toxic or not. We collect a subset of 200 instances with 4,253 features.

**Evaluation Criteria.** For unsupervised feature selection, clustering performance is often used as an evaluation metric [Li *et al.*, 2017a] to assess the quality of selected features. Specifically, we use *Clustering Accuracy (ACC)* and *Normalized Mutual Information (NMI)* to compare the obtained cluster labels with the ground truth cluster labels, and higher values often imply higher quality of selected features. Meanwhile, we use the widely used metrics *Balance* [Li *et al.*, 2020] and define a new fairness metric *Proportion* as a compliment since *Balance* may be too restrict to reflect the distribution of the clustering. These two metrics are used to quantify how well the selected features can eliminate discrimination—the selected features are considered fairer if

they can lead to a more balanced cluster structure toward protected attributes (i.e., higher value of *Balance* and lower value of *Proportion*). These four metrics are defined as follows:

$$ACC = \frac{\sum_{i=1}^n \delta(y_i, \text{map}(\hat{y}_i))}{n}, \quad (8)$$

$$NMI = \frac{\sum_{c \in C} \sum_{c' \in C'} p(c, c') \log(p(c, c') / p(c)p(c'))}{\text{mean}(H(C), H(C'))}, \quad (9)$$

$$Balance = \min_i \frac{\min_g |C_i \cap X_g|}{|C_i|}, \quad (10)$$

$$Proportion = \sum_i \frac{\max_g |C_i \cap X_g|}{|C_i|}, \quad (11)$$

where  $\hat{y}_i$  is the clustering result,  $y_i$  is the true cluster label,  $\text{map}(\cdot)$  is a permutation mapping function that maps  $y_i$  to the equivalent label from the ground truth and  $\delta$  is the indicator function such that  $\delta(x, y) = 1$  if  $x = y$ , and  $\delta(x, y) = 0$  otherwise.  $H$  denotes the entropy for a partition set.  $C$  and  $C'$  denote the obtained clusters and the ground truth, respectively.  $C_i$  and  $X_g$  denote the  $i$ -th cluster and the  $g$ -th protected subgroup regarding the sensitive attribute.

**Competitive Methods and Implementation.** We compare our proposed framework FUFs with the following unsupervised feature selection methods:

- **LapScore** [He *et al.*, 2006]: Laplacian Score selects important features that best align with the local manifold structure of data.
- **MCFS** [Cai *et al.*, 2010]: MCFS selects features that can best preserve the multi-cluster structure of data.
- **UDFS** [Yang *et al.*, 2011]: UDFS is a pseudo-label based approach that exploits local discriminative information and  $\ell_{2,1}$ -norm regularization.
- **NDFS** [Li *et al.*, 2012]: NDFS selects important features by performing joint nonnegative spectral analysis and  $\ell_{2,1}$ -norm regularization.
- **REFS** [Li *et al.*, 2017b]: REFS is a reconstruction based approach that learns the reconstruction function from data automatically for unsupervised feature selection.

We follow the suggestions of the original papers [He *et al.*, 2006; Cai *et al.*, 2010; Yang *et al.*, 2011; Li *et al.*, 2012; Li *et al.*, 2017b] to specify the hyperparameters for these baseline methods. For our proposed FUFs framework, we set the hyperparameters as  $\alpha = 1, \beta = 0.1$  on CRIME and GOOGLE+ while  $\alpha = 0.01, \beta = 10$  on ADOLESCENT and TOXICITY. The original distribution of the protected groups in CRIME and GOOGLE+ is more unbalanced thus a larger value of  $\alpha$  is necessary to eliminate discrimination. Whereas ADOLESCENT and TOXICITY have more features and a larger value of  $\beta$  is necessary for unsupervised feature selection. Besides, we specify the kernel function in FUFs as the RBF kernel. For all the compared methods, we first apply unsupervised feature selection to select the top- $k$  ranked features and employ K-means clustering algorithm on the selected features. The clustering results and the ground truth cluster labels are compared and the values on the aforementioned four evaluation metrics can then be

<sup>2</sup><http://archive.ics.uci.edu/ml/datasets/Communities+and+Crime+Unnormalized>

<sup>3</sup><https://www.thearda.com/>

<sup>4</sup><http://snap.stanford.edu/data/ego-Gplus.html>

<sup>5</sup><https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification/>

Table 2: Results on CRIME w.r.t. cluster validity and fairness.

Method	ACC	NMI	Balance	Proportion
LapScore	0.644 (35)	0.024 (35)	0.192 (10)	1.492 (35)
NDFS	0.627 (20)	0.021 (30)	0.201 (10)	1.527 (15)
UDFS	0.728 (30)	0.150 (30)	0.107 (40)	1.456 (25)
REFS	0.774 (15)	0.082 (15)	0.208 (25)	1.552 (40)
MCFS	0.683 (25)	0.101 (20)	0.182 (20)	1.511 (25)
FUFS (ours)	0.758 (15)	0.141 (35)	0.204 (35)	1.446 (10)

Table 3: Results on ADOLES. w.r.t. cluster validity and fairness.

Method	ACC	NMI	Balance	Proportion
LapScore	0.555 (10)	0.006 (10)	0.379 (10)	1.163 (10)
NDFS	0.554 (15)	0.006 (15)	0.380 (35)	1.184 (35)
UDFS	0.556 (10)	0.007 (10)	0.359 (15)	1.184 (15)
REFS	0.544 (10)	0.004 (10)	0.380 (10)	1.184 (10)
MCFS	0.562 (10)	0.010 (10)	0.380 (15)	1.184 (15)
FUFS (ours)	0.553 (35)	0.013 (35)	0.407 (10)	1.148 (10)

obtained. Since the results of K-means depend on initialization, we repeat K-means 50 times and report the average results. Choosing the optimal number of selected features is still an open problem, thus we follow conventional settings [Li *et al.*, 2017a] to vary the number of selected features as {10%, 15%, 20%, 25%, 30%, 35%, 40%} of the total number of features and report the best results regarding different evaluation metrics.

## 4.2 Performance Evaluation

We compare FUFS with different baseline methods in terms of feature utility (*ACC* and *NMI*) and fairness metrics (*Balance* and *Proportion*). The experimental results are shown in Tables 2-5. The number in parentheses denotes the percentage of features when the best performance is achieved. Values in red cell indicates the highest result, and blue cell indicates the second highest one. We make the following observations:

- FUFS significantly outperforms the baseline methods in terms of *Balance* and *Proportion* with the best performance in almost all cases and the second best performance in terms of *Balance* on CRIME. Existing unsupervised feature selection methods often do not have the fairness considerations and deliver the unfair results, while our proposed FUFS framework can obtain the most balanced clustering results across different protected subgroups.
- FUFS achieves a good balance between feature utility and feature fairness. While achieving a good performance w.r.t. different fairness metrics, the feature utility is also well maintained as the clustering performance on the selected features is not jeopardized. For example, on CRIME and TOXICITY, FUFS achieves the second best performance in terms of *ACC* and *NMI* while on ADOLESCENT and GOOGLE+, FUFS achieves the best *NMI* values and does not have obvious difference w.r.t. *ACC* compared with the best baseline method.
- The proposed FUFS framework can achieve great performance in terms of fairness with a small number of features. Specifically, on ADOLESCENT and GOOGLE+, FUFS achieves the best results in terms of *Balance* and *Proportion* compared with the baseline methods with

Table 4: Results on GOOGLE+ w.r.t. cluster validity and fairness.

Method	ACC	NMI	Balance	Proportion
LapScore	0.723 (40)	0.114 (15)	0.004 (40)	1.865 (10)
NDFS	0.724 (40)	0.113 (40)	0.000 (10)	1.885 (15)
UDFS	0.723 (30)	0.115 (20)	0.000 (15)	1.881 (10)
REFS	0.724 (35)	0.114 (20)	0.004 (20)	1.886 (15)
MCFS	0.719 (40)	0.109 (15)	0.228 (10)	1.412 (35)
FUFS (ours)	0.721 (10)	0.164 (15)	0.301 (10)	1.308 (10)

Table 5: Results on TOXICITY w.r.t. cluster validity and fairness.

Method	ACC	NMI	Balance	Proportion
LapScore	0.803 (30)	0.012 (30)	0.009 (40)	1.568 (10)
NDFS	0.675 (40)	0.007 (40)	0.240 (20)	1.327 (15)
UDFS	0.663 (40)	0.006 (30)	0.284 (10)	1.309 (10)
REFS	0.674 (40)	0.007 (35)	0.334 (40)	1.579 (40)
MCFS	0.650 (35)	0.006 (35)	0.285 (10)	1.391 (15)
FUFS (ours)	0.701 (40)	0.008 (25)	0.409 (15)	1.136 (15)

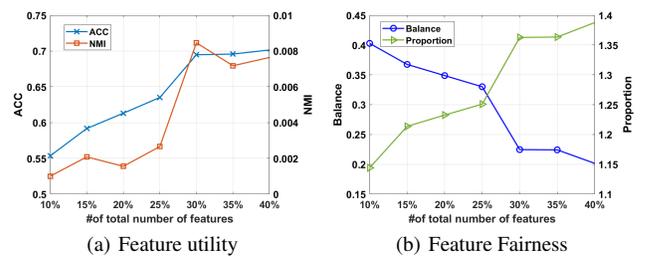


Figure 2: Utility and fairness performance variation on TOXICITY w.r.t. different numbers of selected features.

merely 10% of the total number of features. On TOXICITY, FUFS achieves the best results of fairness with 15% of the total number of features.

## 4.3 In-Depth Exploration of FUFS

**Effects of the Number of Selected Features.** Choosing an optimal number of features is still an open problem in unsupervised feature selection research, thus we vary the number of selected features as {10%, 15%, 20%, 25%, 30%, 35%, 40%} of the total feature number and investigate how the feature utility and feature fairness performance change. We only show the results on TOXICITY (Fig. 2) as we have similar observations on other datasets. As we can see, the clustering results (w.r.t. *ACC* and *NMI*) first increase and then keep stable when the number of selected features increase. Meanwhile, the fairness performance (w.r.t. *Balance* and *Proportion*) is the best when only 10% of features are selected (it should be noted that lower values of *Proportion* denotes fairer results). The fairness performance gradually decreases when the number of selected features continuously increases, the reason is that more features that are correlated with sensitive features could be included in the selected feature subset.

**Effects of the Decomposition Indicator Vector  $\mathbf{g}$ .** In order to investigate the effect of the decomposition indicator vector  $\mathbf{g}$ , we remove it from our framework and compare its performance with the original FUFS framework. The results on CRIME and GOOGLE+ shown in Fig. 3 imply that the introduction of the decomposition indicator vector  $\mathbf{g}$  is necessary and improves both the utility and fairness performance.

Table 6: Fairness results (w.r.t. *Balance* and *Proportion*) comparison based on the top-ranked features in  $\mathbf{m}$  and  $\mathbf{g}$ .

Dataset	Top- $k$ ranked features in $\mathbf{m}$		Top- $k$ ranked features in $\mathbf{g}$	
	<i>Balance</i>	<i>Proportion</i>	<i>Balance</i>	<i>Proportion</i>
CRIME	0.204 (35)	1.446 (10)	0.188 (40)	1.468 (15)
ADOLESCENT	0.407 (10)	1.148 (10)	0.308 (40)	1.179 (30)
GOOGLE+	0.301 (10)	1.308 (10)	0.000 (20)	1.863 (20)
TOXICITY	0.409 (15)	1.136 (15)	0.063 (15)	1.629 (40)

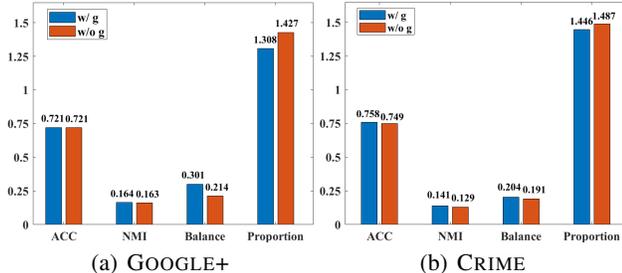


Figure 3: Clustering performance w/ and w/o indicator vector  $\mathbf{g}$ .

We also compare the fairness performance based on the top-ranked features in  $\mathbf{m}$  and  $\mathbf{g}$  and the results are shown in Table 6. The number in parentheses denotes the percentage of features when the best performance is achieved. It is obvious that the clustering results based on the top-ranked features in the vector  $\mathbf{m}$  are more fair than those in the vector  $\mathbf{g}$ . It shows the effectiveness of introducing the decomposition indicator vector  $\mathbf{g}$ , which can help eliminate the sensitive information in the selected feature subset.

**Parameter Study.** The proposed framework FUSF has two important hyperparameters. The first parameter  $\alpha$  controls how strong we would like to enforce the fairness of unsupervised feature selection. The other parameter  $\beta$  controls the sparsity of the proposed model. We first fix  $\beta = 0.1$  and then vary  $\alpha$  among  $\{0.001, 0.01, 0.1, 1, 10, 100, 1000\}$ . Next, We first fix  $\alpha = 1$  and then vary  $\beta$  among  $\{0.001, 0.01, 0.1, 1, 10, 100, 1000\}$ . The performance on GOOGLE+ is shown in Fig. 4. It should be noted that as the X-axis is plotted in a log scale, we do not expect to see a smooth curve. Due to space limit, we only show the parameter study results on GOOGLE+ in terms of *ACC* and *Balance*. The results imply that the clustering performance is relatively stable when  $\alpha = 1, \beta \in [0.001, 0.1]$  or  $\alpha \in [0.001, 0.1], \beta = 0.1$ . When the parameter  $\alpha$  increases, the algorithm becomes more partial to the fairness consideration with decreasing *ACC* and increasing *Balance*. It can be observed that the fairness performance decreases a lot if  $\beta$  is specified as a very large value.

## 5 Related Work

**Unsupervised Feature Selection.** Due to the expensive annotation cost, unsupervised feature selection has sparked great interests in recent years. To quantify the importance of features, unsupervised methods often rely on alternative evaluation criteria based on different characteristics of data. Specifically, similarity based methods [He *et al.*, 2006; Zhao and Liu, 2007] select features that can best preserve the local manifold structure of data. Another family of methods aim to select important features that can best reconstruct [Farahat *et al.*, 2011; Li *et al.*, 2017b] or max-

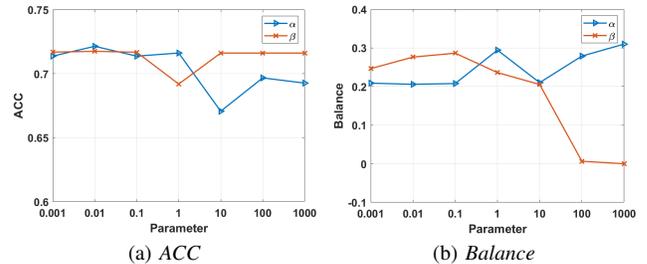


Figure 4: Performance variation on GOOGLE+ w.r.t. different parameter settings. X-axis is not in a linear scale.

imally preserve information embedded in the original features [Wei *et al.*, 2016]. A number of studies learn the pseudo label from data by exploiting local/global discriminative information and select features to predict these pseudo labels with  $\ell_{2,1}$ -norm based regression [Li *et al.*, 2012; Li *et al.*, 2018]. Recently, data reconstruction [Li *et al.*, 2017b; Farahat *et al.*, 2011] emerged as a new criterion to evaluate feature relevance, which evaluates the capability of features in approximating the original data through a reconstruction function.

**Fairness of Unsupervised Learning Methods.** To our best knowledge, we are the first to study the fairness issue of unsupervised feature selection. Here we review some related fairness topics in terms of clustering and representation learning. The initial work [Chierichetti *et al.*, 2017] defines fair variants of classical clustering problems such as  $k$ -center and  $k$ -median and proposes the concepts of fairlets and fairlet decomposition, which is further extended to  $k$ -means++ algorithm by [Schmidt *et al.*, 2018]. Other related works focus on scalable fair clustering [Backurs *et al.*, 2019], fair spectral clustering [Kleindessner *et al.*, 2019], and deep fair clustering [Li *et al.*, 2020]. Another family of work aims to learn fair representations. Fair PCA [Creager *et al.*, 2019] is a two-step algorithm for dimension reduction. Fair Autoencoders [Louizos *et al.*, 2015; Moyer *et al.*, 2018] encourage independence between sensitive and latent factors of variation for representation learning. Extended work [Creager *et al.*, 2019] learns general-purpose flexible fair representations regarding multiple sensitive attributes.

## 6 Conclusion

Unsupervised feature selection plays an essential role in preparing high-dimensional and unlabeled data for various learning tasks and has been increasingly used in high-stake applications. Despite its fundamental importance, the fairness of unsupervised feature selection has largely remained nascent. In this paper, we addressed a novel problem of fairness-aware unsupervised feature selection and developed a principled framework FUSF. The proposed framework leverages the technique of kernel alignment to select high-quality features that achieve a good balance between improving downstream learning tasks and eliminating sensitive information that is highly correlated with protected attributes. These two desiderata were modeled together in a joint optimization framework. Experimental evaluations on real-world datasets demonstrated the superiority of the proposed FUSF framework in terms of feature utility and feature fairness.

## References

- [Aronszajn, 1950] Nachman Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, 1950.
- [Backurs *et al.*, 2019] Arturs Backurs, Piotr Indyk, Krzysztof Onak, Baruch Schieber, Ali Vakilian, and Tal Wagner. Scalable fair clustering. *arXiv preprint arXiv:1902.03519*, 2019.
- [Cai *et al.*, 2010] Deng Cai, Chiyuan Zhang, and Xiaofei He. Unsupervised feature selection for multi-cluster data. In *KDD*, 2010.
- [Chierichetti *et al.*, 2017] Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, and Sergei Vassilvitskii. Fair clustering through fairlets. In *NeurIPS*, 2017.
- [Chouldechova and Roth, 2018] Alexandra Chouldechova and Aaron Roth. The frontiers of fairness in machine learning. *arXiv preprint arXiv:1810.08810*, 2018.
- [Creager *et al.*, 2019] Elliot Creager, David Madras, Jörn-Henrik Jacobsen, Marissa A Weis, Kevin Swersky, Toniann Pitassi, and Richard Zemel. Flexibly fair representation learning by disentanglement. *arXiv preprint arXiv:1906.02589*, 2019.
- [Cristianini *et al.*, 2006] Nello Cristianini, Jaz Kandola, Andre Elisseeff, and John Shawe-Taylor. On kernel target alignment. In *Innovations in Machine Learning*. 2006.
- [Du *et al.*, 2020] Mengnan Du, Fan Yang, Na Zou, and Xia Hu. Fairness in deep learning: A computational perspective. *IEEE Intelligent Systems*, 2020.
- [Farahat *et al.*, 2011] A. K. Farahat, A. Ghodsi, and M. S. Kamel. An efficient greedy method for unsupervised feature selection. In *ICDM*, pages 161–170, 2011.
- [Grgić-Hlača *et al.*, 2018] Nina Grgić-Hlača, Muhammad Bilal Zafar, Krishna P Gummadi, and Adrian Weller. Beyond distributive fairness in algorithmic decision making: Feature selection for procedurally fair learning. In *AAAI*, 2018.
- [Guyon and Elisseeff, 2003] Isabelle Guyon and André Elisseeff. An introduction to variable and feature selection. *Journal of Machine Learning Research*, 2003.
- [He *et al.*, 2006] Xiaofei He, Deng Cai, and Partha Niyogi. Laplacian score for feature selection. In *NeurIPS*, 2006.
- [Inbarani *et al.*, 2014] H Hannah Inbarani, Ahmad Taher Azar, and G Jothi. Supervised hybrid feature selection based on pso and rough sets for medical diagnosis. *Computer Methods and Programs in Biomedicine*, 2014.
- [Kallus *et al.*, 2019] Nathan Kallus, Xiaojie Mao, and Angela Zhou. Assessing algorithmic fairness with unobserved protected class using data combination. *arXiv preprint arXiv:1906.00285*, 2019.
- [Kleindessner *et al.*, 2019] Matthäus Kleindessner, Samira Samadi, Pranjal Awasthi, and Jamie Morgenstern. Guarantees for spectral clustering with fairness constraints. *arXiv preprint arXiv:1901.08668*, 2019.
- [Li *et al.*, 2012] Zechao Li, Yi Yang, Jing Liu, Xiaofang Zhou, and Hanqing Lu. Unsupervised feature selection using nonnegative spectral analysis. In *AAAI*, 2012.
- [Li *et al.*, 2017a] Jundong Li, Kewei Cheng, Suhang Wang, Fred Morstatter, Robert P Trevino, Jiliang Tang, and Huan Liu. Feature selection: a data perspective. *ACM Computing Surveys*, 2017.
- [Li *et al.*, 2017b] Jundong Li, Jiliang Tang, and Huan Liu. Reconstruction-based unsupervised feature selection: An embedded approach. In *IJCAI*, 2017.
- [Li *et al.*, 2018] Jundong Li, Liang Wu, Harsh Dani, and Huan Liu. Unsupervised personalized feature selection. In *AAAI*, 2018.
- [Li *et al.*, 2020] Peizhao Li, Han Zhao, and Hongfu Liu. Deep fair clustering for visual learning. In *CVPR*, 2020.
- [Liang *et al.*, 2015] Deron Liang, Chih-Fong Tsai, and Hsin-Ting Wu. The effect of feature selection on financial distress prediction. *Knowledge-Based Systems*, 2015.
- [Louizos *et al.*, 2015] Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. The variational fair autoencoder. *arXiv preprint arXiv:1511.00830*, 2015.
- [Mehrabi *et al.*, 2019] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *arXiv preprint arXiv:1908.09635*, 2019.
- [Moyer *et al.*, 2018] Daniel Moyer, Shuyang Gao, Rob Brekelmans, Aram Galstyan, and Greg Ver Steeg. Invariant representations without adversarial training. In *NeurIPS*, 2018.
- [Schmidt *et al.*, 2018] Melanie Schmidt, Chris Schwiegelshohn, and Christian Sohler. Fair coresets and streaming algorithms for fair k-means clustering. *arXiv preprint arXiv:1812.10854*, 2018.
- [Sobnath *et al.*, 2020] Drishty Sobnath, Tobiasz Kaduk, Ikram Ur Rehman, and Olufemi Isiaq. Feature selection for uk disabled students’ engagement post higher education: a machine learning approach for a predictive employment model. *IEEE Access*, 2020.
- [Wei *et al.*, 2016] Xiaokai Wei, Bokai Cao, and Philip S Yu. Nonlinear joint unsupervised feature selection. In *SDM*, 2016.
- [Yang *et al.*, 2011] Yi Yang, Heng Tao Shen, Zhigang Ma, Zi Huang, and Xiaofang Zhou.  $\ell_{2,1}$ -norm regularized discriminative feature selection for unsupervised learning. In *IJCAI*, 2011.
- [Zhang *et al.*, 2016] Lu Zhang, Yongkai Wu, and Xintao Wu. A causal framework for discovering and removing direct and indirect discrimination. *arXiv preprint arXiv:1611.07509*, 2016.
- [Zhao and Liu, 2007] Zheng Zhao and Huan Liu. Spectral feature selection for supervised and unsupervised learning. In *ICML*, 2007.