
Fair Sparse Regression with Clustering: An Invox Relaxation for a Combinatorial Problem

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we study the problem of fair sparse regression on a biased dataset
2 where bias depends upon a hidden binary attribute. The presence of a hidden
3 attribute adds an extra layer of complexity to the problem by combining sparse
4 regression and clustering with unknown binary labels. The corresponding opti-
5 mization problem is combinatorial, but we propose a novel relaxation of it as an
6 *invex* optimization problem. To the best of our knowledge, this is the first *invex*
7 relaxation for a combinatorial problem. We show that the inclusion of the debi-
8 asing/fairness constraint in our model has no adverse effect on the performance.
9 Rather, it enables the recovery of the hidden attribute. The support of our re-
10 covered regression parameter vector matches exactly with the true parameter vec-
11 tor. Moreover, we simultaneously solve the clustering problem by recovering the
12 exact value of the hidden attribute for each sample. Our method uses carefully
13 constructed primal dual witnesses to provide theoretical guarantees for the combi-
14 natorial problem. To that end, we show that the sample complexity of our method
15 is logarithmic in terms of the dimension of the regression parameter vector.

16 1 Introduction

17 In modern times, machine learning algorithms are used in a wide variety of applications, many of
18 which are decision making processes such as hiring (Hoffman et al., 2018), predicting human behav-
19 ior (Subrahmanian & Kumar, 2017), COMPAS (Correctional Offender Management Profiling for
20 Alternative Sanctions) risk assessment (Brennan et al., 2009), among others. These decisions have
21 large impacts on society (Kleinberg et al., 2018). Consequently, researchers have shown interest
22 in developing methods that can mitigate unfair decisions and avoid bias amplification. Several fair
23 algorithms have been proposed for machine learning problems such as regression (Agarwal, 2019;
24 Berk, 2017; Calders, 2013), classification (Agarwal et al., 2018; Donini et al., 2018; Dwork et al.,
25 2012; Feldman et al., 2015; Hardt et al., 2016; Huang & Vishnoi, 2019; Pedreshi et al., 2008; Zafar
26 et al., 2019; Zemel et al., 2013) and clustering (Backurs et al., 2019; Bera et al., 2019; Chen et al.,
27 2019; Chierichetti et al., 2017; Huang et al., 2019). A common thread in the above literature is that
28 performance is only viewed in terms of risks, e.g., misclassification rate, false positive rate, false
29 negative rate, mean squared error.

30 In the literature, fairness is discussed in the context of discrimination based on membership to a
31 particular group (e.g. race, religion, gender) which is considered a sensitive attribute. Fairness is
32 generally modeled explicitly by adding a fairness constraint or implicitly by incorporating it in the
33 model itself. There have been several notions of fairness studied in linear regression. Berk (2017)
34 proposed notions of individual fairness and group fairness, and modeled them as penalty functions.
35 Calders (2013) proposed the fairness notions of equal means and balanced residuals by modeling
36 them as explicit constraints. Agarwal (2019), Fitzsimons (2019) and Chzhen et al. (2020) studied

Table 1: Comparison to prior work. Notation: s is the number of non-zero entries in the regression parameter vector and d is its dimension. The terms independent of s and d are not shown in the order notation.

Paper	Hidden sensitive attribute	Modeling type	Sample complexity
Calders (2013); Agarwal (2019); Fitzsimons (2019)	No	Explicit constraint	Not provided
Berk (2017)	No	Penalty function	Not provided
Chzhen et al. (2020)	No	Implicit	Not provided
Our paper	Yes	Implicit	$\Omega(s^3 \log d)$

37 demographic parity. While Agarwal (2019), Fitzsimons (2019) modeled it as an explicit constraint,
 38 Chzhen et al. (2020) included it implicitly in their proposed model.

39 All the above work assume access to the sensitive attribute in the training samples and provide a
 40 framework which are inherently fair. Our work fundamentally differs from these work as we do
 41 not assume access to the sensitive attribute. Without knowing the sensitive attribute, it becomes
 42 difficult to ascertain bias, even for linear regression. In this work, we focus on identifying unfairly
 43 treated members/samples. This adds an extra layer of complexity to linear regression. We solve the
 44 linear regression problem while simultaneously solving a clustering problem where we identify two
 45 clusters – one which is positively biased and the other which is negatively biased. Table 1 shows a
 46 consolidated comparison of our work with the existing literature.

47 Once one identifies bias (positive or negative) for each sample, one could perform debiasing which
 48 would lead to the fairness notion of equal means (Calders, 2013) among the two groups (See Figure
 49 1). It should be noted that identifying groups with positive or negative bias may not be same as
 50 identifying the sensitive attribute. The reason is that there may be multiple attributes that are highly
 51 correlated with the sensitive attribute. In such a situation, these correlated attributes can facilitate
 52 indirect discrimination even if the sensitive attribute is identified and removed. This is called the
 red-lining effect (Calders, 2010). Our model avoids this by directly identifying biased groups.

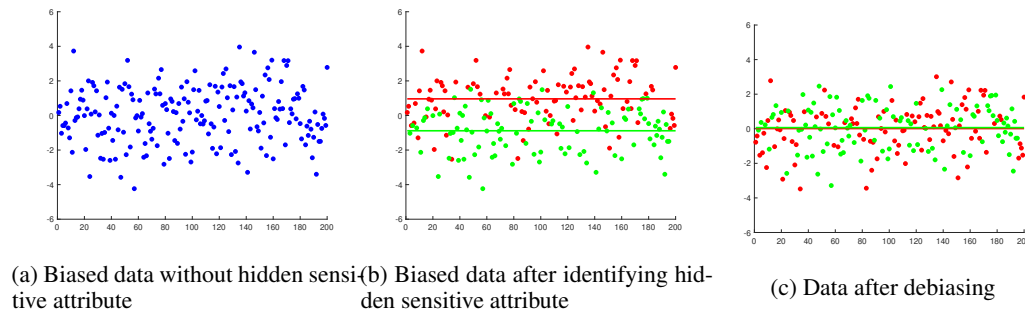


Figure 1: Data before debiasing and after debiasing. Notice how means for two groups (shown as horizontal lines) become almost equal after debiasing.

53
 54 While the standard algorithms solving the sparse/LASSO problem in this setting do provide an esti-
 55 mate of the regression parameter vector, they do not fit the model accurately as they fail to consider
 56 any fairness criteria in their formulations. It is natural then to think about including the hidden at-
 57 tribute in LASSO itself. However, this breaks the convexity of the loss function which makes the
 58 problem intractable by the standard LASSO algorithms. The resulting problem is a combinatorial
 59 version of sparse linear regression with added clustering according to the hidden attribute. In this
 60 work, we propose a novel technique to tackle the combinatorial LASSO problem with a hidden at-
 61 tribute and provide theoretical guarantees about the quality of the solution given a sufficient number
 62 of samples. Our method provably detects unfairness in the system. It should be noted that observing
 63 unfairness does not always imply that the designer of the system intended for such inequalities to
 64 arise. In such cases, our method acts as a check to detect and remove such unintended discrimination.
 65 While the current belief is that there is a trade-off between fairness and performance (Corbett-Davies
 66 et al., 2017; Kleinberg et al., 2017; Pleiss et al., 2017; Zliobaite, 2015; Zhao & Gordon, 2019), our

67 theoretical and experimental results show evidence on the contrary. Our theoretical results allow for
 68 a new understanding of fairness, as an “enabler” instead of as a “constraint”.

69 **Contribution.** Broadly, we can categorize our contribution in the following points:

- 70 • **Defining the problem:** We formulate a novel combinatorial version of sparse linear regres-
 71 sion which takes fairness/bias into the consideration. The addition of clustering comes at
 72 no extra cost in terms of the performance.
- 73 • **Inex relaxation:** Most of the current methods solve convex optimization problems as it
 74 makes the solution tractable. We propose a novel relaxation of the combinatorial problem
 75 and formally show that it is inex. To the best of our knowledge, this is the first inex
 76 relaxation for a combinatorial problem.
- 77 • **Theoretical Guarantees:** Our method can detect bias in the system. In particular, our
 78 method recovers the exact hidden attributes for each sample and thus provides an exact
 79 measure of bias between two different groups. Our method solves linear regression and
 80 clustering simultaneously with theoretical guarantees. To that end, we recover the true
 81 clusters (hidden attributes) and a regression parameter vector which is correct up to the
 82 sign of entries with respect to the true parameter vector. On a more technical side, we
 83 provide a primal-dual witness construction for our inex problem and provide theoretical
 84 guarantees for recovery. The sample complexity of our method varies logarithmically with
 85 respect to dimension of the regression parameter vector.

86 2 Notation and Problem Definition

87 In this section, we collect all the notations used throughout the paper. We also formally introduce
 88 our novel problem. We consider a problem where we have a binary hidden attribute, and where
 89 fairness depends upon the hidden attribute. Let $y \in \mathbb{R}$ be the response variable and $X \in \mathbb{R}^d$ be the
 90 observed attributes. Let $z^* \in \{-1, 1\}$ be the *hidden* attribute and $\gamma \in \mathbb{R}_{>0}$ be the amount of bias due
 91 to the hidden attribute. The response y is generated using the following mechanism:

$$y = X^\top w^* + \gamma z^* + e \quad (1)$$

92 where e is an independent noise term. For example, y could represent the market salary of a new
 93 candidate, X could represent the candidate’s skills and z could represent the population group the
 94 candidate belongs to (e.g., majority or minority). While the group of the candidate is not public
 95 knowledge, a bias associated with the candidate’s group may be present in the underlying data. For
 96 our problem, we will assume that an estimate of the bias $\gamma \in \mathbb{R}_{>0}$ is available. In practice, even a
 97 rough estimate ($\pm 25\%$) of γ also works well (See Appendix J).

98 Let $[d]$ denote the set $\{1, 2, \dots, d\}$. We assume $X \in \mathbb{R}^d$ to be a zero mean sub-Gaussian random
 99 vector (Hsu et al., 2012) with covariance $\Sigma \in \mathbb{S}_+^d$, i.e., there exists a $\rho > 0$, such that for all $\alpha \in \mathbb{R}^d$
 100 the following holds: $\mathbb{E}(\exp(\alpha^\top X)) \leq \exp(\frac{\|\alpha\|_2^2 \rho^2}{2})$. By simply taking $\alpha_i = r$ and $\alpha_k = 0, \forall k \neq i$,
 101 it follows that each entry of X is sub-Gaussian with parameter ρ . In particular, we will assume that
 102 $\forall i \in [d], \frac{X_i}{\sqrt{\Sigma_{ii}}}$ is a sub-Gaussian random variable with parameter $\sigma > 0$. It follows trivially that
 103 $\max_{i \in [d]} \sqrt{\Sigma_{ii}} \sigma \leq \rho$. We will further assume that e is zero mean independent sub-Gaussian noise
 104 with variance σ_e . We assume that as the number of samples increases, the noise in the model gently
 105 decreases. We model this by taking $\sigma_e = \frac{k}{\sqrt{\log n}}$ for some $k > 0$. Our setting works with a variety
 106 of random variables as the class of sub-Gaussian random variable includes for instance Gaussian
 107 variables, any bounded random variable (e.g., Bernoulli, multinomial, uniform), any random vari-
 108 able with strictly log-concave density, and any finite mixture of sub-Gaussian variables. Notice that
 109 for the group with $z = +1$, $\mathbb{E}(y) = \gamma$ and for the group with $z = -1$, $\mathbb{E}(y) = -\gamma$. This means
 110 that after correctly identifying groups, one could perform debiasing by subtracting or adding γ for
 111 $z = +1$ and -1 respectively. After debiasing, the expected value of both groups would match (and
 112 be equal to 0). This complies with the notion of equal mean fairness proposed by Calders (2013).

113 The parameter vector $w^* \in \mathbb{R}^d$ is s -sparse, i.e., at most s entries of w^* are non-zero. We receive n
 114 i.i.d. samples of $X \in \mathbb{R}^d$ and $y \in \mathbb{R}$ and collect them in $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$ respectively. Thus,
 115 in the finite-sample setting,

$$\mathbf{y} = \mathbf{X}w^* + \gamma \mathbf{z}^* + \mathbf{e}, \quad (2)$$

116 where $\mathbf{z}^* \in \{-1, 1\}^n$ and $\mathbf{e} \in \mathbb{R}^n$ both collect n independent realizations of $z^* \in \{-1, 1\}$ and $e \in \mathbb{R}$.
 117 Our goal is to recover w^* and \mathbf{z}^* using the samples (\mathbf{X}, \mathbf{y}) .

118 We denote a matrix $A \in \mathbb{R}^{p \times q}$ restricted to the columns and rows in $P \subseteq [p]$ and $Q \subseteq [q]$ re-
 119 spectively as A_{PQ} . Similarly, a vector $v \in \mathbb{R}^p$ restricted to entries in P is denoted as v_P . We use
 120 $\text{eig}_i(A)$ to denote the i -th eigenvalue (1st being the smallest) of matrix A . Similarly, $\text{eig}_{\max}(A)$
 121 denotes the maximum eigenvalue of matrix A . We use $\text{diag}(A)$ to denote a vector containing
 122 the diagonal element of matrix A . By overriding the same notation, we use $\text{diag}(v)$ to denote a
 123 diagonal matrix with its diagonal being the entries in vector v . We denote the inner product be-
 124 tween two matrices A and B by $\langle A, B \rangle$, i.e., $\langle A, B \rangle = \text{trace}(A^\top B)$, where trace denotes the
 125 trace of a matrix. The notation $A \geq B$ denotes that $A - B$ is a positive semidefinite matrix.
 126 Similarly, $A > B$ denotes that $A - B$ is a positive definite matrix. For vectors, $\|v\|_p$ denotes
 127 the ℓ_p -vector norm of vector $v \in \mathbb{R}^d$, i.e., $\|v\|_p = (\sum_{i=1}^d |v_i|^p)^{\frac{1}{p}}$. If $p = \infty$, then we define
 128 $\|v\|_\infty = \max_{i=1}^d |v_i|$. For matrices, $\|A\|_p$ denotes the induced ℓ_p -matrix norm for matrix $A \in \mathbb{R}^{p \times q}$.
 129 In particular, $\|A\|_2$ denotes the spectral norm of A and $\|A\|_\infty \triangleq \max_{i \in [p]} \sum_{j=1}^q |A_{ij}|$. A function
 130 $f(x)$ is of order $\Omega(g(x))$ and denoted by $f(x) = \Omega(g(x))$, if there exists a constant $C > 0$ such
 131 that for big enough x_0 , $f(x) \geq Cg(x), \forall x \geq x_0$. Similarly, a function $f(x)$ is of order $\mathcal{O}(g(x))$
 132 and denoted by $f(x) = \mathcal{O}(g(x))$, if there exists a constant $C > 0$ such that for big enough x_0 ,
 133 $f(x) \leq Cg(x), \forall x \geq x_0$. For brevity in our notations, we treat any quantity independent of d, s and
 134 n as constant. Detailed proofs for lemmas and theorems are available in the supplementary material.

135 3 Our New Optimization Problem and Inconvexity

136 In this section, we introduce our novel combinatorial problem and propose an invex relaxation. To
 137 the best of our knowledge, this is the first invex relaxation for a combinatorial problem. Without any
 138 information about the hidden attribute \mathbf{z}^* in Equation (2), the following LASSO formulation could
 139 be incorrectly and unsuccessfully used to estimate the parameter w^* .

Definition 1 (Standard LASSO).

$$\min_w \frac{1}{n}(\mathbf{X}w - \mathbf{y})^\top(\mathbf{X}w - \mathbf{y}) + \lambda_n \|w\|_1 \quad (3)$$

140 However, without including \mathbf{z}^* , standard LASSO does not provide accurate estimation of w^* in
 141 Equation (2). We provide the following novel formulation of LASSO which fits our goals of esti-
 142 mating both w^* and \mathbf{z}^* :

Definition 2 (Combinatorial Fair LASSO).

$$\min_{w, \mathbf{z}} \frac{1}{n}(\mathbf{X}w + \gamma \mathbf{z} - \mathbf{y})^\top(\mathbf{X}w + \gamma \mathbf{z} - \mathbf{y}) + \lambda_n \|w\|_1, \quad \text{such that } \mathbf{z}_i \in \{-1, 1\}, \forall i \in [n], \quad (4)$$

143 where $\lambda_n > 0$ is the regularization level which depends on n .

144 In its current form, optimization problem (4) is a non-convex mixed integer quadratic program
 145 (MIQP). Solving MIQP is NP-hard (See Appendix B). Next, we will provide a continuous but still
 146 non-convex relaxation of (4). For ease of notation, we define the following quantities:

$$l(w) \triangleq \frac{1}{n}(\mathbf{X}w - \mathbf{y})^\top(\mathbf{X}w - \mathbf{y}), \quad \mathbf{Z} \triangleq \begin{bmatrix} 1 & \mathbf{z}^\top \\ \mathbf{z} & \mathbf{z}\mathbf{z}^\top \end{bmatrix}, \quad \mathbf{M}(w) \triangleq \begin{bmatrix} l(w) & \frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})^\top \\ \frac{\gamma}{n}(\mathbf{X}w - \mathbf{y}) & \frac{\gamma^2}{n} \mathbf{I}_{n \times n} \end{bmatrix}, \quad (5)$$

147 where \mathbf{I} is an $n \times n$ identity matrix. We provide the following invex relaxation to the optimization
 148 problem (4).

Definition 3 (Invex Fair LASSO).

$$\min_{w, \mathbf{z}} \langle \mathbf{M}(w), \mathbf{Z} \rangle + \lambda_n \|w\|_1, \quad \text{such that } \text{diag}(\mathbf{Z}) = \mathbf{1}, \mathbf{Z} \geq \mathbf{0}_{n+1 \times n+1}. \quad (6)$$

149 Note that optimization problem (6) is continuous and convex with respect to w and \mathbf{Z} separately but
 150 it is not jointly convex (See Appendix N for details). Specifically, for a fixed w , the matrix $\mathbf{M}(w)$
 151 becomes a constant and problem (6) resembles a semidefinite program. For a fixed \mathbf{Z} , problem (6)
 152 resembles a standard LASSO. Unfortunately, problem (6) is not jointly convex on w and \mathbf{Z} , and
 153 thus, it might still remain difficult to solve. Next, we will provide arguments that despite being
 154 non-convex, optimization problem (6) belongs to a particular class of non-convex functions namely
 155 ‘‘invex’’ functions. We define ‘‘invexity’’ of functions, as a generalization of convexity (Hanson,
 156 1981).

157 **Definition 4** (Invex function). Let $\phi(t)$ be a function defined on a set C . Let η be a vector valued
158 function defined in $C \times C$ such that $\eta(t_1, t_2)^\top \nabla \phi(t_2)$, is well defined $\forall t_1, t_2 \in C$. Then, $\phi(t)$ is a
159 η -invex function if $\phi(t_1) - \phi(t_2) \geq \eta(t_1, t_2)^\top \nabla \phi(t_2)$, $\forall t_1, t_2 \in C$.

160 Note that convex functions are η -invex for $\eta(t_1, t_2) = t_1 - t_2$. Hanson (1981) showed that if the
161 objective function and constraints are both η -invex with respect to same η defined in $C \times C$, then
162 Karush-Kuhn-Tucker (KKT) conditions are sufficient for optimality, while it is well-known that
163 KKT conditions are necessary. Ben-Israel & Mond (1986) showed a function is invex if and only if
164 each of its stationarity point is a global minimum.

165 In the next lemma, we show that the relaxed optimization problem (6) is indeed η -invex for a par-
166 ticular η defined in $C \times C$ and a well defined set C . Before that, we will reformulate it into an
167 equivalent optimization problem. Note that in the optimization problem (6), $\text{diag}(\mathbf{Z}) = \mathbf{1}$. Thus,
168 $\langle \mathbf{I}, \mathbf{Z} \rangle$ is a constant equal to $n + 1$. Using this, we can rewrite the optimization problem as:

$$\min_{w, \mathbf{Z}} \langle \mathbf{M}(w), \mathbf{Z} \rangle + \lambda_n \|w\|_1 + \langle \mathbf{I}, \mathbf{Z} \rangle, \quad \text{such that} \quad \text{diag}(\mathbf{Z}) = \mathbf{1}, \mathbf{Z} \geq \mathbf{0}_{n+1 \times n+1}, \quad (7)$$

169 Let $C = \{(w, \mathbf{Z}) \mid w \in \mathbb{R}^d, \text{diag}(\mathbf{Z}) = \mathbf{1}, \mathbf{Z} \geq \mathbf{0}_{n+1 \times n+1}\}$. We take $\mathbf{M}'(w) = \mathbf{M}(w) + \mathbf{I}$ and the
170 corresponding optimization problem becomes: $\min_{(w, \mathbf{Z}) \in C} \langle \mathbf{M}'(w), \mathbf{Z} \rangle + \lambda_n \|w\|_1$. We will show
171 that $\forall (w, \mathbf{Z}) \in C$, $\langle \mathbf{M}'(w), \mathbf{Z} \rangle + \lambda_n \|w\|_1$ is an invex function. Note that by definition of the ℓ_1 -norm,
172 $\|w\|_1 = \sup_{\|a\|_\infty=1} \langle a, w \rangle$. Thus, it suffices to show that $\forall a \in \mathbb{R}^d$, $\langle \mathbf{M}'(w), \mathbf{Z} \rangle$ and $\langle a, w \rangle$ are invex
173 for the same $\eta(w, \bar{w}, \mathbf{Z}, \bar{\mathbf{Z}})$.

174 **Lemma 1.** For $(w, \mathbf{Z}) \in C$, the functions $f(w, \mathbf{Z}) = \langle \mathbf{M}'(w), \mathbf{Z} \rangle$ and $g(w, \mathbf{Z}) = \langle a, w \rangle$ are η -
175 invex for $\eta(w, \bar{w}, \mathbf{Z}, \bar{\mathbf{Z}}) \triangleq \left[\begin{array}{c} w - \bar{w} \\ \mathbf{M}'(\bar{w})^{-1} \mathbf{M}'(w)(\mathbf{Z} - \bar{\mathbf{Z}}) \end{array} \right]$, where we abuse the vector/matrix notation
176 for clarity of presentation, and avoid the vectorization of matrices.

177 Now that we have established that optimization problem (6) is invex, we are ready to discuss our
178 main results in the next section.

179 4 Our Theoretical Analysis

180 In this section, we show that our Invex Fair Lasso formulation correctly recovers the hidden attributes
181 and the regression parameter vector. More formally, we want to achieve the two goals by solving
182 optimization problem (6) efficiently. First, we want to correctly and uniquely determine the hidden
183 sensitive attribute for each data point, i.e., $\mathbf{z}^* \in \{-1, 1\}^n$. Second, we want to recover regression
184 parameter vector which is close to the true parameter vector $w^* \in \mathbb{R}^d$ in ℓ_2 -norm. Let \tilde{w} and $\tilde{\mathbf{Z}}$ be
185 the solution to optimization problem (6). Then, we will prove that \tilde{w} and w^* have the same support
186 and $\tilde{\mathbf{z}}$ constructed from $\tilde{\mathbf{Z}}$ is exactly equal to \mathbf{z}^* . We define $\Delta \triangleq (\tilde{w} - w^*)$.

187 4.1 KKT conditions

188 We start by writing the KKT conditions for optimization problem (6). Let $\mu \in \mathbb{R}^{n+1}$ and
189 $\Lambda \geq \mathbf{0}_{n+1 \times n+1}$ be the dual variables for optimization problem (6). For a fixed λ_n , the Lagrangian
190 $L(w, \mathbf{Z}; \mu, \Lambda)$ can be written as $L(w, \mathbf{Z}; \mu, \Lambda) = \langle \mathbf{M}(w), \mathbf{Z} \rangle + \lambda_n \|w\|_1 + \langle \text{diag}(\mu), \mathbf{Z} \rangle - \mathbf{1}^\top \mu -$
191 $\langle \Lambda, \mathbf{Z} \rangle$. Using this Lagrangian, the KKT conditions at the optimum can be written as:

192 1. Stationarity conditions:

$$\frac{\partial \langle \mathbf{M}(w), \mathbf{Z} \rangle}{\partial w} + \lambda_n \mathbf{g} = \mathbf{0}_{d \times 1}, \quad (8)$$

193 where \mathbf{g} is an element of the subgradient set of $\|w\|_1$, i.e., $\mathbf{g} \in \frac{\partial \|w\|_1}{\partial w}$ and $\|\mathbf{g}\|_\infty \leq 1$.

$$\mathbf{M}(w) + \text{diag}(\mu) - \Lambda = \mathbf{0}_{n+1 \times n+1} \quad (9)$$

194 2. Complementary Slackness condition:

$$\langle \Lambda, \mathbf{Z} \rangle = 0 \quad (10)$$

195 3. Dual Feasibility condition:

$$\Lambda \geq \mathbf{0}_{n+1 \times n+1} \quad (11)$$

196 4. Primal Feasibility conditions:

$$w \in \mathbb{R}^d, \text{diag}(\mathbf{Z}) = \mathbf{1}, \mathbf{Z} \geq \mathbf{0}_{n+1 \times n+1} \quad (12)$$

197 Next, we will provide a setting for primal and dual variables which satisfies all the KKT conditions.
198 But before that, we will describe a set of technical assumptions which will help us in our analysis.

199 4.2 Assumptions

200 Let S denote the support of w^* , i.e., $S = \{i \mid w_i^* \neq 0, i \in [d]\}$. Similarly, we define the complement
201 of support S as $S^c = \{i \mid w_i^* = 0, i \in [d]\}$. Let $|S| = s$ and $|S^c| = d - s$. For ease of notation, we
202 define $\mathbf{H} \triangleq \mathbb{E}(X X^\top)$ and $\hat{\mathbf{H}} \triangleq \frac{1}{n} \mathbf{X}^\top \mathbf{X}$. As the first assumption, we need the minimum eigenvalue
203 of the population covariance matrix of X restricted to rows and columns in S to be greater than
204 zero. Later, we will show that this assumption is needed to uniquely recover w in the optimization
205 problem (6).

206 **Assumption 1** (Positive Definiteness of Hessian). $\mathbf{H}_{SS} > \mathbf{0}_{s \times s}$ or equivalently $\text{eig}_{\min}(\mathbf{H}_{SS}) =$
207 $C_{\min} > 0$.

208 In practice, we only deal with finite samples and not populations. In the next lemma, we will
209 show that with a sufficient number of samples, a condition similar to Assumption 1 holds with high
210 probability in the finite-sample setting.

211 **Lemma 2.** *If Assumption 1 holds and $n = \Omega(\frac{s+\log d}{C_{\min}^2})$, then $\text{eig}_{\min}(\hat{\mathbf{H}}_{SS}) \geq \frac{C_{\min}}{2}$ with probability*
212 *at least $1 - \mathcal{O}(\frac{1}{d})$.*

213 As the second assumption, we will need to ensure that the variates outside the support of w^* do not
214 exert lot of influence on the variates in the support of w^* . This sort of technical condition, known as
215 the mutual incoherence condition, has been previously used in many problems related to regularized
216 regression such as compressed sensing (Wainwright, 2009), Markov random fields (Ravikumar et al.,
217 2010), non-parametric regression (Ravikumar et al., 2007), diffusion networks (Daneshmand et al.,
218 2014), among others. We formally present this technical condition in what follows.

219 **Assumption 2** (Mutual Incoherence). $\|\mathbf{H}_{S^c S} \mathbf{H}_{SS}^{-1}\|_\infty \leq 1 - \alpha$ for some $\alpha \in (0, 1]$.

220 Again, we will show that with a sufficient number of samples, a condition similar to Assumption 2
221 holds in the finite-sample setting with high probability.

222 **Lemma 3.** *If Assumption 2 holds and $n = \Omega(\frac{s^3(\log s + \log d)}{\tau(C_{\min}, \alpha, \sigma, \Sigma)})$, then $\|\hat{\mathbf{H}}_{S^c S} \hat{\mathbf{H}}_{SS}^{-1}\|_\infty \leq 1 - \frac{\alpha}{2}$ with*
223 *probability at least $1 - \mathcal{O}(\frac{1}{d})$ where $\tau(C_{\min}, \alpha, \sigma, \Sigma)$ is a constant independent of n, d and s .*

224 In Appendix L, we experimentally show that Assumption 1 is easier to hold (i.e., $n \in \Omega(s + \log d)$)
225 than Assumption 2 (i.e., $n \in \Omega(s^3 \log d)$). Eventually, both assumptions hold as the number of
226 samples increases.

227 4.3 Construction of Primal and Dual Witnesses

228 In this subsection, we will provide a construction of primal and dual variables which satisfies the
229 KKT conditions for optimization problem (6). To that end, we provide our first main result.

230 **Theorem 1** (Primal Dual Witness Construction). *If Assumptions 1 and 2 hold, $\lambda_n \geq \frac{128\rho k \sqrt{\log d}}{\alpha n}$*
231 *and $n = \Omega(\frac{s^3 \log d}{\tau_0(C_{\min}, \alpha, \sigma, \Sigma, \rho, k, \gamma)})$, then the following setting of primal and dual variables*

$$\begin{aligned} \text{Primal Variables: } \tilde{w} &= (\tilde{w}_S, \mathbf{0}_{d-s \times 1}) \\ \text{where, } \tilde{w}_S &= \arg \min_{w_S} \frac{1}{n} (\mathbf{X}_{\cdot, S} w_S + \gamma \mathbf{z}^* - \mathbf{y})^\top (\mathbf{X}_{\cdot, S} w_S + \gamma \mathbf{z}^* - \mathbf{y}) + \lambda_n \|w_S\|_1 \\ \mathbf{Z} &= \mathbf{Z}^* \triangleq \begin{bmatrix} \mathbf{1} & \mathbf{z}^{*\top} \\ \mathbf{z}^* & \mathbf{z}^* \mathbf{z}^{*\top} \end{bmatrix} \\ \text{Dual Variables: } \mu &= -\text{diag}(M(\tilde{w}) \mathbf{Z}^*), \quad \Lambda = M(\tilde{w}) - \text{diag}(M(\tilde{w}) \mathbf{Z}^*) \end{aligned} \quad (13)$$

232 satisfies all the KKT conditions for optimization problem (6) with probability at least $1 - \mathcal{O}(\frac{1}{n})$,
 233 where $\tau_0(C_{\min}, \alpha, \sigma, \Sigma, \rho, k, \gamma)$ is a constant independent of s, d and n and thus, the primal vari-
 234 ables are a globally optimal solution for (6). Furthermore, the above solution is also unique.

235 **Proof Sketch.** The main idea behind our proofs is to verify that the setting of primal and dual
 236 variables in Theorem 1 satisfies all the KKT conditions described in subsection 4.1. We do this by
 237 proving multiple lemmas in subsequent subsections. The outline of the proof is as follows:

- 238 • It can be trivially verified that the primal feasibility condition (12) holds. Similarly, the
 239 second stationarity condition (9) holds by construction of Λ .
- 240 • In subsection 4.4, we use Lemmas 4 and 11 to verify that the stationarity condition (8)
 241 holds.
- 242 • In subsection 4.5, we use Lemma 5 to verify the complementary slackness condition (10).
- 243 • In subsection 4.6, we show that the dual feasibility condition (11) is satisfied using results
 244 from Lemmas 6, 7, 8 and 12.
- 245 • Finally, in subsection 4.7, we show that our proposed solution is also unique.

246 4.4 Verifying the Stationarity Condition (8)

247 In this subsection, we will show that the setting of \tilde{w} and \mathbf{Z}^* satisfies the first stationarity condition
 248 (8) by proving the following lemma.

249 **Lemma 4.** *If Assumptions 1 and 2 hold, $\lambda_n \geq \frac{128\rho k \sqrt{\log d}}{\alpha n}$ and $n = \Omega(\frac{s^3 \log d}{\tau_1(C_{\min}, \alpha, \sigma, \Sigma, \rho)})$, then the
 250 setting of w and \mathbf{Z} from equation (13) satisfies the stationarity condition (8) with probability at least
 251 $1 - \mathcal{O}(\frac{1}{d})$, where $\tau_1(C_{\min}, \alpha, \sigma, \Sigma, \rho)$ is a constant independent of d, s or n .*

252 4.5 Verifying the Complementary Slackness (10)

253 Next, we will show that the setting of Λ and \mathbf{Z} in (13) satisfies the complementary slackness condi-
 254 tion (10). To this end, we will show the following:

255 **Lemma 5.** *Let Λ be defined as in equation (13), then $\zeta^* \triangleq \begin{bmatrix} 1 \\ \mathbf{z}^* \end{bmatrix}$ is an eigenvector of Λ correspond-
 256 ing to the eigenvalue 0. Furthermore, $\langle \Lambda, \mathbf{Z}^* \rangle = 0$.*

257 *Proof.* We will show that $\Lambda \zeta^* = \mathbf{0}_{n+1 \times 1}$. Note that,

$$\Lambda = M(w) - \text{diag}(M(w)\mathbf{Z}^*) = \begin{bmatrix} -\frac{\gamma}{n}(Xw - y)^T \mathbf{z}^* & \frac{\gamma}{n}(Xw - y)^T \\ \frac{\gamma}{n}(Xw - y) & -\text{diag}(\frac{\gamma}{n}(Xw - y)\mathbf{z}^{*T}) \end{bmatrix}$$

258 Multiplying the above matrix with ζ^* gives us $\mathbf{0}_{n+1 \times 1}$. Now $\langle \Lambda, \mathbf{Z}^* \rangle = \text{trace}(\Lambda^T \mathbf{Z}^*) =$
 259 $\text{trace}(\Lambda \zeta^* \zeta^{*T}) = 0$ as $\Lambda \zeta^* = \mathbf{0}_{n+1 \times 1}$. \square

260 4.6 Verifying the Dual Feasibility (11)

261 We have already shown that Λ has 0 as one of its eigenvalues. To verify that it satisfies the dual
 262 feasibility condition (11), we show that second minimum eigenvalue of Λ is greater than zero with
 263 high probability. At this point, it might not be clear why strict positivity is necessary, but this will
 264 be argued later in subsection 4.7. Now, note that:

$$\mathbb{P}(\text{eig}_2(\Lambda) > 0) \geq \mathbb{P}(\text{eig}_2(\Lambda) > 0, \|\Delta\|_2 \leq h(n)) \geq \mathbb{P}(\text{eig}_2(\Lambda) > 0 \|\Delta\|_2 \leq h(n)) \mathbb{P}(\|\Delta\|_2 \leq h(n)) \quad (14)$$

265 where $h(n)$ is a function of n . We bound $\mathbb{P}(\text{eig}_2(\Lambda) > 0)$ in two parts. First, we bound $\mathbb{P}(\text{eig}_2(\Lambda) >$
 266 $0)$ given that $\|\Delta\|_2 \leq h(n)$ and then we bound the probability of $\|\Delta\|_2 \leq h(n)$.

267 **Lemma 6.** *Given that $\|\Delta\|_2 \leq h(n)$, the second minimum eigenvalue of Λ as defined in equation
 268 (13) is strictly greater than 0 with probability at least $1 - \exp(-\frac{\gamma^2}{8(\rho^2 h(n)^2 + \sigma^2)}) + \log n$.*

269 *Proof.* As the first step, we invoke Haynesworth's inertia additivity formula (Haynsworth, 1968) to
 270 prove our claim. Let R be a block matrix of the form $R = \begin{bmatrix} A & B \\ B^T & C \end{bmatrix}$, then inertia of matrix R ,

271 denoted by $\text{In}(R)$, is defined as the tuple $(\pi(R), \nu(R), \delta(R))$ where $\pi(R)$ is the number of positive
 272 eigenvalues, $\nu(R)$ is the number of negative eigenvalues and $\delta(R)$ is the number of zero eigenvalues
 273 of matrix R . Haynesworth's inertia additivity formula is given as:

$$\text{In}(R) = \text{In}(C) + \text{In}(A - B^\top C^{-1} B) \quad (15)$$

274 Note that,

$$\Lambda = \mathbf{M}(w) - \text{diag}(\mathbf{M}(w)\mathbf{Z}^*) = \begin{bmatrix} -\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})^\top \mathbf{z}^* & \frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})^\top \\ \frac{\gamma}{n}(\mathbf{X}w - \mathbf{y}) & -\text{diag}(\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})\mathbf{z}^{*\top}) \end{bmatrix}$$

275 Then the following holds true by applying equation (15):

$$\text{In}(\Lambda) = \text{In}(-\text{diag}(\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})\mathbf{z}^{*\top})) + \text{In}(-\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})^\top \mathbf{z}^* - \frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})^\top (-\text{diag}(\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})\mathbf{z}^{*\top})^{-1} \frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})))$$

276 Notice that the term $-\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})^\top \mathbf{z}^* - \frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})^\top (-\text{diag}(\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})\mathbf{z}^{*\top})^{-1} \frac{\gamma}{n}(\mathbf{X}w - \mathbf{y}))$
 277 evaluates to 0. Thus, it has 0 positive eigenvalue, 0 negative eigenvalue and 1 zero eigenvalue. We
 278 have also shown in Lemma 5 that Λ has at least 1 zero eigenvalue. It follows that

$$\begin{aligned} \pi(\Lambda) &= \pi(-\text{diag}(\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})\mathbf{z}^{*\top})), \quad \nu(\Lambda) = \nu(-\text{diag}(\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})\mathbf{z}^{*\top})) \\ \delta(\Lambda) &= \delta(-\text{diag}(\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})\mathbf{z}^{*\top})) + 1 \end{aligned} \quad (16)$$

279 Next, we will show that $-\text{diag}(\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})\mathbf{z}^{*\top})$ has all of its eigenvalues being positive.

280 **Lemma 7.** For a given $\|\Delta\|_2 \leq h(n)$, all eigenvalues of $-\text{diag}(\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})\mathbf{z}^{*\top})$ are strictly greater
 281 than 0 with probability at least $1 - \exp(-\frac{\gamma^2}{8(\rho^2 h(n)^2 + \sigma_e^2)}) + \log n$.

282 *Proof.* Using equation (1), we can expand the term $-\text{diag}(\frac{\gamma}{n}(\mathbf{X}w - \mathbf{y})\mathbf{z}^{*\top})$ as $-\text{diag}(\frac{\gamma}{n}(\mathbf{X}(w -$
 283 $w^*) - \gamma\mathbf{z}^* - \mathbf{e})\mathbf{z}^{*\top})$. Since eigenvalues of a diagonal matrix are its diagonal elements, we focus on
 284 the i -th diagonal element of $-\text{diag}(\frac{\gamma}{n}(\mathbf{X}\Delta - \gamma\mathbf{z}^* - \mathbf{e})\mathbf{z}^{*\top})$ which is $\frac{\gamma^2}{n} - \frac{\gamma}{n}z_i^*(\mathbf{X}_i^\top \Delta + \mathbf{e}_i)$. Note
 285 that $(\mathbf{X}_i^\top \Delta + \mathbf{e}_i)$ is a sub-Gaussian random variable with parameter $\rho^2 \|\Delta\|_2^2 + \sigma_e^2$. Using the tail
 286 inequality for sub-Gaussian random variables, for some $t > 0$, we can write:

$$\mathbb{P}((\mathbf{X}_i \Delta + \mathbf{e}_i) \geq t) \leq \exp(-\frac{t^2}{2(\rho^2 \|\Delta\|_2^2 + \sigma_e^2)})$$

287 We take union bound across all the diagonal elements and replace $t = \frac{\gamma}{2}$ and $\|\Delta\|_2 \leq h(n)$ to
 288 complete the proof, i.e.,

$$\mathbb{P}(\exists i \in [n], (\mathbf{X}_i \Delta + \mathbf{e}_i) \geq t) \leq n \exp(-\frac{\gamma^2}{8(\rho^2 h(n)^2 + \sigma_e^2)}). \quad (17)$$

289 □

290 The result of Lemma 6 follows directly from Lemma 7 and equation (16). □

291 Now, we are ready to bound $\|\Delta\|_2$. Due to our primal dual construction, $\|\Delta\|_2$ is simply equal to
 292 $\|\Delta_S\|_2$. We provide a bound on Δ_S in the following lemma:

293 **Lemma 8.** If Assumptions 1 and 2 hold, $\lambda_n \geq \frac{128\rho k\sqrt{\log d}}{cn}$ and $n = \Omega(\frac{s^3 \log d}{\tau_2(C_{\min}, \rho, k)})$, then $\|\Delta_S\|_2 \leq$
 294 $\frac{2\lambda_n \sqrt{s}}{C_{\min}}$ with probability at least $1 - \mathcal{O}(\frac{1}{d})$ where $\tau_2(C_{\min}, \rho, k)$ is a constant independent of s, d or
 295 n .

296 By taking $h(n) = \frac{2\lambda_n \sqrt{s}}{C_{\min}}$ in (14), we get the following: $\mathbb{P}(\text{eig}_2(\Lambda) > 0) \geq 1 - \mathcal{O}(\frac{1}{n})$, as long as
 297 $n = \Omega(\frac{s^3 \log d}{\tau_3(C_{\min}, \rho, k, \alpha, \gamma)})$, where $\tau_3(C_{\min}, \rho, k, \alpha, \gamma)$ is a constant independent of s, d and n . The
 298 above results combined with the property that optimization problem (6) is invex ensure that the
 299 setting of primal and dual variables in Theorem 1 is indeed the globally optimal solution to the
 300 problem (6). It remains to show that this solution is also unique.

301 4.7 Uniqueness of the Solution

302 First, we prove that \mathbf{Z}^* is a unique solution. Suppose there is another solution $\bar{\mathbf{Z}}$ which satisfies all
 303 KKT conditions and is optimal. Then, $\bar{\mathbf{Z}} \geq \mathbf{0}_{n+1 \times n+1}$ and $\langle \Lambda, \bar{\mathbf{Z}} \rangle = 0$. Since, $\Lambda \geq \mathbf{0}_{n+1 \times n+1}$ and
 304 $\text{eig}_2(\Lambda) > 0$, ζ^* spans all of its null space. This enforces that $\bar{\mathbf{Z}}$ is a multiple of \mathbf{Z}^* . But primal
 305 feasibility dictates that $\text{diag}(\bar{\mathbf{Z}}) = \mathbf{1}$. It follows that $\bar{\mathbf{Z}} = \mathbf{Z}^*$. To show that \tilde{w} is unique, it suffices
 306 to show that \tilde{w}_S is unique. After substituting $\mathbf{Z} = \mathbf{Z}^*$, we observe that the Hessian of optimization
 307 problem (6) with respect to w and restricted to rows and columns in S , i.e., $\hat{\mathbf{H}}_{SS}$ is positive definite.
 308 This ensures that \tilde{w} is a unique solution.

309 The setting of primal and dual variables in Theorem 1 not only solves the optimization problem (6)
 310 but also gives rise to the following results:

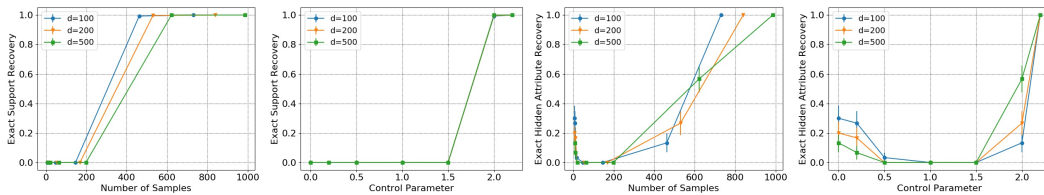
311 **Corollary 1.** *If Assumptions 1 and 2 hold, $\lambda_n \geq \frac{128\rho k}{\alpha} \frac{\sqrt{\log d}}{n}$ and $n = \Omega\left(\frac{s^3 \log d}{\tau_1(C_{\min}, \alpha, \sigma, \Sigma, \rho)}\right)$, then the*
 312 *following statements are true with probability at least $1 - \mathcal{O}\left(\frac{1}{n}\right)$:*

- 313 1. *The solution \mathbf{Z} correctly recovers hidden attribute for each sample, i.e., $\mathbf{Z} = \mathbf{Z}^* = \zeta^* \zeta^{*\top}$.*
- 314 2. *The support of recovered regression parameter \tilde{w} matches exactly with the support of w^* .*
- 315 3. *If $\min_{i \in S} |w_i^*| \geq \frac{4\lambda_n \sqrt{s}}{C_{\min}}$ then for all $i \in [d]$, \tilde{w}_i and w_i^* match up to their sign.*

316 5 Experimental Validation

317 **Synthetic Experiments.** We validate our theoretical result in Theorem 1 and Corollary 1 by
 318 conducting experiments on synthetic data. We show that for a fixed s , we need $n = 10^\beta \log d$
 319 samples for recovering the exact support of w^* and exact hidden attributes \mathbf{Z}^* , where $\beta \equiv$
 320 $\beta(s, C_{\min}, \alpha, \sigma, \Sigma, \rho, \gamma, k)$ is a control parameter which is independent of d . We draw $\mathbf{X} \in \mathbb{R}^{n \times d}$
 321 and $\mathbf{e} \in \mathbb{R}^n$ from Gaussian distributions. We randomly generate $w^* \in \mathbb{R}^d$ with $s = 10$ non-zero
 322 entries. Regarding the hidden attribute $\mathbf{z}^* \in \{-1, 1\}^n$, we set $\frac{n}{2}$ entries as $+1$ and the rest as -1 .
 323 The response $\mathbf{y} \in \mathbb{R}^n$ is generated according to (1). According to Theorem 1, the regularizer λ_n is
 324 chosen to be equal to $\frac{128\rho k}{\alpha} \frac{\sqrt{\log d}}{n}$. We solve optimization problem (6) by using an alternate opti-
 325 mization algorithm that converges to the optimal solution (See Appendix K for details). Figure 2a
 326 shows that our method recovers the true support as we increase the number of samples. Similarly,
 327 Figure 2c shows that as the number of samples increase, our recovered hidden attributes are 100%
 328 correct. Curves line up perfectly in Figure 2b and 2d when plotting with respect to the control
 329 parameter $\beta = \log \frac{n}{\log d}$. This validates our theoretical results (Details in Appendix M).

330 **Real World Experiments.** We show applicability of our method by identify groups with bias in
 331 the Communities and Crime data set (Redmond, 2002) and the Student Performance data set (Cortez,
 332 2008). In both cases, our method is able to recover groups with bias (Details in Appendix O).



(a) Recovery of S vs n (b) Recovery of S vs β (c) Recovery of \mathbf{Z}^* vs n (d) Recovery of \mathbf{Z}^* vs β

Figure 2: Left two: Exact support recovery of w^* across 30 runs. Right two: Exact hidden attribute recovery of \mathbf{Z}^* across 30 runs.

333 References

- 334 Agarwal, A., Beygelzimer, A., Dudík, M., Langford, J., and Wallach, H. A reductions approach to
 335 fair classification. *International Conference on Machine Learning*, 2018.
- 336 Agarwal, A., Dudík, M., and Wu, Z. S. Fair regression: Quantitative definitions and reduction-based
 337 algorithms. *International Conference on Machine Learning*, 2019.

- 338 Backurs, A., Indyk, P., Onak, K., Schieber, B., Vakilian, A., and Wagner, T. Scalable fair clustering.
339 In *International Conference on Machine Learning*, pp. 405–413. PMLR, 2019.
- 340 Ben-Israel, A. and Mond, B. What is invexity? *The ANZIAM Journal*, 28(1):1–9, 1986.
- 341 Bera, S. K., Chakrabarty, D., Flores, N. J., and Negahbani, M. Fair algorithms for clustering. *Neural*
342 *Information Processing Systems*, 2019.
- 343 Berk, R., Heidari, H., Jabbari, S., Joseph, M., Kearns, M., Morgenstern, J., Neel, S., and Roth, A. A
344 convex framework for fair regression. *ACM International Conference on Knowledge Discovery*
345 *and Data Mining, Workshop on Fairness, Accountability, and Transparency in Machine Learning*,
346 2017.
- 347 Brennan, T., Dieterich, W., and Ehret, B. Evaluating the predictive validity of the compas risk and
348 needs assessment system. *Criminal Justice and Behavior*, 36(1):21–40, 2009.
- 349 Calders, T., Karim, A., Kamiran, F., Ali, W., and Zhang, X. Controlling attribute effect in linear
350 regression. *2013 IEEE 13th international conference on data mining*, pp. 71–80, 2013.
- 351 Chen, X., Fain, B., Lyu, L., and Munagala, K. Proportionally fair clustering. In *International*
352 *Conference on Machine Learning*, pp. 1032–1041. PMLR, 2019.
- 353 Chierichetti, F., Kumar, R., Lattanzi, S., and Vassilvitskii, S. Fair clustering through fairlets. *Neural*
354 *Information Processing Systems*, 2017.
- 355 Chzhen, E., Denis, C., Hebiri, M., Oneto, L., and Pontil, M. Fair regression with wasserstein
356 barycenters. *Neural Information Processing Systems*, 2020.
- 357 Corbett-Davies, S., Pierson, E., Feller, A., Goel, S., and Huq, A. Algorithmic decision making and
358 the cost of fairness. In *Proceedings of the 23rd acm sigkdd international conference on knowledge*
359 *discovery and data mining*, pp. 797–806, 2017.
- 360 Daneshmand, H., Gomez-Rodriguez, M., Song, L., and Schoelkopf, B. Estimating Diffusion Net-
361 work Structures: Recovery Conditions, Sample Complexity & Soft-Thresholding Algorithm. In
362 *International Conference on Machine Learning*, pp. 793–801, 2014.
- 363 Donini, M., Oneto, L., Ben-David, S., Shawe-Taylor, J. S., and Pontil, M. Empirical risk minimiza-
364 tion under fairness constraints. pp. 2791–2801, 2018.
- 365 Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness through awareness. In *Pro-*
366 *ceedings of the 3rd innovations in theoretical computer science conference*, pp. 214–226, 2012.
- 367 Feldman, M., Friedler, S. A., Moeller, J., Scheidegger, C., and Venkatasubramanian, S. Certifying
368 and removing disparate impact. In *proceedings of the 21th ACM SIGKDD international confer-*
369 *ence on knowledge discovery and data mining*, pp. 259–268, 2015.
- 370 Hanson, M. A. On sufficiency of the kuhn-tucker conditions. *Journal of Mathematical Analysis and*
371 *Applications*, 80(2):545–550, 1981.
- 372 Hardt, M., Price, E., and Srebro, N. Equality of opportunity in supervised learning. In *Advances in*
373 *neural information processing systems*, pp. 3315–3323, 2016.
- 374 Haynsworth, E. V. Determination of the inertia of a partitioned hermitian matrix. *Linear algebra*
375 *and its applications*, 1(1):73–81, 1968.
- 376 Hoffman, M., Kahn, L. B., and Li, D. Discretion in hiring. *The Quarterly Journal of Economics*,
377 133(2):765–800, 2018.
- 378 Horn, R. A. and Johnson, C. R. *Matrix analysis*. Cambridge university press, 2012.
- 379 Hsu, D., Kakade, S., Zhang, T., et al. A tail inequality for quadratic forms of subgaussian random
380 vectors. *Electronic Communications in Probability*, 17, 2012.
- 381 Huang, L. and Vishnoi, N. K. Stable and fair classification. *International Conference on Machine*
382 *Learning*, 2019.

- 383 Huang, L., Jiang, S. H.-C., and Vishnoi, N. K. Coresets for clustering with fairness constraints.
384 *Neural Information Processing Systems*, 2019.
- 385 Kleinberg, J., Mullainathan, S., and Raghavan, M. Inherent trade-offs in the fair determination of
386 risk scores. *Innovations in Theoretical Computer Science*, 2017.
- 387 Kleinberg, J., Lakkaraju, H., Leskovec, J., Ludwig, J., and Mullainathan, S. Human decisions and
388 machine predictions. *The quarterly journal of economics*, 133(1):237–293, 2018.
- 389 Knutson, A. and Tao, T. Honeycombs and sums of hermitian matrices. *Notices Amer. Math. Soc*, 48
390 (2), 2001.
- 391 Pedreshi, D., Ruggieri, S., and Turini, F. Discrimination-aware data mining. In *Proceedings of*
392 *the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp.
393 560–568, 2008.
- 394 Pleiss, G., Raghavan, M., Wu, F., Kleinberg, J., and Weinberger, K. Q. On fairness and calibration.
395 *Neural Information Processing Systems*, 2017.
- 396 Ravikumar, P., Liu, H., Lafferty, J., and Wasserman, L. Spam: Sparse Additive Models. In *Pro-*
397 *ceedings of the 20th International Conference on Neural Information Processing Systems*, pp.
398 1201–1208. Curran Associates Inc., 2007.
- 399 Ravikumar, P., Wainwright, M. J., Lafferty, J. D., et al. High-dimensional izing model selection
400 using ℓ_1 -regularized logistic regression. *The Annals of Statistics*, 38(3):1287–1319, 2010.
- 401 Ravikumar, P., Wainwright, M. J., Raskutti, G., Yu, B., et al. High-dimensional covariance estima-
402 tion by minimizing ℓ_1 -penalized log-determinant divergence. *Electronic Journal of Statistics*, 5:
403 935–980, 2011.
- 404 Subrahmanian, V. and Kumar, S. Predicting human behavior: The next frontiers. *Science*, 355
405 (6324):489–489, 2017.
- 406 Vershynin, R. How close is the sample covariance matrix to the actual covariance matrix? *Journal*
407 *of Theoretical Probability*, 25(3):655–686, 2012.
- 408 Wainwright, M. J. Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -
409 constrained quadratic programming (lasso). *IEEE transactions on information theory*, 55(5):
410 2183–2202, 2009.
- 411 Wainwright, M. J. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge
412 University Press, 2019.
- 413 Zafar, M. B., Valera, I., Gomez-Rodriguez, M., and Gummadi, K. P. Fairness constraints: A flexible
414 approach for fair classification. *Journal of Machine Learning Research*, 20(75):1–42, 2019.
- 415 Zemel, R., Wu, Y., Swersky, K., Pitassi, T., and Dwork, C. Learning fair representations. In
416 *International Conference on Machine Learning*, pp. 325–333, 2013.
- 417 Zhao, H. and Gordon, G. J. Inherent tradeoffs in learning fair representations. *Neural Information*
418 *Processing Systems*, 2019.
- 419 Zliobaite, I. On the relation between accuracy and fairness in binary classification. *Interna-*
420 *tional Conference on Machine Learning, Workshop on Fairness, Accountability, and Transparency in*
421 *Machine Learning*, 2015.
- 422 Calders, Toon and Karim, Asim and Kamiran, Faisal and Ali, Wasif and Zhang, Xiangliang Con-
423 trolling attribute effect in linear regression. *IEEE 13th international conference on data mining*,
424 2013.
- 425 Berk, Richard and Heidari, Hoda and Jabbari, Shahin and Joseph, Matthew and Kearns, Michael
426 and Morgenstern, Jamie and Neel, Seth and Roth, Aaron A convex framework for fair regression.
427 *Fairness, Accountability, and Transparency in Machine Learning*, 2017.

- 428 Agarwal, Alekh and Dudik, Miroslav and Wu, Zhiwei Steven Fair regression: Quantitative defini-
429 tions and reduction-based algorithms. *International Conference on Machine Learning*, 2019.
- 430 Fitzsimons, Jack and Al Ali, AbdulRahman and Osborne, Michael and Roberts, Stephen A general
431 framework for fair regression. *Entropy*, 2019.
- 432 Calders, Toon and Verwer, Sicco Three naive bayes approaches for discrimination-free classification.
433 *Data Mining and Knowledge Discovery*, 2010.
- 434 Billionnet, Alain and Elloumi, Sourour Using a mixed integer quadratic programming solver for the
435 unconstrained quadratic 0-1 problem. *Mathematical Programming*, 2007.
- 436 Redmond, Michael and Baveja, Alok A data-driven software tool for enabling cooperative informa-
437 tion sharing among police departments. *European Journal of Operational Research*, 2002.
- 438 Cortez, Paulo and Silva, Alice Maria Gonçalves Using data mining to predict secondary school
439 student performance. *EUROSIS-ETI*, 2008.

440 Checklist

- 441 1. For all authors...
- 442 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
443 contributions and scope? [Yes]
- 444 (b) Did you describe the limitations of your work? [Yes] See Assumptions 1, 2
- 445 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 446 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
447 them? [Yes]
- 448 2. If you are including theoretical results...
- 449 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Assump-
450 tions 1, 2
- 451 (b) Did you include complete proofs of all theoretical results? [Yes]
- 452 3. If you ran experiments...
- 453 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
454 mental results (either in the supplemental material or as a URL)? [Yes] Code and data
455 are not included but instructions to reproduce them has been included
- 456 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
457 were chosen)? [Yes] See Section 5 and Appendix M
- 458 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
459 ments multiple times)? [Yes] See Figure 2
- 460 (d) Did you include the total amount of compute and the type of resources used (e.g., type
461 of GPUs, internal cluster, or cloud provider)? [N/A]
- 462 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 463 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 464 (b) Did you mention the license of the assets? [N/A]
- 465 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 466
- 467 (d) Did you discuss whether and how consent was obtained from people whose data
468 you're using/curating? [N/A]
- 469 (e) Did you discuss whether the data you are using/curating contains personally identifi-
470 able information or offensive content? [N/A]
- 471 5. If you used crowdsourcing or conducted research with human subjects...
- 472 (a) Did you include the full text of instructions given to participants and screenshots, if
473 applicable? [N/A]
- 474 (b) Did you describe any potential participant risks, with links to Institutional Review
475 Board (IRB) approvals, if applicable? [N/A]
- 476 (c) Did you include the estimated hourly wage paid to participants and the total amount
477 spent on participant compensation? [N/A]