# Removing the influence of a group variable in high-dimensional predictive modelling

Emanuele Aliverti

*Department of Statistical Sciences, University of Padova*

E-mail: aliverti@stat.unipd.it

Kristian Lum

*Human Rights Data Analysis Group, San Francisco*

James E. Johndrow

*Department of Statistics, Stanford University*

David B. Dunson

*Department of Statistical Science, Duke University*

**Summary**. In many application areas, predictive models are used to support or make important decisions. There is increasing awareness that these models may contain spurious or otherwise undesirable correlations. Such correlations may arise from a variety of sources, including batch effects, systematic measurement errors, or sampling bias. Without explicit adjustment, machine learning algorithms trained using these data can produce poor out-of-sample predictions which propagate these undesirable correlations. We propose a method to pre-process the training data, producing an adjusted dataset that is statistically independent of the nuisance variables with minimum information loss. We develop a conceptually simple approach for creating an adjusted dataset in high-dimensional settings based on a constrained form of matrix decomposition. The resulting dataset can then be used in any predictive algorithm with the guarantee that predictions will be statistically independent of the group variable. We develop a scalable algorithm for implementing the method, along with theory support in the form of independence guarantees and optimality. The method is illustrated on some simulation examples and applied to two case studies: removing machine-specific correlations from brain scan data, and removing race and ethnicity information from a dataset used to predict recidivism. That the motivation for removing undesirable correlations is quite different in the two applications illustrates the broad applicability of our approach.

*Keywords*: Constrained optimization; Criminal justice; Neuroscience; Orthogonal predictions; Predictive modelling; Singular value decomposition.

## 1. Introduction

As machine learning models are deployed in high-stakes application areas, there has been growing recognition that machine learning algorithms can reproduce unwanted associations encoded in the data upon which they were trained (e.g. Dunson, 2018; Zech et al., 2018). Such potentially problematic associations are typically between the outcome of interest $Y$, which is specific to the scientific application under investigation, and one or more group variables $Z$. A model that retains correlations between $Y$ and $Z$ may be undesirable for different reasons depending on the application area. This article focuses on a method for adjusting training data to prevent models fit to that data from learning associations between undesirable features and the outcome of interest. Without proper adjustment, the association between the outcome of interest and unwanted group information can obscure scientifically relevant signals leading to biased results and misleading conclusions, or perpetuate socially undesirable stereotypes or practices.

In this article we focus on two challenging research application areas where this problem is particularly salient. The first is the field of "omic" data, where there is strong interest in using algorithms to relate images, scans, and medical records to outcomes of interest, such as disease and subject-specific phenotypes. A recent example highlights how unwanted associations can manifest in these settings. Zech et al. (2018) showed that radiographic image data encoded information on the specific hospital system from which the data were collected, since different systems tended to use different imaging equipment. The systematic differences in the image data produced by the different equipment were not readily apparent to the human eye.

However, when the data were used to train a machine learning model for pneumonia screening, the model learned to associate these equipment-specific characteristics with the outcome of interest. That is, the model leveraged the fact that some hospital systems saw higher rates of disease and also used equipment that left a specific signature in the imaging data. The model's predictions were then partially based on correlations between the signature of the machines used in the higher-prevalence hospital systems and the disease. The model's reliance on such associations jeopardizes the generalizability of the results, as there is no reason to believe that other higher prevalence hospitals will use similar equipment in the future. This issue is also detrimental for in-sample evaluation and inferential purposes, since regarding such associations as risk factors for disease is clinically misleading.

In the field of medical imaging and biostatistics, removal of spurious effects in high dimensional data has been an important problem for at least the last decade. Various other tools for dealing with this issue have been developed, with methods based on matrix factorisation being very popular in genomics (e.g. Wall et al., 2003). For example, Alter et al. (2000) use a singular value decomposition (SVD) of the original expression matrix to filter out the singular vector associated with systematic and non-experimental variations, and reconstruct an adjusted expression matrix with the remaining singular vectors. A similar strategy was pursued in Bylesjö et al. (2007) and Leek and Storey (2007), where matrix decompositions are used to separate the latent structures of the expression matrix that are not affected by batch assignment from the nuisance batch variable. Model-based methods have also been successful. For example, distance discrimination based on support vector machines (Benito et al., 2004) and multilevel models estimated via empirical Bayes (Johnson et al., 2007) have been developed to deal with groups with few observations. For a more comprehensive review of approaches for batch correction and array normalization see Luo et al. (2010), Lazar et al. (2012), and references therein.

The focus of these methods has mostly been on pre-processing high-dimensional data in order to remove the batch effects and then conducting different analysis on the adjusted data; for example, visualising the data to locate clusters of similarly expressed genes (e.g. Lazar et al., 2012). Data visualisation is often used as the main evaluation tool to determine the success in removing batch-effects, while assessment and evaluation for out-of-sample data has not received much interest (Luo et al., 2010). We take an approach that is similar to the methods based on matrix decomposition, with the focus being on providing a low-dimensional approximation of the data matrix which guarantees that batch effects are removed with minimum information loss. Unlike the approaches popular in the biostatistics literature, our approach explicitly focuses on quantitatively addressing the problem of removing unwanted associations by imposing a form of statistical independence among the adjusted data and the group membership variable and explicitly enforcing such constraints in the matrix decomposition. In addition, our approach focuses on the generalization problem, guaranteeing that predictions for *future* data will be statistically independent from group variables.

Our second motivating example comes from algorithmically automated or assisted decision-making, where algorithms are deployed to aid or replace human decision-making. Examples include algorithms for selecting short-listed candidates for job positions or college admission. If the measurement error varies with group membership status, associations between the group variable and the outcome of

interest may obscure correlations with scientific factors of genuine interest. Moreover, when groups are defined by socially sensitive attributes – such as race or gender – a model that learns these associations may lead to unfair or discriminatory predictions. For example, there has been much recent attention in the "algorithmic fairness" literature on the use of criminal risk assessment models, many of which use demographic, criminal history, and other information to predict whether someone who has been arrested will be re-arrested in the future. These predictions then inform decisions on pre-trial detention, sentencing, and parole. In many cases the model seeks to predict re-arrest, which may be influenced by racially biased policing (Bridges and Crutchfield, 1988; Rudovsky, 2001; Simoiu et al., 2017). When risk assessment models are trained using these data, the end result is that individuals in racial minority groups tend to be systematically assigned to higher risk categories on average (Johndrow et al., 2019; Angwin et al., 2016). Even when information on race is not explicitly included as a covariate, this phenomenon persists, since other variables that are included in the model are strongly associated with race.

Models for which one group is more likely to be flagged as "high risk" are in conflict with one notion of algorithmic fairness– statistical parity. However, many definitions are available in the literature. See Berk et al. (2018); Mitchell et al. (2018); Corbett-Davies and Goel (2018) for recent overviews of mathematical notions of fairness. Arguably, each definition proposed in the literature can be regarded as a valid concept of algorithmic fairness, though which notion is most sensible depends heavily on the data and application. Some methods for achieving notions of fairness have focused on altering the objective function used in estimation to, for example, balance the false positive or negative rate across groups (e.g. Hardt et al., 2016; Zafar et al., 2017). Others have focused on defining metrics that more fairly characterize relevant differences between individuals to enforce notions of individual fairness. The recent work of Zhang et al. (2018a) achieves a variety of notions of fairness through adversarial learning— to achieve demographic parity, the adversary is tasked with inferring group membership status from model's predictions. This method is applied to recidivism prediction in Wadsworth et al. (2018).

Other approaches to creating "fair" models modify the training data rather than the objective function, with some focusing on modifications to $Y$ (Kamiran and Calders, 2009). The approach we take most closely follows Johndrow et al. (2019), Adler et al. (2018) and Feldman et al. (2015) in pre-processing the covariates $X$, guaranteeing that any model estimated on the "adjusted" data will produce out-of-sample predictions with formal guarantees of statistical independence. Our algorithm has significant advantages over existing approaches in scalability and ability to handle high-dimensional data efficiently.

Motivated by datasets from these two distinct application areas, in this article we propose a simple method to adjust high-dimensional datasets in order to remove associations between covariates and a nuisance or otherwise undesirable variable. Our procedure creates an adjusted dataset with guarantees that algorithms estimated using the adjusted data produce predictions which will be statistically independent from the nuisance variable. The main advantage of our contribution is its simplicity and scalability to a large number of covariates ($p$), including when $p$ is greater than the number of observations $n$. In this sense, it is particularly well-suited for applications like brain imaging and "omic" data, in which the observed covariates are high-dimensional. It also has significant advantages in the case of highly collinear predictors, which is very common in applications. Moreover, the solution we propose has theoretical guarantees, both in terms of optimal dimension reduction and the ability to achieve independence from the undesirable variables under a linearity condition. We also provide guarantees of minimal information loss during the pre-processing.

It is worth mentioning that the goal of achieving statistical independence is not without controversy. In medical imaging predictions, one could imagine a scenario where hospitals use a similar equipment type and also treat similar patient populations. Scrubbing the data of dependence on imaging equip-

ment may also remove some information that is useful for prediction and is not captured by other covariates in the dataset. A generous read of a common critique of creating race-independent data in the recidivism prediction setting is that race might be an important risk factor, since the correlation between race and the likelihood of re-arrest reflects the reality that people of color are more likely to experience conditions leading to criminality, such as poverty. Creating race-independent predictions may *under*-estimates the likelihood of re-arrest, even after accounting for systematic bias in the process by which individuals are arrested. In reality, it is difficult to know how much of the difference in re-arrest rate can be attributed to differential patterns of enforcement relative to differential participation in criminal activity. In this setting, it may be ethically appealing to choose equal group-wise treatment, which is achieved by enforcing independence between predictions and race. Furthermore, even if a reasonable ground truth were available, one might still want to create race-independent predictions to avoid further exacerbating the existing racial disparities in incarceration rates.

## 2.   Data description

The first case study discussed in this article comes from medical imaging data, and is drawn from a study in neuroscience conducted by the Human Connectome Project (HCP) on $n = 1056$ adult subjects with no history of neurological disorder (Glasser et al., 2016, 2013). The study provides, for each individual, information on the structural interconnections among the 68 brain regions characterizing the Desikan atlas (Desikan et al., 2006), measured through a combination of diffusion and structural magnetic resonance imaging; see Zhang et al. (2018b) for additional details. Many different features are also available, covering a wide range of biographical, physiological and behavioural information at the individual level and technical information related to the specific session in which brain scans were collected. For an extended description of the Humane Connectome Project, the tools involved in the collection process and the aims of the study see Zhang et al. (2018b) and Glasser et al. (2016, 2013). For our purposes, it is enough to characterize the outcomes of interest as physiological and behavioural traits and the covariates as data on the presence and strength of connections between the 68 brain regions.

Recent developments in neuroscience have stimulated considerable interest in analysing the relationship among brain structure or activity and subject-specific traits, with the main focus being on detecting if variations in the brain structure are associated with variation in phenotypes (e.g. Genovese et al., 2002; Zhang et al., 2018b; Durante and Dunson, 2018). There is evidence for the existence of brain differences across subjects with severe drug addictions, both in terms of functional connectivity (Wilcox et al., 2011; Kelly et al., 2011) and volumes of brain regions (Beck et al., 2012; Goldstein et al., 2009). As discussed in the Introduction, a fundamental problem with observational medical data is the presence of spurious associations such as batch effects. In neuroscience, it is well known that subject motion, eye movements, different protocols and hardware-specific features, among many others can complicate data analysis (e.g. Basser and Jones, 2002; Sandrini et al., 2011). Here, we use the HCP data to investigate the relationship between connectivity and drug use while removing batch effects from the identity of the imaging equipment used for each subject.

The second case study discussed in this article comes from algorithmic criminal risk assessment. We will focus here on the COMPAS dataset, a notable example in the fairness literature which includes detailed information on criminal history for more than 6000 defendants in Broward County, Florida. For each individual, several features on criminal history are available, such as the number of past felonies, misdemeanors, and juvenile offenses; additional demographic information include the sex, age and race of each defendant. Defendants are followed over two years after their release, and the data contain an indicator for whether each defendant is re-arrested within this time range. See Larson et al. (2016) for details on the data collection process and additional information on the dataset. The

focus of this example is on predicting two-year recidivism as a function of defendant's demographic information and criminal history. As discussed in the Introduction, algorithms for making such predictions are routinely used in courtrooms to advise judges, and concerns about the fairness of such tools with respect to race of the defendants were recently raised, mainly using this specific dataset as an example (Angwin et al., 2016). Therefore, it is of interest to develop novel methods to produce predictions while avoiding disparate treatment on the basis of race.

## 3. Generating data orthogonal to groups

### 3.1. Notation and setup

Let $X$ denote an $n \times p$ data matrix of $p$ features measured over $n$ subjects, and let $Z$ denote an additional nuisance variable; for example, brain scanner in the neuroscience case study or race in the recidivism example. We focus for simplicity on a scalar $Z$, but the methods directly generalize to multivariate nuisance variables. We seek to estimate $\widetilde{X}$, an $n \times p$ reconstructed version of the data matrix that is orthogonal to $Z$ with minimal information loss. In our setting, the reconstructed version is used to produce a prediction rule $\hat{Y}(\widetilde{X})$ that returns a prediction $\hat{Y}$ of $Y$ for any input $\tilde{X}$. Our aim is to guarantee that $\hat{Y}(\widetilde{X})$ is statistically independent of $Z$.

We will focus on statistical models linear in the covariates, such as generalised linear models or support vector machines with linear kernels. It is easy to check that when $\hat{Y}(\widetilde{X})$ is a linear function of $\widetilde{X}$, $\text{cov}(\widetilde{X}, Z) = 0$ implies $\text{cov}(\hat{Y}, Z) = 0$. Our procedure transforms $X$ by imposing orthogonality between $\widetilde{X}$ and $Z$, while attempting to preserve as much of the information in $X$ as possible. The former requirement is geometrically analogous to requiring that $Z$ is in the null space of $\widetilde{X}$, so that it is not possible to predict $Z$ using the transformed variables in a statistical model which is linear in the covariates. Non-linear dependencies could still be present in the transformed matrix $\widetilde{X}$, so that predictions of models which are not linear in the covariates would still depend on $Z$. One potential solution to this problem is to include interactions in the matrix $X$ before transformation; we use this approach in the recidivism case study in Section 5.2. Another possible solution is to attempt to remove nonlinear dependence as well, as in Johndrow et al. (2019). We do not pursue the latter here, favoring the simplicity and computational scalability of imposing a zero covariance constraint. Indeed, one of our two motivating applications is a large $p$, small $n$ problem, for which fast computation is critical.

Another desirable property of our procedure is dimensionality reduction. In high-dimensional settings, it is often assumed that covariate matrices have approximately a low-rank representation. We express a reduced rank approximation of $X$ as $\widetilde{X} = SU^T$, where $U$ is a $p \times k$ matrix of $k$ linear orthonormal basis vectors and $S$ is the $n \times k$ matrix of associated scores. The problem of preprocessing the data to ensure Orthogonality to Groups (henceforth OG) can be expressed as a minimization of the Frobenius distance between the original data and the approximated version, $\|X - \widetilde{X}\|_F^2$, under the constraint $\langle \widetilde{X}, Z \rangle = 0$. With $\widetilde{X} = SU^T$, this leads to the following optimization problem.

$$\operatorname*{arg\,min}_{S,U} \ \|X - SU^T\|_F^2, \quad \text{subject to } \langle SU^T, Z \rangle = 0, \quad U \in \mathcal{G}_{p,k} \tag{1}$$

where $\mathcal{G}_{p,k}$ is the Grassman manifold consisting of orthonormal matrices (James, 1954). While we focus on a univariate categorical $Z$ in this article, the procedure can be used as-is for continuous $Z$, and extension to multivariate $Z$ is straightforward.

Since the constraints are separable, it is possible to reformulate Equation 1 as $p$ distinct constraints, one over each column of $\widetilde{X}$. Moreover, since any column of $\widetilde{X}$ is a linear combination of the $k$ columns of $S$, and $U$ is orthonormal, the $p$ constraints over $\widetilde{X}$ can be equivalently expressed as $k$ constraints over the columns of the score matrix $S$. The matrix $U$ is forced to lie on the Grassman manifold to prevent degeneracy, such as basis vectors being identically zero or solutions with double multiplicity.

The optimization problem admits an equivalent formulation in terms of Lagrange multipliers,

$$\underset{S,U}{\arg\min} \left\{ \frac{1}{n} \sum_{i=1}^{n} \|x_i - \sum_{j=1}^{k} s_{ij} u_j^T\|^2 + \frac{2}{n} \sum_{j=1}^{k} \lambda_j \sum_{i=1}^{n} s_{ij} z_i \right\}, \tag{2}$$

with the introduction of the factor $2/n$ for ease of computation.

### 3.2. *Theoretical support*

The following Lemma characterizes the solution of the OG optimization problem, which can be interpreted as the residual of a multivariate regression among left singular values and a group variable. Let $V_k \Sigma_k U_k^T$ denote the rank-$k$ Singular Values Decomposition (SVD) of $X$.

LEMMA 1. *The problem stated in Equation 1 can be solved exactly, and admits an explicit solution in terms of singular values. The solution is equal to $\widetilde{X} = (I_n - P_z)V_k\Sigma_k U_k^T$, with $P_Z = Z(Z^T Z)^{-1}Z^T$.*

Proofs are given in Appendix A.1. The computational cost of the overall procedure is dominated by the cost of the truncated SVD, which can be computed with modern methods in $O(nk^2)$ (Golub and Van Loan, 2012). The procedure outlined in Lemma 1 is simple and only involves matrix decomposition and least squares theory; hence we can fully characterise the solution and its properties. In the univariate case, the expression in Lemma 1 is easily interpretable since $P_z$ and $(I_n - P_z)$ are $n \times n$ matrices with elements

$$[P_z]_{ij} = \frac{1}{n} + \frac{z_i z_j}{\sum_{i=1}^{n} z_i^2}, \qquad [(I_n - P_z)]_{ij} = \mathbf{I}[i = j] - \frac{1}{n} - \frac{z_i z_j}{\sum_{i=1}^{n} z_i^2}, \quad i,j = 1,\ldots,n.$$

The following Lemma guarantees that the solution $\widetilde{X}$ of the OG algorithm optimally preserves information in $X$.

LEMMA 2. *The solution $\widetilde{X}$ of the OG algorithm is the best rank-$k$ approximation, in Frobenius norm, of the data matrix $X$ under the OG constraint.*

The SVD achieves the minimum error in Frobenius distance among all matrices of rank-$k$ (e.g., Golub and Van Loan, 2012). Naturally, the introduction of additional constraints reduces the accuracy of the approximation relative to the SVD. The following result bounds the additional error analytically.

LEMMA 3. *Let $\widetilde{X}_k = V_k D_k U_k^T$ denote the best rank-$k$ approximation of the matrix $X$ obtained from the truncated SVD of rank $k$. The reconstruction error of the OG algorithm is lower bounded by the optimal error rate of $\widetilde{X}_k$, and the amount of additional error is equal to $\|P_z V_k D_k\|_F^2$.*

The additional reconstruction error can be interpreted as a measure of the collinearity between the subspace spanned by $Z$ and the left singular vectors of the data $X$. The more correlated the singular vectors are with the group variable, the greater the additional error relative to the solution without the OG constraint. When $Z$ is in the null space of $X$, the solution is identical to the truncated singular value decomposition and the reconstruction achieves the minimum error. Therefore, if the correlation between the group variable and the observed data is negligible, the procedure achieves a loss of information which is close to the SVD.

### 3.3. *Sparse OG procedure*

Estimation of low-dimensional structure from high-dimensional data can be challenging when the number of features $p$ is larger than the number of observations $n$. In this setting, common methods in multivariate analysis impose constraints on the elements of a matrix decomposition, usually through

an $\ell_1$-norm penalty to favour sparsity and improve numerical estimation (e.g. Zou et al., 2006; Jolliffe et al., 2003; Witten et al., 2009).

To make the OG problem tractable and stable when the number of features is very large – potentially larger than the number of observations – we introduce additional constraints in the algorithm. We will build our method on a standard procedure to perform sparse matrix decomposition (e.g. Hastie et al., 2015, Chapter 8), and adapt the computations to introduce the orthogonality constraint. We define the Sparse Orthogonal to Group (SOG) optimization problem as follows.

$$\underset{S,U}{\arg\min} \left\| X - SU^T \right\|_F^2$$
$$\text{subject to} \quad \|u_j\|_2 \le 1, \|u_j\|_1 \le t, \|s_j\|_2 \le 1, s_j^T s_l = 0, \ s_j^T Z = 0, \tag{3}$$

for $j = 1, \dots, k$ and $l \ne j$. The problem in Equation 3 includes sparsity constraints over the vectors $u_j$ and imposes orthogonality constraints among the score vectors $s_j$ and the group variable, since the reconstruction $\widetilde{X} = SU^T$ is a linear combination of the vectors $s_j$.

To fix ideas, we focus initially on a rank-1 approximation. Adapting the results in Witten et al. (2009), it is possible to show that the solutions in $s$ and $u$ for Equation 3 when $k = 1$ also solve

$$\underset{s,u}{\arg\max} \ s^T X u \quad \text{subject to} \quad \|u\|_2 \le 1, \|u\|_1 \le t, \ \|s\|_2 \le 1, \ s^T Z = 0. \tag{4}$$

Although the minimisation in Equation 4 is not jointly convex in $s$ and $u$, it can be solved with an iterative algorithm. Since the additional orthogonality constraints do not involve the vector $u$, when $s$ is fixed the minimisation step is mathematically equivalent to a sparse matrix decomposition with constraints on the right singular vectors. This takes the form

$$\underset{u}{\arg\max} \ bu \quad \text{subject to} \quad \|u\|_2 \le 1, \|u\|_1 \le t, \tag{5}$$

with $b = s^T X$ and solution

$$u = g(b, \theta) = \frac{\mathcal{S}_\theta(b)}{\|\mathcal{S}_\theta(b)\|_2},$$

where $\theta$ is a penalty term in the equivalent representation of the constrained problem in (4) in terms of Langrange multipliers, and $\mathcal{S}_\theta(x) := \text{sign}(x)(|x| - \theta)\mathbb{I}(|x| \ge \theta)$ is the soft threshold operator applied over every element separately. The value of $\theta$ is 0 if $\|b\|_1 \le t$, and otherwise $\theta > 0$ is selected such that $\|g(b, \theta)\|_1 = t$ (Hastie et al., 2015; Witten et al., 2009).

When $u$ is fixed, the solution is similar to that described in Section 3.2, which can be seen by rearranging Equation 4 to obtain

$$\underset{s}{\arg\max} \ s^T a \quad \text{subject to} \quad \|s\|_2 \le 1, \ s^T Z = 0, \tag{6}$$

with $a = Xu$ (Witten et al., 2009). The solution to Equation 6 is directly related to the method outlined in Section 3.2, and is given by the following expression.

$$s = \frac{a - \beta Z}{\|a - \beta Z\|_2},$$

with $\beta = (Z^T Z)^{-1} Z^T a$.

Solutions with rank greater than 1 are obtained by consecutive univariate optimization. For the $j$-th pair $(u_j, s_j)$, $j = 2, \dots, k$, the vector $a$ in Equation 6 is replaced with $P_{k-1} X u_j^T$, where $P_{k-1} = I_{n \times n} - \sum_{l=1}^{k-1} s_l s_l^T$ projects $X u_j^T$ onto the complement of the orthogonal subspace $\text{span}(s_l, \dots, s_{k-1})$, thus guaranteeing orthogonality among the vectors $s_j$, $j = 1, \dots, k$. A more detailed description of the algorithm outlined above is given in Appendix A.2.

## 4. Simulation study

We conduct a simulation study to evaluate the empirical performance of the proposed algorithms and compare them with different competitors. The focus of the simulations is on assessing fidelity in reconstructing a high-dimensional data matrix, success in removing the influence of the group variable from predictions for future subjects, and evaluation of the goodness of fit of predictions. We also compare our methods with two popular approaches developed in biostatistics; specifically, the COMBAT method of Johnson et al. (2007) and the PSVA approach of Leek and Storey (2007). These methods adjust for batch effects via Empirical Bayes and matrix decomposition, respectively. The R code implementing the methods illustrated in this article is available at github.com/emanuelealiverti/sog, and it also includes a tutorial to reproduce the simulation studies. The approaches used as competitors are available through the R package SVA (Leek et al., 2012).

The first simulation scenario puts $n = 1000$, $p = 200$, and $X$ has rank $k = 10$. The data matrix $X$ is constructed in two steps. We simulate a loading matrix $S$, with size $(n, k)$, and a score matrix $U$ with size $(k, p)$, with entries sampled from independent Gaussian distributions. A group variable $Z$ of length $n$ is sampled from independent Bernoulli distributions with probability equal to 0.5. Each $p$-dimensional row of the $(n \times p)$ data matrix $X$ is drawn from a $p$-variate standard Gaussian distribution with mean vector $\mu_i = (s_i - \lambda z_i)U$, $i = 1, \ldots, n$ and $\lambda$ is sampled from a $k$-variate Gaussian distribution with identity covariance. Lastly, a continuous response $Y$ with elements $y_i$, $i = 1, \ldots, n$ is sampled by first generating the elements of $\beta$ independently from Uniform$(-5, 5)$, then sampling $y_i \sim N((s_i - \lambda z_i)\beta, 1)$. We highlight that in this setting the data matrix $X$ has a low-rank structure, and the response variable $Y$ is a function both of the group variable $Z$ and the low-dimensional embedding of $X$. In the second simulation scenario, the construction of the data matrix $X$ is similar to the first setting, except that the response variable $Y$ does not depend on $Z$. This is achieved by sampling the elements $y_i$ of $Y$ from standard Gaussians with mean vector $\mu_i = s_i \beta$, $i = 1, \ldots, n$. Therefore, the response $Y$ depends on the data only though the low-dimensional embedding of $X$. The third scenario focuses on a "large $p$ - small $n$" setting, in which the dimension of the data matrix $X$ is $n = 200, p = 1000$ with $k = 10$, and its construction follows the first scenario, with dimensions of the score and loading matrix modified accordingly. In the fourth and last scenario, the low-rank assumption is relaxed and we focus on a case with $k = p$, following the same generating processes described above.

In each scenario, data are divided into a training set and a test set, with size equal to 3/4 and 1/4 of the observations, respectively. Therefore, the number of observations in the training and test set is equal to 750 and 250 respectively in the first, second and fourth scenario, and is equal to 150 and 50 in the third. In each scenario, the proposed procedures and the competitors are applied separately over the training set and the test set. Then, linear regressions are estimated on the adjusted training sets and used to provide predictions $\hat{Y}$ over the different test sets, where the performance is evaluated by comparing predictions with the truth.

Table 1 reports the root mean square error (RMSE), mean absolute error (MAE) and median absolute error (MDAE) for the adjusted predictions from the competitors (COMBAT, PSVA), the proposed methods (OG, SOG) and the unadjusted case (SVD), in which linear regressions are estimated on the left singular vectors of $X$. In order to reduce simulation error in the results from splitting the data into training and test sets, results are averaged across 50 different splits, with standard deviations across splits reported in brackets. Empirical results suggest that models estimated on datasets adjusted with the competitor methods have comparable performances, with PSVA providing more accurate predictions than COMBAT in most settings. Both OG and SOG outperform the competitors in all the settings considered. For example, in the third scenario involving a $p \gg n$ setting, the RMSE of the PSVA is 22.93 compared to 16.15 for SOG, which also provides better results for the other metrics considered. It is worth highlighting that even when $X$ is not low rank, as in the fourth scenario, our methods

**Table 1.** Out-of sample predictions for $Y$. Lower values correspond to more accurate predictions. Values are averaged over $50$ random splits, with standard deviations reported in brackets. Best performance is highlighted in boldface.

|  |  | COMBAT | PSVA | OG | SOG | SVD |
|---|---|---|---|---|---|---|
| *Scenario 1* | RMSE | 20.73 (3.78) | 17.49 (5.81) | 15.31 (3.30) | **13.19** (2.64) | 21.48 (4.44) |
|  | MAE | 16.54 (3.00) | 13.63 (4.49) | 12.13 (2.54) | **10.45** (2.09) | 17.00 (3.40) |
|  | MDAE | 13.92 (2.67) | 11.39 (3.65) | 10.13 (2.14) | **8.69** (1.84) | 14.08 (2.77) |
| *Scenario 2* | RMSE | 15.37 (2.49) | 13.99 (3.63) | 11.54 (1.85) | **10.36** (1.71) | 15.75 (2.34) |
|  | MAE | 12.28 (1.96) | 10.84 (2.95) | 9.25 (1.51) | **8.28** (1.38) | 12.60 (1.89) |
|  | MDAE | 10.47 (1.74) | 9.16 (2.75) | 7.85 (1.41) | **6.97** (1.19) | 10.71 (1.72) |
| *Scenario 3* | RMSE | 25.78 (4.54) | 22.93 (8.65) | 18.87 (3.96) | **16.15** (2.91) | 26.38 (4.74) |
|  | MAE | 20.87 (3.78) | 18.34 (7.08) | 15.14 (3.31) | **12.99** (2.45) | 21.15 (4.01) |
|  | MDAE | 18.14 (3.88) | 15.54 (6.58) | 12.64 (3.34) | **11.01** (2.47) | 17.86 (4.10) |
| *Scenario 4* | RMSE | 44.74 (2.55) | 45.70 (2.29) | 42.20 (1.98) | **41.93** (2.03) | 45.12 (2.27) |
|  | MAE | 35.96 (2.18) | 36.61 (2.01) | 33.92 (1.65) | **33.76** (1.72) | 36.28 (1.91) |
|  | MDAE | 30.75 (2.45) | 31.20 (2.37) | 29.25 (2.30) | **29.06** (2.03) | 31.00 (2.29) |

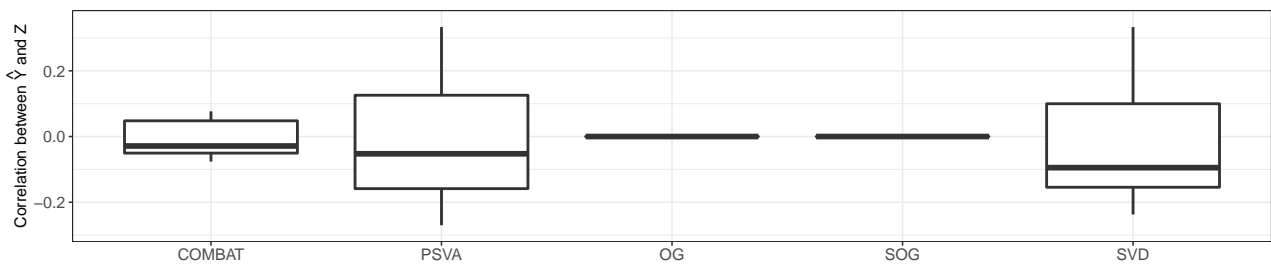provide more accurate out-of-sample predictions than the competitors.

We also compare the accuracy of each method in recovering the data matrix $X$ in simulation scenario 1 with increasing values of the rank $k$ of the approximation $\widetilde{X}$. The rank of the reconstruction is directly specified for OG, SOG and PSVA during estimation. Batch-adjustment via COMBAT, instead, does not utilize a low-rank representation. To create a fair basis for comparison, we use a rank-$k$ SVD decomposition of the batch-corrected matrix obtained from COMBAT in the comparison of reconstruction accuracy. Table 2 illustrates the Frobenius norms between the training set of the first scenario and its approximations for increasing values of the rank of each approximation. Results are averaged over 50 random splits, and standard deviations across splits are reported in brackets.

For each value of the rank $k$, the best reconstruction is, as expected, obtained by the SVD decomposition, which returns a perfect reconstruction for ranks greater or equal than the true value $k = 10$, but of course SVD does not impose orthogonality to $Z$, and thus is reported only to provide a basis for comparison of the reconstruction error of the other methods. Among the procedures that do seek to remove the influence of $Z$, the best reconstruction (minimum error) is obtained by OG for every value of $k$. This empirical result is consistent with Lemma 2, which illustrates the optimality of the method in terms of reconstruction error. For values of $k$ greater than or equal to 10, the error is stable and corresponds to the quantity derived in Lemma 3. The error achieved with the SOG algorithm is considerably higher than that of the OG algorithm. This is expected since the true singular vectors are not sparse in our simulation scenarios. Despite this, SOG performs better than OG in out of sample prediction, suggesting that the additional regularization is useful in prediction. Among the competitors, COMBAT is more accurate in reconstructing $X$ than PSVA, which encounters estimation issues with values of the rank greater than 8, as illustrated by the increasing standard deviation of the reconstruction error with large $k$.

Finally, we analyse success in removing information about $Z$ from predictions based on adjusted data. Figure 1 compares the empirical correlation between the out-of-sample predictions for $\hat{Y}$ and the group variable $Z$, over 50 random splits, for the first simulation scenario. Results from Figure 1 suggest a general tendency of all the adjustment methods to reduce the correlation between the out-of-sample predictions and group variable, confirming the ability of the proposed algorithms and the competitors to remove the linear effect of group information from predictions. Results from the two competitor approaches suggest that predictions from COMBAT have smaller correlation than PSVA on average, and results are also less variable across different splits. The high variability observed with

**Table 2.** Reconstruction error of the data matrix in the first scenario, measured in Frobenius norm. Results are averaged across 50 random splits, with standard deviations reported in brackets. Lower values represent more accurate reconstructions.

| Rank $k$ | COMBAT | PSVA | OG | SOG | SVD |
|---|---|---|---|---|---|
| 2 | 310.48 (1.45) | 347.78 (2.97) | 309.03 (1.48) | 338.01 (1.56) | 305.65 (1.31) |
| 3 | 282.80 (1.36) | 326.58 (3.25) | 281.07 (1.38) | 318.90 (1.49) | 277.32 (1.23) |
| 4 | 256.61 (1.24) | 304.34 (4.70) | 254.53 (1.27) | 300.93 (1.38) | 250.28 (1.07) |
| 5 | 228.90 (1.35) | 286.01 (15.83) | 226.37 (1.37) | 282.60 (1.71) | 221.56 (1.06) |
| 6 | 202.38 (1.33) | 264.54 (18.54) | 199.27 (1.33) | 266.01 (1.50) | 193.74 (1.00) |
| 7 | 174.52 (1.31) | 247.20 (36.40) | 170.67 (1.33) | 249.55 (1.42) | 164.13 (1.04) |
| 8 | 143.73 (1.41) | 231.63 (58.65) | 138.83 (1.41) | 233.48 (1.56) | 130.62 (1.13) |
| 9 | 107.80 (1.62) | 291.11 (106.68) | 100.86 (1.63) | 217.29 (1.69) | 89.19 (1.11) |
| 10 | 60.97 (2.85) | 139.36 (187.70) | 47.09 (3.16) | 201.88 (1.92) | 0.00 (0.00) |
| 11 | 62.49 (2.82) | 162.44 (192.83) | 47.09 (3.16) | 201.64 (3.21) | 0.00 (0.00) |



**Fig. 1.** Correlation between $\hat{Y}$ and $Z$ across 50 random splits into training and test set in the first scenario.

PSVA, and its very modest effect at reducing correlation with $Z$, might be due to the convergence issues described previously. As expected, predictions from OG and SOG always have zero correlation across all splits, since this constraint is explicitly imposed in the algorithm during optimization.

## 5.  Application

### 5.1.  *Human connectome project*

To apply OG and SOG to the imaging case study, the brains scans were first vectorised into a $n \times p$ matrix $X$, with $n = 1065$ subjects and $p = 2278$ features corresponding to the strength of connection among all pairs of brain regions. The outcome of interest is coded as a binary variable $Y$ indicating a positive result to a drug test for at least one among Cocaine, Opiates, Amphetamines, MethAmphetamine and Oxycontine. The nuisance variable $Z$ indicates the machine on which the data were gathered. Data are randomly divided into a training and a test set, with size equal to 798 and 267 observations, respectively. In order to reduce simulation error from data splitting, results are averaged over 50 different splits into training and test. Table 3 illustrates the averages and the standard deviations of Accuracy (ACC), Area Under the Roc Curve (AUC), True Positive Rates (TPR), and True Negative Rates (TNR) for the out-of-sample predictions across splits. The threshold for predicting observations as positive corresponds to the proportion of positive observations in the training data, which is very small (around 3% for both scanners).

To provide a basis for comparison, analyses are conducted using the original covariates without any adjustment. The first columns of Table 3 represent results for a logistic regression using sparse PCA with $k = 30$ components as covariates. The second and third columns compare predictive performance for Lasso (LASSO) and Random Forest (RF), using all the unadjusted available covariates. The right half of Table 3 shows results for the adjusted procedures, with the columns OG and SOG giving

**Table 3.** Predictive performance on the HCP dataset.  Higher values correspond to more accurate predictions. Best performance is highlighted in boldface.

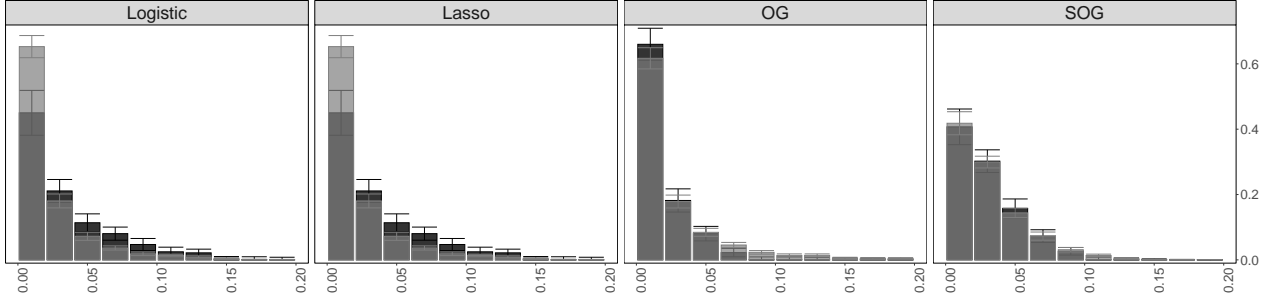|     | LOGISTIC | LASSO | RF | OG | SOG | LASSO, OG | RF, OG |
|-----|----------|-------|-----|-----|-----|-----------|--------|
| ACC | 0.60 (0.06) | **0.62** (0.06) | 0.62 (0.06) | 0.72 (0.32) | **0.79** (0.23) | 0.53 (0.06) | 0.58 (0.07) |
| AUC | **0.54** (0.08) | 0.52 (0.08) | 0.51 (0.10) | 0.51 (0.06) | 0.53 (0.08) | **0.57** (0.10) | 0.52 (0.10) |
| TNR | 0.61 (0.06) | **0.63** (0.06) | 0.63 (0.07) | 0.27 (0.34) | 0.20 (0.25) | **0.43** (0.07) | 0.37 (0.07) |
| TPR | 0.29 (0.18) | **0.40** (0.16) | 0.33 (0.19) | 0.55 (0.33) | **0.57** (0.28) | 0.42 (0.18) | 0.53 (0.18) |



**Fig. 2.** Histograms for $\hat{Y}$ under the four approaches. Light gray corresponds to Scanner 1, black to Scanner 2. Regions where the two histograms overlap are shown as dark gray. The more extensive overlap between the histograms by scanner in the right two panels indicates success of the method at removing scanner information from the data.

results from a logistic regression estimated on $\tilde{X}$ obtained from our two methods, while LASSO,OG and RF,OG give results for predictions from Lasso and Random Forest estimated after OG adjustment. Empirical findings suggest that predictive performance is not noticeably worse for models estimated on reconstructed data $\tilde{X}$.  In some cases, performance improves, while in other cases it declines. A similar phenomenon is observed in the recidivism example, suggesting that the loss of predictive performance after performing our pre-processing procedure can be minimal in applications. Figure 2 provides additional results illustrating the empirical distribution of $\hat{Y}$ under the four approaches. Results suggests that predictions produced by LOGISTIC and LASSO are different across two scans, with predicted probabilities from the first Scanner (black histogram) being on average greater. After pre-processing using our procedures, the empirical distribution of predicted probabilities of drug consumption becomes more similar, with results from predictions formed on SOG-preprocessed data being close to equal across scanners.  We also conducted sensitivity analysis for different values of the approximation rank ranging in $\{10, 50, 100\}$. The results were consistent with the main empirical findings discussed above.

### 5.2.   COMPAS *recidivism data*

We apply our methods to the recidivism dataset by first constructing a design matrix which includes the available features and all the interaction terms, for a total of $p = 64$ variables. The proportion of individuals with a positive outcome was roughly 40% for Caucasian and 50% for non-Caucasian. Although the number of features is not particularly large in this case study, and considerably smaller than in the previous example, methods for pre-processing $X$ that require manual interaction with the statistician, such as Johndrow et al. (2019), are burdensome for even 64 covariates, since each covariate must be modeled separately conditional on others. Moreover, the inclusion of every interaction term induces substantial collinearity in the design matrix, so that low-rank approximation of the data matrix is likely to be both accurate and improve computational and predictive performance for models estimated on the data. Data are randomly divided into a training set and a test set, with size 5090 and
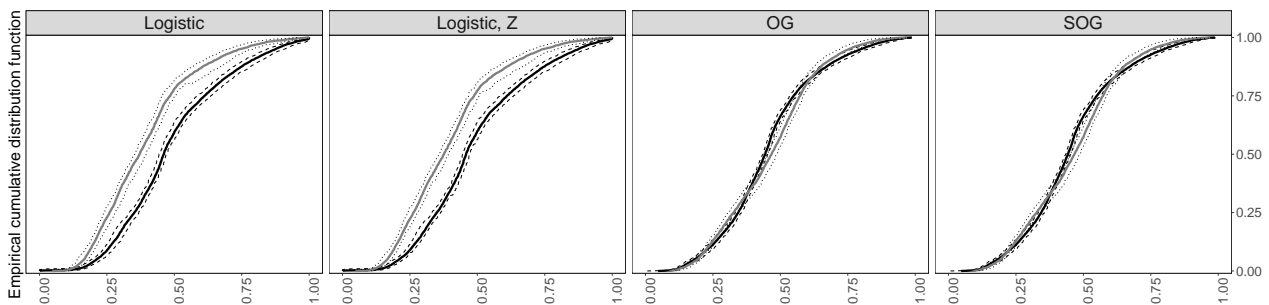
**Table 4.** Predictive performance on the COMPAS dataset. Higher values correspond to more accurate performances. Best performance is highlighted in boldface.

|  | LOGISTIC | LOGISTIC, Z | RF | OG | SOG | RF, OG |
|---|---|---|---|---|---|---|
| ACC | 0.669 (0.02) | **0.671** (0.01) | 0.671 (0.01) | **0.654** (0.01) | 0.654 (0.01) | 0.606 (0.01) |
| AUC | 0.712 (0.02) | **0.714** (0.02) | 0.712 (0.01) | **0.708** (0.01) | 0.708 (0.01) | 0.642 (0.01) |
| TNR | 0.719 (0.02) | 0.716 (0.04) | **0.766** (0.02) | 0.659 (0.02) | **0.661** (0.02) | 0.608 (0.01) |
| TPR | 0.609 (0.05) | **0.617** (0.03) | 0.558 (0.02) | **0.649** (0.01) | 0.647 (0.01) | 0.602 (0.02) |
| PPV | 0.644 (0.02) | 0.645 (0.02) | **0.665** (0.02) | 0.613 (0.01) | **0.614** (0.01) | 0.562 (0.02) |
| NPV | 0.689 (0.02) | **0.692** (0.01) | 0.675 (0.01) | **0.692** (0.01) | 0.691 (0.01) | 0.647 (0.01) |

1697 observations, respectively. Table 4 reports Accuracy (ACC), Area Under the Roc Curve (AUC), True Positive Rates (TPR), True Negative Rates (TNR), Positive Predicted Values (PPV) and Negative Predicted Values (NPV) of the out-of-sample predictions. Results are averaged over 50 random splits of the data, with standard deviations reported in brackets.

The first and second columns of Table 4 represent, respectively, results for a logistic regression using all the available original covariates plus all interaction terms, and all the covariates and interaction terms plus race, while the third column corresponds to a Random Forest trained on the original data. The second part of 4 illustrates results for the adjusted procedures, with OG and SOG reporting results from a logistic regression estimated on the adjusted data and RF,OG from a Random Forest estimated on data preprocessed using OG. As observed in the case of the brain scan data, predictive performance on the whole is quite similar when data are first pre-processed using our procedures, indicating that in some applications there is an almost "free lunch;" information on $Z$ can be removed with very little effect on the overall predictive performance on average.

Figure 3 assesses success in removing dependence on race from $\hat{Y}$ by comparing the empirical cumulative distribution function (CDF) of the predicted probabilities with and without pre-processing of $X$. The first two panels of Figure 3 represent the empirical CDF of $\hat{Y}$ for the models denoted as LOGISTIC and LOGISTIC,Z in Table 4. Results are averaged across splits, with bands corresponding to the 0.975 and 0.025 quantiles reported as dotted and dashed lines to illustrate variability of the curves across splits. Without adjustment, white defendants (gray curves) are assigned lower probabilities of recidivism. This issue affects predictions both excluding the race variable (first panel) and including it is a predictor (second panel). The third and fourth panels correspond, respectively, to predictions obtained from logistic regressions estimated on the data pre-processed with the OG and SOG procedures. The gap between the empirical CDFs is substantially reduced, both with the standard OG and the sparse implementation, leading to distributions of the predictive probabilities which are more similar across different racial groups. This indicates that the procedure was largely successful at rendering predictions independent of race.



**Fig. 3.** Empirical cumulative distribution functions for $\hat{Y}$ under the four approaches. Gray refers to white ethnicity, black to non-white. Dashed and dotted lines corresponds to confidence bands.

## 6. Discussion

We have proposed a simple approach for adjusting a dataset so that predictive algorithms applied to the adjusted data produce predictions that are independent of a group variable. This is motivated by two very different applications, involving adjusting for batch effects in neuroimaging and removing the influence of race/ethnicity in risk assessments in criminal justice. Although these applications areas are fundamentally different, our adjustment procedure can be used in both cases. An appealing aspect is that the procedure is agnostic to the type of predictive algorithm used, so that a single adjusted dataset can be released without prior knowledge of the types of predictive algorithms that will be applied.

There are several interesting areas for future research building upon the proposed framework. One interesting direction is to consider more complex types of covariates. For example, in neuroimaging one may obtain a tensor- or function-valued predictor. For functional predictors, it is potentially straightforward to modify the proposed framework by leveraging developments in functional principal components analysis (FPCA), combining FPCA factorizations with our proposed algorithm. For tensor predictors, one can similarly rely on tensor PCA-type factorizations. However, in such settings, it is important to think carefully about the ramifications of orthogonality assumptions, as unlike for matrices most tensors are not orthogonally decomposable. Another important direction is to consider more complex data structures; for example, containing longitudinal observations, a network-type dependence structure and potentially censoring. In the risk assessment setting, data are increasingly collected over time and there is dependence across different individuals which should be accommodated.

## Acknowledgments

## A. Appendix

### A.1. Results on OG procedure

PROOF (PROOF OF LEMMA 1). Focusing on the case $k = 1$, the approximation of the original set of data consists of finding the closest rank-1 matrix (vector). Equation 2 can be reformulated as

$$\arg\min_{s_1, u_1}\left\{\frac{1}{n}\sum_{i=1}^{n}\|x_i - s_{i1}u_1^T\|^2 + \frac{2}{n}\lambda_1\sum_{i=1}^{n}s_{i1}z_i\right\}, \tag{7}$$

and some algebra and the orthonormal condition on $u_1$ allows us to express the loss function to be minimized as

$$L(s_1, u_1) = \frac{1}{n}\sum_{i=1}^{n}(x_i - s_{i1}u_1^T)^T(x_i - s_{i1}u_1^T) + \frac{2}{n}\lambda_1\sum_{i=1}^{n}s_{i1}z_i$$

$$= \frac{1}{n}\sum_{i=1}^{n}(x_i^T x_i - 2s_{i1}x_i u_1^T + s_{i1}^2) + \frac{2}{n}\lambda_1\sum_{i=1}^{n}s_{i1}z_i.$$

The function is quadratic, and the partial derivative in $s_{i1}$ leads to

$$\frac{\partial}{\partial s_{i1}}L(s_1, u_1) = \frac{1}{n}(-2x_i u_1^T + 2s_{i1}) + \frac{2}{n}\lambda_1 z_i,$$

with stationary point given by $s_{i1} = x_i u_1^T - \lambda_1 z_i$. The optimal score for the $i$-th subject is obtained by projecting the observed data onto the first basis, and then subtracting $\lambda_1$-times $z$. The constraint does not involve the orthonormal basis $u_1$, hence the solution of Equation 7 for $u_1$ is equivalent to the unconstrained scenario. A standard result of linear algebra states that the optimal $u_1$ for Equation 7 without constraints equivalent to the first right singular vector of $X$, or equivalently to the first eigenvector of the matrix $X^T X$ (e.g., Hastie and Tibshirani, 2009). Plugging in the solution for $u_1$ and setting the derivative with respect to $\lambda_1$ equal to 0 leads to

$$\sum_{i=1}^{n}(x_i u_1^T - \lambda_1 z_i)^T z_i = 0 \qquad \lambda_1 = \frac{\sum_{i=1}^{n} x_i u_1^T z_i}{\sum_{i=1}^{n} z_i^2} = \frac{\langle X u_1^T, z \rangle}{\langle z, z \rangle}, \tag{8}$$

a least squares estimate of $X u_1^T$ over $z$.

Consider now the more general problem formulated in Equation 2. The derivatives with respect to the generic element $s_{ij}$ can be calculated easily due to the constraint on $U$, which simplifies the computation. Indeed, the optimal solution for the generic score $s_{ij}$ is given by

$$s_{ij} = x_i u_j^T - \lambda_j z_i, \tag{9}$$

since $u_i^T u_j = 0$ for all $i \neq j$ and $u_j^T u_j = 1$ for $j = 1, \ldots, k$. The solution has an intuitive interpretation, since it implies that the optimal scores for the $j$-th dimension are obtained projecting the original data over the $j$-th basis, and then subtracting $\lambda_j$-times the observed value of $z$. Moreover, since the OG constraints do not involve any vector $u_j$, the optimization with respect to the basis can be derived from known results in linear algebra. The optimal value for the vector $u_j$, with $j = 1, \ldots, k$, is equal to the first $k$ right singular values of $X$, sorted accordingly to the associated singular values (e.g., Bishop, 2006; Hastie et al., 2015).

The global solution for $\lambda = (\lambda_1, \ldots, \lambda_k)$ can be derived from least squares theory, since we can interpret Equation 9 as a multivariate linear regression where the $k$ columns of the projected matrix $X U^T$ are response variables and $z$ a covariate. The general optimal value for $\lambda_k$ is then equal to the multiple least squares solution

$$\lambda_k = \frac{\langle X u_k^T, z \rangle}{\langle z, z \rangle}.$$

PROOF (PROOF OF LEMMA 2). Since the optimization problem of Equation 1 is quadratic with a linear constraint, any local minima is also a global minima. The solution performed via the singular value decomposition and the least squares constitute a stationary point, that is also global minimum.

PROOF (PROOF OF LEMMA 3). Let $V_k D_k U_k^T$ define the rank-$k$ truncated SVD decomposition of the matrix $X$, using the first $k$ left and right singular vectors, and the first $k$ singular vales. Let $\widetilde{X}_{OG}$ define the approximated reconstruction obtained by the OG algorithm. The reconstruction error between the original data matrix $X$ and its low-rank approximation $\widetilde{X}_{OG}$ can be decomposed as follow.

$$\begin{aligned}
\|X - \widetilde{X}_{OG}\|_F^2 &= \|X - (V_k D_k - Z\lambda)U_k^T\|_F^2 \\
&= \|X - V_k D_k U_k^T + Z\lambda U_k^T\|_F^2 \\
&= \|X - V_k D_k U_k^T\|_F^2 + \|Z\lambda U_k^T\|_F^2 + 2\langle X - V_k D_k U_k^T, Z\lambda U_k^T \rangle_F.
\end{aligned}$$

The Frobenius-inner product term is equal to 0 due to the orthogonality of the singular vectors, and rearranging terms the following expression is obtained.

$$\|X - \widetilde{X}_{OR}\|_F^2 - \|X - V_k D_k U_k^T\|_F^2 = \|Z\lambda U^T\|_F^2 = \|Z\lambda\|_F^2,$$

Since the optimal value for $\lambda$ is equal to the least squares solution of $Z$ over $V_k D_k$, it follows that $\|Z\lambda\|_F^2 = \|Z(Z^T Z)^{-1} Z^T V_k D_k\|_F^2 = \|P_Z V_k D_k\|_F^2$, and the proof is complete.

## A.2. *Sparse* OG *procedure*

The following pseudo-code illustrates the key-steps to implement the SOG procedure illustrated in Section 3.3. An R implementation is available at github.com/emanuelealiverti/sog.

---

**Algorithm 1:** SOG algorithm

**Input:** Data matrix $X$, $n \times p$. Approximation rank $k$.

**for** $j = 1, \ldots, k$ **do**

    **while** *Changes in $u_j$ and $s_j$ are not sufficiently small* **do**

        Compute $\beta_j$ via least squares as

$$\beta_j = (Z^T Z)^{-1} Z^T P_{j-1} X u_j,$$

        with $P_{j-1} = I_{n \times n} - \sum_{l=1}^{j-1} s_l s_l^T$

        Update $s_j \in \mathbb{R}^n$ as

$$s_j = \frac{P_{j-1} X u_j - \beta_j Z}{\|P_{j-1} X u_j - \beta_j Z\|_2}$$

        Update $u_j \in \mathbb{R}^p$ as

$$u_j = \frac{\mathcal{S}_\theta(X^T s_j)}{\|\mathcal{S}_\theta(X^T s_j)\|_2},$$

        where $\mathcal{S}_\theta(x) = \text{sign}(x)(|x| - \theta)\mathbb{I}(|x| \geq \theta)$ and

        **if** $\|X^T s_j\|_1 \leq t$ **then**

          | Set $\theta = 0$

        **else**

          | Set $\theta > 0$ such that $\|u_j\|_1 = t$

**Output:** Set $d_j = s_j^T X u_j$. Let $S$ denote the $n \times k$ matrix with columns $d_j s_j$, $j = 1, \ldots, k$. Let $U$ denote the $p \times k$ sparse matrix with rows $u_j$, $j = 1, \ldots, k$. Return $\widetilde{X} = SU^T$

---

## References

Adler, P., Falk, C., Friedler, S. A., Nix, T., Rybeck, G., Scheidegger, C., Smith, B. and Venkatasubramanian, S. (2018) Auditing black-box models for indirect influence. *Knowledge and Information Systems*, **54**, 95–122.

Alter, O., Brown, P. O. and Botstein, D. (2000) Singular value decomposition for genome-wide expression data processing and modeling. *Proceedings of the National Academy of Sciences*, **97**, 10101–10106.

Angwin, J., Larson, J., Mattu, S. and Kirchner, L. (2016) Machine bias: There's software used across the country to predict future criminals. and it's biased against blacks. *ProPublica*. URL: https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing.

Basser, P. J. and Jones, D. K. (2002) Diffusion-tensor mri: theory, experimental design and data analysis–a technical review. *NMR in Biomedicine: An International Journal Devoted to the Development and Application of Magnetic Resonance In Vivo*, **15**, 456–467.

Beck, A., Wüstenberg, T., Genauck, A., Wrase, J., Schlagenhauf, F., Smolka, M. N., Mann, K. and Heinz, A. (2012) Effect of brain structure, brain function, and brain connectivity on relapse in alcohol-dependent patients. *Archives of general psychiatry*, **69**, 842–852.

Benito, M., Parker, J., Du, Q., Wu, J., Xiang, D., Perou, C. M. and Marron, J. S. (2004) Adjustment of systematic microarray data biases. *Bioinformatics*, **20**, 105–114.

Berk, R., Heidari, H., Jabbari, S., Kearns, M. and Roth, A. (2018) Fairness in criminal justice risk assessments: The state of the art. *Sociological Methods & Research*, 0049124118782533.

Bishop, C. M. (2006) *Pattern recognition and machine learning*. Springer, New York.

Bridges, G. S. and Crutchfield, R. D. (1988) Law, social standing and racial disparities in imprisonment. *Social Forces*, **66**, 699–724.

Bylesjö, M., Eriksson, D., Sjödin, A., Jansson, S., Moritz, T. and Trygg, J. (2007) Orthogonal projections to latent structures as a strategy for microarray data normalization. *BMC bioinformatics*, **8**, 207.

Corbett-Davies, S. and Goel, S. (2018) The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023*.

Desikan, R., Ségonne, F., Fischl, B., Quinn, B., Dickerson, B., Blacker, D., Buckner, R., Dale, A., Maguire, R. and Hyman, B. (2006) An automated labeling system for subdividing the human cerebral cortex on mri scans into gyral based regions of interest. *Neuroimage*, **31**, 968–980.

Dunson, D. B. (2018) Statistics in the big data era: Failures of the machine. *Statistics & Probability Letters*, **136**, 4–9.

Durante, D. and Dunson, D. B. (2018) Bayesian inference and testing of group differences in brain networks. *Bayesian Analysis*, **13**, 29–58.

Feldman, M., Friedler, S. A., Moeller, J., Scheidegger, C. and Venkatasubramanian, S. (2015) Certifying and removing disparate impact. *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 259–268.

Genovese, C. R., Lazar, N. A. and Nichols, T. (2002) Thresholding of statistical maps in functional neuroimaging using the false discovery rate. *Neuroimage*, **15**, 870–878.

Glasser, M. F., Smith, S. M., Marcus, D. S., Andersson, J. L., Auerbach, E. J., Behrens, T. E., Coalson, T. S., Harms, M. P., Jenkinson, M., Moeller, S. et al. (2016) The human connectome project's neuroimaging approach. *Nature Neuroscience*, **19**, 1175–1187.

Glasser, M. F., Sotiropoulos, S. N., Wilson, J. A., Coalson, T. S., Fischl, B., Andersson, J. L., Xu, J., Jbabdi, S., Webster, M., Polimeni, J. R. et al. (2013) The minimal preprocessing pipelines for the human connectome project. *Neuroimage*, **80**, 105–124.

Goldstein, R. Z., Bechara, A., Garavan, H., Childress, A. R., Paulus, M. P., Volkow, N. D. et al. (2009) The neurocircuitry of impaired insight in drug addiction. *Trends in cognitive sciences*, **13**, 372–380.

Golub, G. H. and Van Loan, C. F. (2012) *Matrix computations*, vol. 3. JHU Press.

Hardt, M., Price, E., Srebro, N. et al. (2016) Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, 3315–3323.

Hastie, T. and Tibshirani, R. (2009) The elements of statistical learning; data mining, inference and prediction.

Hastie, T., Tibshirani, R. and Wainwright, M. (2015) *Statistical learning with sparsity: the lasso and generalizations*. CRC press.

James, A. T. (1954) Normal multivariate analysis and the orthogonal group. *The Annals of Mathematical Statistics*, **25**, 40–75.

Johndrow, J. E., Lum, K. et al. (2019) An algorithm for removing sensitive information: application to race-independent recidivism prediction. *The Annals of Applied Statistics*, **13**, 189–220.

Johnson, W. E., Li, C. and Rabinovic, A. (2007) Adjusting batch effects in microarray expression data using empirical bayes methods. *Biostatistics*, **8**, 118–127.

Jolliffe, I. T., Trendafilov, N. T. and Uddin, M. (2003) A modified principal component technique based on the lasso. *Journal of computational and Graphical Statistics*, **12**, 531–547.

Kamiran, F. and Calders, T. (2009) Classifying without discriminating. In *2009 2nd International Conference on Computer, Control and Communication*, 1–6. IEEE.

Kelly, C., Zuo, X.-N., Gotimer, K., Cox, C. L., Lynch, L., Brock, D., Imperati, D., Garavan, H., Rotrosen, J., Castellanos, F. X. et al. (2011) Reduced interhemispheric resting state functional connectivity in cocaine addiction. *Biological psychiatry*, **69**, 684–692.

Larson, J., Mattu, S., Kirchner, L. and Angwin, J. (2016) How we analyzed the COM-PAS recidivism algorithm. *ProPublica*, **9**. URL: https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm.

Lazar, C., Meganck, S., Taminau, J., Steenhoff, D., Coletta, A., Molter, C., Weiss-Solís, D. Y., Duque, R., Bersini, H. and Nowé, A. (2012) Batch effect removal methods for microarray gene expression data integration: a survey. *Briefings in bioinformatics*, **14**, 469–490.

Leek, J. T., Johnson, W. E., Parker, H. S., Jaffe, A. E. and Storey, J. D. (2012) The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics*, **28**, 882–883.

Leek, J. T. and Storey, J. D. (2007) Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS genetics*, **3**, e161.

Luo, J., Schumacher, M., Scherer, A., Sanoudou, D., Megherbi, D., Davison, T., Shi, T., Tong, W., Shi, L., Hong, H. et al. (2010) A comparison of batch effect removal methods for enhancement of prediction performance using maqc-ii microarray gene expression data. *The pharmacogenomics journal*, **10**, 278.

Mitchell, S., Potash, E. and Barocas, S. (2018) Prediction-based decisions and fairness: A catalogue of choices, assumptions, and definitions. *arXiv preprint arXiv:1811.07867*.

Rudovsky, D. (2001) Law enforcement by stereotypes and serendipity: Racial profiling and stops and searches without cause. *U. Pa. J. Const. L.*, **3**, 296.

Sandrini, M., Umiltà, C. and Rusconi, E. (2011) The use of transcranial magnetic stimulation in cognitive neuroscience: a new synthesis of methodological issues. *Neuroscience & Biobehavioral Reviews*, **35**, 516–536.

Simoiu, C., Corbett-Davies, S. and Goel, S. (2017) The problem of infra-marginality in outcome tests for discrimination. *The Annals of Applied Statistics*, **11**, 1193–1216.

Wadsworth, C., Vera, F. and Piech, C. (2018) Achieving fairness through adversarial learning: an application to recidivism prediction. *arXiv preprint arXiv:1807.00199*.

Wall, M. E., Rechtsteiner, A. and Rocha, L. M. (2003) Singular value decomposition and principal component analysis. In *A practical approach to microarray data analysis*, 91–109. Springer.

Wilcox, C. E., Teshiba, T. M., Merideth, F., Ling, J. and Mayer, A. R. (2011) Enhanced cue reactivity and fronto-striatal functional connectivity in cocaine use disorders. *Drug and alcohol dependence*, **115**, 137–144.

Witten, D. M., Tibshirani, R. and Hastie, T. (2009) A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, **10**, 515–534.

Zafar, M. B., Valera, I., Gomez Rodriguez, M. and Gummadi, K. P. (2017) Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *Proceedings of the 26th International Conference on World Wide Web*, 1171–1180. International World Wide Web Conferences Steering Committee.

Zech, J. R., Badgeley, M. A., Liu, M., Costa, A. B., Titano, J. J. and Oermann, E. K. (2018) Confounding variables can degrade generalization performance of radiological deep learning models. *arXiv preprint arXiv:1807.00431.*

Zhang, B. H., Lemoine, B. and Mitchell, M. (2018a) Mitigating unwanted biases with adversarial learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, 335–340. ACM.

Zhang, Z., Allen, G., Zhu, H. and Dunson, D. B. (2018b) Relationships between human brain structural connectomes and traits. *bioRxiv*, 256933.

Zou, H., Hastie, T. and Tibshirani, R. (2006) Sparse principal component analysis. *Journal of computational and graphical statistics*, **15**, 265–286.