

# Robust and flexible learning of a high-dimensional classification rule using auxiliary outcomes

Muxuan Liang<sup>1,\*</sup>, Jaeyoung Park<sup>2</sup>, Qing Lu<sup>1</sup>, Xiang Zhong<sup>3</sup>

<sup>1</sup>Department of Biostatistics, University of Florida, Gainesville, FL 32611, United States, <sup>2</sup>School of Global Health Management and Informatics, University of Central Florida, Orlando, FL 32816, United States, <sup>3</sup>Department of Industrial and Systems Engineering, University of Florida, Gainesville, FL 32611, United States

\*Corresponding author: Muxuan Liang, Department of Biostatistics, University of Florida, 2004 Mowry Road, 5th Floor CTRB, Gainesville, FL 32611, USA ([muxuan.liang@ufl.edu](mailto:muxuan.liang@ufl.edu)).

## ABSTRACT

Correlated outcomes are common in many practical problems. In some settings, one outcome is of particular interest, and others are auxiliary. To leverage information shared by all the outcomes, traditional multi-task learning (MTL) minimizes an averaged loss function over all the outcomes, which may lead to biased estimation for the target outcome, especially when the MTL model is misspecified. In this work, based on a decomposition of estimation bias into two types, within-subspace and against-subspace, we develop a robust transfer learning approach to estimating a high-dimensional linear decision rule for the outcome of interest with the presence of auxiliary outcomes. The proposed method includes an MTL step using all outcomes to gain efficiency and a subsequent calibration step using only the outcome of interest to correct both types of biases. We show that the final estimator can achieve a lower estimation error than the one using only the single outcome of interest. Simulations and real data analysis are conducted to justify the superiority of the proposed method.

**KEYWORDS:** auxiliary outcomes; classification; high-dimensional data; multi-task learning; transfer learning.

## 1 INTRODUCTION

With the adoption of electronic health records and medical information systems, datasets increasingly massive in volume and diverse in variable categories have been leveraged for knowledge discovery and clinical decision support. In some datasets, in addition to the patient outcome of primary interest, multiple relevant health outcomes are available. In this paper, we denote these relevant outcomes as auxiliary outcomes, and our goal is to study how to safely use these auxiliary outcomes to help predict a binary target outcome in a classification framework with high-dimensional linear decision rules.

Our motivating example is predicting whether the improvement in hip joint functions fails to achieve the minimal clinical importance difference (MCID) after total hip arthroplasty (THA). The Hip disability and Osteoarthritis Outcome Scores for Joint Replacement (HOOS JR) survey is a commonly used instrument to measure THA patients' health outcomes. Predicting whether the change of the overall score measured in preoperative and postoperative surveys exceeds the MCID can help inform whether surgery is necessary. However, it is a challenging task due to the large variability of the reported outcomes and the high event rates of achieving MCID (ie, imbalanced labels). In addition to the overall score, the questionnaire also collects disease-specific information that measures the improvement in various perspectives, including pain, sleep, fatigue, and function

(Katakam et al., 2022; Kunze et al., 2020). This motivates us to explore, whether we can leverage these related questionnaire items—auxiliary outcomes to facilitate target outcome prediction (ie, whether the overall score change exceeds MCID).

To model related outcomes jointly, multi-task learning (MTL) has emerged, aiming to exploit commonalities and differences across outcomes (Caruana, 1997). In MTL, it is typically assumed that some parameters are similar across tasks (Ando and Zhang, 2005; Argyriou et al., 2007; Bakker and Heskes, 2003; Maurer et al., 2013; Titsias and Lázaro-Gredilla, 2011; Yu et al., 2005; Zhang et al., 2008; Zhu et al., 2011), or these tasks bear a shared sparsity structure (Gong et al., 2013; 2014; Hernández-Lobato and Hernández-Lobato, 2013; Lounici et al., 2009; Obozinski et al., 2008; Rao et al., 2013; Wang et al., 2016; Yang et al., 2009). Subsequently, a common feature representation can be learned through MTL, and this approach has been widely applied in many fields (Li et al., 2014; Liu et al., 2017; 2015; Mrkšić et al., 2015; Shinohara, 2016; Zhang et al., 2016, 2014). In MTL, since outcomes are equally important, the objective function to be minimized is the averaged loss across all tasks. Different from MTL, we only address the performance of predicting the target outcome. The decision rule learned in MTL driven by the averaged loss might be biased towards predicting the auxiliary outcomes rather than the target outcome, ie, the jointly learned decision rule may not perform well when

predicting the target outcome. Thus, our objective is to develop a robust learning approach that is capable of exploiting commonalities and differences across outcomes with guaranteed performance in target outcome prediction.

Focusing on the performance of predicting the target outcome, a commonly used approach is transfer learning (Olivas et al., 2009). Transfer learning aims to improve the performance of target learners on target domains by transferring the knowledge contained in different but related source domains (Zhuang et al., 2020). Recently, Li et al. (2022) and Bastani (2021) addressed transfer learning problems in high-dimensional linear regressions; Tian and Feng (2023) addressed transfer learning problems in high-dimensional generalized linear models. In their proposed procedures, they (1) adopt a common working model for all auxiliary outcomes; (2) assume the contrast between the parameters in the target model and those in the auxiliary models are sufficiently close in  $l_1$  or  $l_0$  norm.

However, these assumptions are easily violated in many practical settings such as our motivating example. First, the auxiliary outcomes are related but different, and thus, they are not likely to share the same model. Second, the requirement regarding the contrast between the parameters in the target model and those in the auxiliary models can be restrictive for classification problems. For instance, considering both the target and auxiliary outcomes follow logistic regression models, if the parameters in the target model are twice as large as those in the auxiliary models, the contrast of the two sets of parameters is not necessarily small in  $l_1$  or  $l_0$  norm. However, from the perspective of classification problems, the optimal decision boundaries are identical for the target and auxiliary outcomes. Thus, there is a need for a more flexible learning approach that efficiently utilizes the possible similarity between decision boundaries, rather than focusing on the contrast of parameters, especially for classification problems.

In this work, we develop a robust and flexible learning approach using auxiliary outcomes to aid the estimation of a high-dimensional linear decision rule for the target outcome. Specifically, we propose a two-stage procedure. In the first stage, a common linear representation of the covariates is learned with all auxiliary outcomes using MTL to gain efficiency by borrowing relevant information from auxiliary outcomes. In the second stage, a calibration procedure is performed to reduce or correct the bias induced in the first stage to ensure the robustness of the estimator for the target outcome prediction. Compared with the existing literature, our contributions are the following. In the first stage, different from Li et al. (2022), Tian and Feng (2023), where the working models for auxiliary outcomes share similar coefficients and intercepts, we posit different decision rules (or models) for different outcomes to accommodate possible heterogeneity. In the second stage, instead of assuming that the contrast between the parameters in models for auxiliary outcomes and the target outcome enjoys a small  $l_1$  norm or a sparse  $l_0$  norm, we define a novel concept of within-subspace bias and against-subspace bias, and we only assume that the minimal against-subspace bias is sparse in  $l_0$  norm or small in  $l_1$  norm, which is a weaker condition than those in Bastani (2021), Li et al. (2022), Tian and Feng (2023). Theoretically, we show that the proposed estimator always has an estimation error comparable

to that of using only the target outcome, even if the conditions in Li et al. (2022), Tian and Feng (2023) are violated. Especially, we show that with the presence of many weakly dependent outcomes, our proposed method can also lead to a convergence rate faster than the derived rate in Bastani (2021), Li et al. (2022), Tian and Feng (2023) and faster than using only the target outcome.

The rest of the paper is organized as follows. Section 2 introduces the proposed method. In Section 3, we investigate the theoretical properties of the proposed method. In Section 4, we conduct simulations to compare our method with other methods, especially MTL and methods in Li et al. (2022). In Section 5, we apply the proposed method to the motivating study for THA patients. We present a discussion and concluding remarks in Section 6.

## 2 LEARNING USING AUXILIARY OUTCOMES UNDER HETEROGENEOUS MODELS

Let  $\mathbf{X} \in \mathbb{R}^p$  be a  $p$ -dimensional covariate vector excluding the intercept and  $Y_0 \in \{1, -1\}$  be a univariate target outcome. We assume that some auxiliary outcomes are available along with the target outcome  $Y_0$ . We denote the auxiliary outcomes as  $Y_1, Y_2, \dots, Y_J \in \{1, -1\}$ , where  $J$  is the number of auxiliary outcomes.

In our motivating example, the target outcome and auxiliary outcomes are available in the same dataset. There are other scenarios where the target outcome and auxiliary outcomes are not in the same dataset. For example, we may have a separate dataset containing only the auxiliary outcomes and covariates, denoted as the source-only dataset. To accommodate this scenario, we assume that we observe  $n$  samples in the target dataset where both the target outcome and auxiliary outcomes are available, ie,  $\{(\mathbf{X}_i, Y_{0,i}, Y_{1,i}, \dots, Y_{J,i})\}_{i=1}^n$ ; in addition, we observe  $N - n$  samples in the source-only dataset where only auxiliary outcomes are available, ie,  $\{(\mathbf{X}_i, Y_{1,i}, \dots, Y_{J,i})\}_{i=n+1}^N$ . We use  $R_i = 0$  to indicate samples coming from the target dataset, and  $R_i = 1$ , from the source-only dataset. In this work, we consider a high-dimensional setting where  $p > n$ .

Learning a linear decision rule to predict the target outcome  $Y_0$  using covariate vector  $\mathbf{X}$  entails a classification problem. Empirical risk minimization (ERM) is often used to learn such a linear decision rule. Specifically, ERM minimizes a convex surrogate of the loss function, ie,

$$\min_{\theta_0} \ell(\theta_0) := \mathbb{E} [\phi \{Y_0(\mathbf{X}^\top \boldsymbol{\beta}_0 + c_0)\} \mid R = 0], \quad (1)$$

where  $\phi(\cdot)$  is a surrogate loss,  $\theta_0 = (\boldsymbol{\beta}_0, c_0)^\top$ ,  $\boldsymbol{\beta}_0$  indicates the linear direction of  $\mathbf{X}$  and  $c_0$  indicates the intercept for predicting  $Y_0$ . By solving optimization problem (1), the decision rule,  $d_0^*(\mathbf{X})$ , with the form  $d_0^*(\mathbf{X}) = \text{sgn}(\mathbf{X}^\top \boldsymbol{\beta}_0^* + c_0^*)$ , can be used for prediction purposes, where  $\theta_0^* = (\boldsymbol{\beta}_0^*, c_0^*)^\top$  is the minimizer of optimization problem (1). Our goal is to use the auxiliary outcomes to improve the estimation  $\theta_0^*$ .

### 2.1 Step one: learn a linear representation using MTL

In this section, we introduce our proposed method, which consists of two steps. The first step is to learn a linear representation

using MTL incorporating the auxiliary outcomes. Denote the index set of auxiliary outcomes as  $\mathcal{J} = \{1, 2, \dots, J\}$ .

In this work, we consider the following MTL method. We obtain a linear representation  $\hat{\mathbf{w}}_{\mathcal{J}}$  by solving

$$\min_{\mathbf{w}, \{c_j\}_{j \in \mathcal{J}}} \hat{\mathbb{E}}_N \left[ \sum_{j \in \mathcal{J}} \phi \{Y_j(\mathbf{X}^\top \mathbf{w} + c_j)\} \right] + \lambda_N \|\mathbf{w}\|_1, \quad (2)$$

where  $\lambda_N$  is a tuning parameter and  $\hat{\mathbb{E}}_N[\cdot]$  is the empirical expectation of both the target and source-only datasets. In this procedure, we estimate  $J$  decision rules for  $\{Y_j\}_{j \in \mathcal{J}}$ , simultaneously. These decision rules are structured to learn a common parameter  $\mathbf{w}$ , which is the direction shared by all outcomes. In addition, the intercept represented by  $c_j$ 's can be different for each outcome to accommodate possible heterogeneity. Leveraging information from auxiliary outcomes (and/or the source-only dataset), the estimator  $\hat{\mathbf{w}}_{\mathcal{J}}$  can approach  $\mathbf{w}_{\mathcal{J}}^*$  with a low estimation error, where  $\mathbf{w}_{\mathcal{J}}^*$  is the minimizer of  $\min_{\mathbf{w}, \{c_j\}_{j \in \mathcal{J}}} \mathbb{E} \left[ \sum_{j \in \mathcal{J}} \phi \{Y_j(\mathbf{X}^\top \mathbf{w} + c_j)\} \right]$ . Although the first step takes advantage of shared information across multiple outcomes, the estimator  $\hat{\mathbf{w}}_{\mathcal{J}}$  may be biased w.r.t.  $\beta_0^*$ , especially when  $\mathbf{w}_{\mathcal{J}}^*$  is biased w.r.t.  $\beta_0^*$ .

**Remark 1** In our proposed MTL step, we primarily specify different intercepts to accommodate possible heterogeneity. Note that, we can allow any low-dimensional sub-vector of the coefficients to be different to accommodate heterogeneous effects. More detailed discussions can be found in the online [Supplementary Materials](#).

## 2.2 Step two: a novel calibration step

In this section, we present how to de-bias  $\hat{\mathbf{w}}_{\mathcal{J}}$  and construct an improved estimator for  $\beta_0^*$  through a novel calibration step. To start with, we decompose the bias of  $\mathbf{w}_{\mathcal{J}}^*$ ,  $\text{bias}(\mathbf{w}_{\mathcal{J}}^*) := \mathbf{w}_{\mathcal{J}}^* - \beta_0^* = (1 - \gamma)\mathbf{w}_{\mathcal{J}}^* - \delta$ , where  $\delta := \beta_0^* - \gamma\mathbf{w}_{\mathcal{J}}^*$ , ie,  $\beta_0^* = \gamma\mathbf{w}_{\mathcal{J}}^* + \delta$ . The first term in this decomposition,  $(1 - \gamma)\mathbf{w}_{\mathcal{J}}^*$ , is along the direction of  $\mathbf{w}_{\mathcal{J}}^*$ , and thus, we refer to it as the *within-subspace bias*; the remaining term  $\delta$  is referred to as the *against-subspace bias*. Notably,  $\beta_0^*$  is unknown. To determine appropriate  $\gamma$  and  $\delta$ , we leverage the fact that  $\beta_0^*$  minimizes  $\mathbb{E}[\phi\{Y_0(\mathbf{X}^\top \beta_0 + c_0)\}]$ . Thus, we replace  $\beta_0$  by  $\gamma\mathbf{w}_{\mathcal{J}}^* + \delta$  and propose to solve

$$\min_{\delta, \gamma, c_0} \mathbb{E}[\phi\{Y_0(\mathbf{X}^\top \delta + \gamma\mathbf{X}^\top \mathbf{w}_{\mathcal{J}}^* + c_0)\}]. \quad (3)$$

The loss function in (3) incorporates two adjustments to  $\mathbf{w}_{\mathcal{J}}^*$ , which corresponds to the within-subspace bias and against-subspace bias. First, we calibrate the scaling parameter  $\gamma$  along the subspace generated by  $\mathbf{w}_{\mathcal{J}}^*$ . The term  $\gamma$  identifies the within-subspace bias. For instance, if  $\mathbf{w}_{\mathcal{J}}^* = 2\beta_0^*$ , then, setting  $\gamma = 1/2$  can eliminate such a bias. Second, we calibrate the subspace generated by  $\mathbf{w}_{\mathcal{J}}^*$  using  $\delta$ . This calibration accounts for the against-subspace bias. If  $\mathbf{w}_{\mathcal{J}}^* = \beta_0^* - \mathbf{e}$ , then setting  $\delta = \mathbf{e}$  can account for such a bias, where  $\mathbf{e} = (1, 0, \dots, 0)^\top$ .

The decomposition of  $\text{bias}(\mathbf{w}_{\mathcal{J}}^*)$ , ie,  $\text{bias}(\mathbf{w}_{\mathcal{J}}^*) = (1 - \gamma)\mathbf{w}_{\mathcal{J}}^* - \delta$  provides multiple options to adjust for possible bias. For each choice of  $\gamma$ , we can obtain a corresponding  $\delta$  that

leads to a unique decomposition of the bias. For example, when  $\gamma = 1$ , the corresponding  $\delta = \beta_0^* - \mathbf{w}_{\mathcal{J}}^*$ ; when  $\gamma = 1/2$ , the corresponding  $\delta = \beta_0^* - \mathbf{w}_{\mathcal{J}}^*/2$ . Both  $\gamma = 1$  and  $\gamma = 1/2$  lead to a specification of  $\delta$  such that  $\beta_0^* = \gamma\mathbf{w}_{\mathcal{J}}^* + \delta$ . However, under different choices of  $\gamma$ , the  $\delta$ 's may be different in terms of their  $l_0$  and  $l_1$  norms, resulting in different levels of difficulties in estimating them. For example, the contrast,  $\mathbf{w}_{\mathcal{J}}^* - \beta_0^*$ , may not be sparse in  $l_0$  norm nor small in  $l_1$  norm. In this case, the contrast  $\mathbf{w}_{\mathcal{J}}^* - \beta_0^*$  may not be easy to estimate. Among all possible decompositions, the  $\gamma$ 's that can lead to a sparse ( $l_0$  norm) or a small ( $l_1$  norm) against-subspace bias,  $\delta$ , are preferable. For ease of exposition, we focus on the  $\delta$  with the least  $l_1$  norm. The results under  $l_0$  norm can be found in the online [Supplementary Materials](#).

Denote the set of  $\delta$ 's with the least  $l_1$  norm as  $\delta^*$ . To pin down the  $\gamma$  such that  $\delta \in \delta^*$ , we propose a special treatment: we first separate the space of  $\delta$  into several domains such that in each domain, the solution is unique; then, we select the final estimator through a cross-fitting procedure. Below we introduce how these domains are defined, and show that, at least one solution to (3) in these domains satisfies that  $\delta \in \delta^*$ .

**Remark 2** If we only focus on the  $\delta$  with the least  $l_1$  norm, we can directly solve (3) with a lasso penalty and this special treatment is not required. However, if we focus on the  $\delta$  with the least  $l_0$  norm, we need to solve (3) with a  $l_0$  penalty, which is not trivial. The proposed procedure provides a unified approach with theoretical guarantees regardless of how sparsity or scale of against-subspace bias is defined.

We construct the following domain  $\Gamma_k = \{\delta = (\delta_1, \delta_2, \dots, \delta_p)^\top : \delta_k = 0\}$ , where  $k = 1, \dots, p$ . Let  $S_{\mathcal{J}}^*$  be the set of indexes of the non-zero coefficients of  $\mathbf{w}_{\mathcal{J}}^*$ . Due to the strict convexity of  $\phi$  and the assumption that the coordinates of  $\mathbf{X}$  are not linearly dependent, for any  $k \in S_{\mathcal{J}}^*$ , there exists a unique  $\gamma$  such that  $\beta_0^* - \gamma\mathbf{w}_{\mathcal{J}}^* \in \Gamma_k$ . This implies that the objective function in (3) on each  $\Gamma_k$  has a unique minimizer, for any  $k \in S_{\mathcal{J}}^*$ .

Lemma 1 further implies that to determine the  $\gamma$  such that  $\delta \in \delta^*$ , we only need to solve the optimization problem (3) within each domain.

**Lemma 1** There exists a  $k \in S_{\mathcal{J}}^*$  such that the minimizer of the optimization problem (3) in the domain  $\Gamma_k$  is the minimizer of the optimization problem (3) with  $\delta \in \delta^*$ .

Motivated by this, we consider a set of optimization problems

$$\min_{\delta \in \Gamma_k, \gamma, c_0} \hat{\mathbb{E}}_n [\phi\{Y_0(\mathbf{X}^\top \delta + \gamma\mathbf{X}^\top \hat{\mathbf{w}}_{\mathcal{J}} + c_0)\}] + \tilde{\lambda}_n \|\delta\|_1 \quad (4)$$

for  $k \in S$ , where  $\tilde{\lambda}_n$  is a tuning parameter and  $S$  is a set of pre-specified indices. The objective function in (4) is the empirical version of that in (3), and the domain of  $\delta$  is constrained to a set  $\Gamma_k$ . In each  $\Gamma_k$ , the solution of optimization (4) is unique for  $k \in S$ ; when  $S_{\mathcal{J}}^* \subset S$ , the optimization is guaranteed to identify  $\delta$  such that  $\delta \in \delta^*$ . In our implementation,  $\tilde{\lambda}_n$  is chosen via cross-validation and  $S$  is chosen as the index set of nonzero coefficients in  $\hat{\mathbf{w}}_{\mathcal{J}}$ .

To select the final estimator among the different domains of  $\delta$ , we propose a cross-fitting procedure. First, we split the entire

target dataset into  $M$  folds. Denote the index set of the  $m$ -th fold as  $\mathcal{I}_m$  and the dataset excluding the  $m$ -th fold as  $\mathcal{I}_m^c$ . For each fold  $m \in \{1, \dots, M\}$  and each  $k \in S$ , we calculate the minimizer of optimization (4) using the data in  $\mathcal{I}_m^c$ . Denote its minimizer as  $\hat{\delta}_{\mathcal{I}_m^c}(k)$ ,  $\hat{\gamma}_{\mathcal{I}_m^c}(k)$ ,  $\hat{c}_{\mathcal{I}_m^c}(k)$ . Subsequently, we have  $\hat{\beta}_{\mathcal{I}_m^c}(k) = \hat{\delta}_{\mathcal{I}_m^c}(k) + \hat{\gamma}_{\mathcal{I}_m^c}(k)\hat{\mathbf{w}}_{\mathcal{J}}$ . Then, we calculate the loss of  $\hat{\beta}_{\mathcal{I}_m^c}(k)$  and  $\hat{c}_{\mathcal{I}_m^c}(k)$  using the data in  $\mathcal{I}_m$  and denote the calculated loss as  $L_{\mathcal{I}_m}(k)$ , ie,  $L_{\mathcal{I}_m}(k) = \widehat{\mathbb{E}}_{\mathcal{I}_m} \left[ \phi \left\{ Y_0(\mathbf{X}^\top \hat{\beta}_{\mathcal{I}_m^c}(k) + \hat{c}_{\mathcal{I}_m^c}(k)) \right\} \right]$  and  $\widehat{\mathbb{E}}_{\mathcal{I}_m}[\cdot]$  is the empirical average of the data in  $\mathcal{I}_m$ . We repeat this procedure for each  $m \in \{1, \dots, M\}$  and each  $k \in S$ . Finally, we calculate the averaged loss  $L(k) = \sum_{m=1}^M L_{\mathcal{I}_m}(k)/M$ , and choose the  $k^*$  that minimizes  $L(k)$  among all  $k \in S$ . Then, we calculate the final estimator  $\hat{\beta}_0$  and  $\hat{c}_0$  by  $\hat{\beta}_0 = \sum_{m=1}^M \hat{\beta}_{\mathcal{I}_m^c}(k^*)/M$ , and  $\hat{c}_0 = \sum_{m=1}^M \hat{c}_{\mathcal{I}_m^c}(k^*)/M$ . In our simulation and real data analysis, we choose  $M = 2$  for ease of computation. A summary of the entire algorithm can be found in the online [Supplementary Materials](#).

### 3 THEORETICAL PROPERTIES

To provide theoretical support for the proposed method, we investigate the convergence rate of the proposed estimator. The proof of all the lemmas, theorems, and corollaries can be found in the online [Supplementary Materials](#).

We first introduce some additional notations. We denote the index set of the non-zero coordinates of  $\beta_0^*$  as  $S^*$ . The cardinality of  $S^*$  is denoted as  $s^*$ . We also define  $\alpha$  and  $h_b$  such that

$$\begin{aligned} \text{var} \left[ \sum_{j \in \mathcal{J}} Y_j \phi' \left\{ Y_j(\mathbf{X}^\top \mathbf{w}_{\mathcal{J}}^* + c_j^*) \right\} X_k \mid \mathbf{X} \right] &\leq C J^\alpha, \\ \sup_{j \in \mathcal{J}} \left| E \left[ Y_j \phi' \left\{ Y_j(\mathbf{X}^\top \mathbf{w}_{\mathcal{J}}^* + c_j^*) \right\} X_k \mid \mathbf{X} \right] \right| &\leq C h_b, \end{aligned}$$

for all  $k \in \{1, 2, \dots, p\}$ , where  $\phi'(\cdot)$  is the first order derivative of  $\phi(\cdot)$ ,  $X_k$  is the  $k$ -th covariate, and  $C$  is a sufficiently large constant. Assuming  $\phi'(\cdot)$  is bounded, these inequalities automatically hold with  $\alpha = 2$  and  $h_b = O(1)$ . Under certain assumptions, these inequalities may hold with  $\alpha < 2$  and  $h_b \rightarrow 0$  (or  $h_b = 0$ ). For example, when  $Y_j$ 's are mutually independent conditional on  $\mathbf{X}$ , we can take  $\alpha = 1$ ; when

$$\begin{aligned} P(Y_j = 1 \mid \mathbf{X}) / \{1 - P(Y_j = 1 \mid \mathbf{X})\} \\ = \phi'(-\mathbf{X}^\top \mathbf{w}_{\mathcal{J}}^* - c_j^*) / \phi'(\mathbf{X}^\top \mathbf{w}_{\mathcal{J}}^* + c_j^*), \end{aligned} \quad (5)$$

we can take  $h_b = 0$ . Note that, when  $\phi(\cdot)$  is a logistic loss, model (5) is equivalent to logistic models with the same coefficients and different intercepts for auxiliary outcomes. From these examples, we can see that  $\alpha$  controls the mutual dependence between  $Y_j$ 's, conditional on  $\mathbf{X}$ , and  $h_b$  controls the bias of  $P(Y_j = 1 \mid \mathbf{X})$  w.r.t. model (5). Thus, by incorporating  $\alpha$  and  $h_b$ , our theoretical results can accommodate dependent  $Y_j$ 's and model misspecification w.r.t. model (5). A detailed discussion can be found in the online [Supplementary Materials](#). Furthermore, we define

$$h \equiv \inf_{\gamma} \sup_{j \in \mathcal{J}} \left\| \gamma \mathbf{w}_j^* - \beta_0^* \right\|_1, \quad (6)$$

where  $\mathbf{w}_j^*$  and  $c_j^*$  are the minimizer of  $\min_{\mathbf{w}, c_j} \mathbb{E} \left[ \phi \left\{ Y_j(\mathbf{X}^\top \mathbf{w} + c_j) \right\} \right]$ .

**Remark 3** The definition (6) implies a relationship for the cosine angle between  $\mathbf{w}_j^*$  and  $\beta_0^*$  (decision boundaries differ up to an intercept). Specifically, we have

$$\begin{aligned} \sup_{j \in \mathcal{J}} \left| (\beta_0^*)^\top \mathbf{w}_j^* \right| / \left( \|\beta_0^*\|_2 \|\mathbf{w}_j^*\|_2 \right) \\ \geq (\|\beta_0^*\|_2 - h) / (\|\beta_0^*\|_2 + h). \end{aligned}$$

To investigate the theoretical property of  $\hat{\beta}_0$ , we introduce the following assumptions.

**Assumption 1** There is a constant  $C_1$  such that  $\|\mathbf{X}\|_\infty$ ,  $\sup_{\mathbf{X}, j} |\mathbf{X}^\top \mathbf{w}_j^*|$ , and  $|c_j^*|$ 's are upper bounded by  $C_1$  with probability 1.

**Assumption 2** Define  $\tilde{\mathbf{X}} = (1, \mathbf{X})$ . There is a positive constant  $\lambda_{\min}$  such that the smallest eigenvalue of  $\mathbb{E}(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top)$  is lower bounded by  $\lambda_{\min}$ .

**Assumption 3** We assume that  $\sup_{\mathbf{X}} |\mathbf{X}^\top \beta_0^*| \leq C_2$ , and  $|c_0^*| \leq C_2$  with probability 1, where  $C_2$  is a constant. We also require that the ratios between the  $k$ -th coefficients in  $\beta_0^*$  and  $\mathbf{w}_{\mathcal{J}}^*$  are bounded for any  $k \in S$ . We also assume that  $s^* \log p/n \rightarrow 0$ .

**Assumption 4** Define  $\tilde{\mathbf{X}}_k = (\mathbf{X}^\top \mathbf{w}_{\mathcal{J}}^*, \mathbf{X}_{-k}, 1)$ , where  $\mathbf{X}_{-k}$  is the vector of covariates  $\mathbf{X}$  excluding the  $k$ th covariate. We assume that there is a positive constant  $\tilde{\lambda}_{\min}$  such that the smallest eigenvalue of  $\mathbb{E}(\tilde{\mathbf{X}}_k \tilde{\mathbf{X}}_k^\top)$  is lower bounded by  $\tilde{\lambda}_{\min}$  for all  $k \in S$ .

Assumptions 1 and 3 impose a uniform upper bound on the design matrix for technical simplicity. Assumptions 2 and 4 impose a uniform lower bound for the eigenvalues of the design matrix to ensure the identifiability of  $\beta_0^*$ ,  $c_0^*$ ,  $\delta$ , and  $\gamma$  for all  $k \in S$ . In Assumption 4, we also assume that the ratios between the  $k$ -th coefficients in  $\beta_0^*$  and  $\mathbf{w}_{\mathcal{J}}^*$  are bounded, for any  $k \in S$ . To achieve this, we could specify  $S$  such that extremely small coefficients in  $\hat{\mathbf{w}}_{\mathcal{J}}$  are excluded. In addition, we assume that  $s^* \log p/n \rightarrow 0$ . When  $s^* = O(n^\kappa)$  and  $N = O(n^{1/(1-\kappa)})$  with  $0 \leq \kappa < 1$ , this assumption requires that  $p = o\{\exp(N^{1-2\kappa})\}$ .

**Theorem 1** Under Assumptions 1–4, taking  $\lambda_N \gg |\mathcal{J}| \sqrt{\log |\mathcal{J}|/N}$ , and

$$\begin{aligned} \lambda_N \geq \sqrt{2(|\mathcal{J}|^\alpha + |\mathcal{J}|^2 h_b^2) \log p/N} \vee |\mathcal{J}| \log p/N, \quad \text{and} \\ \tilde{\lambda}_n \asymp \sqrt{\log p/n}, \text{ we have} \end{aligned}$$

$$\begin{aligned} \max \left( \left\| \hat{\beta}_0 - \beta_0^* \right\|_2^2, \left\| \hat{c}_0 - c_0^* \right\|_2^2 \right) \\ \lesssim (s^* \lambda_N^2 / |\mathcal{J}|^2 + \lambda_N C_\Sigma h / |\mathcal{J}|) \wedge (C_\Sigma h)^2 \\ + (\tilde{\lambda}_n^2 + \tilde{\lambda}_n h_\delta^* \wedge (h_\delta^*)^2) + \log(|S| \vee n)/n \end{aligned}$$

with probability approaching to 1, where  $h_\delta^*$  is the minimizer of  $\min_{\gamma} \|\beta_0^* - \gamma \mathbf{w}_{\mathcal{J}}^*\|_1$  s.t.  $\beta_0^* - \gamma \mathbf{w}_{\mathcal{J}}^* \in \cup_{k \in S} \Gamma_k$ , and  $C_\Sigma$  is defined in Lemma 2 in online [Supplementary Materials](#).

The resultant rate in Theorem 1 is structured as the sum of three terms. The first term,  $(s^* \lambda_N^2 / |\mathcal{J}|^2 + \lambda_N C_\Sigma h / |\mathcal{J}|) \wedge (C_\Sigma h)^2$ , is related to the estimation error of  $\hat{\mathbf{w}}_{\mathcal{J}}$ . The  $C_\Sigma$  reflects the heterogeneity in the design matrix in (generalized) linear models, which is assumed to be a constant in Tian and Feng (2023) by their Assumption 4 or Li et al. (2022) by their Condition 4. The second term,  $\tilde{\lambda}_n^2 + \tilde{\lambda}_n h_\delta^* \wedge (h_\delta^*)^2$ , is associated with the minimal against-subspace bias of  $\mathbf{w}_{\mathcal{J}}^*$ . The third term,  $\log(|S| \vee n) / n$ , accounts for the variability of selecting  $k^*$  in the algorithm.

Compared with the convergence rate of using only  $n$  samples and the target label  $Y_0$ , ie,  $O_p(s^* \log p / n)$ , the convergence rate shown in Theorem 1 can be faster. For example, when  $N \gg n$  and  $C_\Sigma h \ll s^* \sqrt{N \log p / n^2}$ , the first term is smaller than  $s^* \log p / n$ . For the second term, when  $h_\delta^* \ll s^* \sqrt{\log p / n}$ , we have  $\tilde{\lambda}_n^2 + \tilde{\lambda}_n h_\delta^* \ll s^* \log p / n$ . The third term is always negligible compared with  $s^* \log p / n$ . Hence, when  $N$  is sufficiently large compared with  $n$ ,  $h_\delta^* \ll s^* \sqrt{\log p / n}$  can lead to a convergence rate faster than  $O_p(s^* \log p / n)$ . The convergence rate shown in Theorem 1 can also be faster than  $O_p(s^* \log p / n)$  even if  $N = n$ . A detailed discussion can be found in the online [Supplementary Materials](#).

**Remark 4** To achieve the requirement that  $C_\Sigma h \ll s^* \sqrt{N \log p / n^2}$  and  $h_\delta^* \ll s^* \sqrt{\log p / n}$ , the choice of  $\mathcal{J}$  is important. Without an appropriate selection of  $\mathcal{J}$ , the convergence rate of  $\hat{\beta}_0$  is not necessarily faster than the convergence rate obtained using only the target label  $Y_0$  (with sample size  $n$ ); this phenomenon is referred to as the negative transfer (Tian and Feng, 2023). The transferable source detection algorithm proposed in Tian and Feng (2023) can also be applied in our proposed method to avoid a possible negative transfer (see the online [Supplementary Materials](#) for a detailed discussion).

**Remark 5** Theorem 1 is established under definition (6). The definition of  $h$  enables an upper bound on the  $l_1$  norm of  $\gamma \mathbf{w}_{\mathcal{J}}^* - \beta_0^*$ . In the online [Supplementary Materials](#), we derive another convergence rate using the  $l_0$  norm of  $\gamma \mathbf{w}_{\mathcal{J}}^* - \beta_0^*$ .

#### 4 SIMULATIONS

In this section, we conduct simulations to compare the performance of our proposed method (Calibrated) with other existing approaches (eg, MTL approaches and other transfer learning approaches). One of the comparison methods, referred to as the baseline approach, uses solely the target outcome and directly solves  $\min_{\beta_0, c_0} \mathbb{E}_n [\phi \{Y_0(X^\top \beta_0 + c_0)\}] + \lambda_n \|\beta_0\|_1$ , where the logistic loss is chosen for  $\phi(\cdot)$  and  $\lambda_n$  is tuned by cross-validation. The other approaches for comparison include a direct transfer learning approach and two MTL approaches. The direct transfer learning approach implements a modified algorithm, where one fixes  $\gamma = 1$ , and

$c_j$ 's in Step one are assumed to be the same. This modified algorithm can be considered as an extension of the oracle Trans-Lasso Algorithm (TransferDirect) proposed in Li et al. (2022). The MTL approach 1 (MultiTask1) extends the algorithm proposed in Obozinski et al. (2008) using a logistic loss with a grouped lasso penalty. The MTL approach 2 (MultiTask2) solves  $\min_{\mathbf{w}, \{c_j\}_{j \in \mathcal{J} \cup 0}} \mathbb{E}_n \left[ \sum_{j \in \mathcal{J}} \phi \{Y_j(X^\top \beta + c_j)\} \right] + \lambda_n \|\beta\|_1$ . MultiTask2 shares a similar loss function as the MTL used in Step One for our proposed approach. Comparing the proposed method with the baseline approach, we can examine the performance gained from using the auxiliary outcome. Comparing the proposed method with direct transfer learning, we can see the benefit of the proposed method over the existing transfer learning approaches. By comparing the two MTL approaches, we can examine the difference between transfer learning and MTL approaches when focusing on target label prediction.

Let  $\beta_{U_0}$  be the coefficients related to the latent variable  $U_0$  for the target outcome. We generate experimental data following the simulation scenarios below:

- (i) We set  $n = N$ . Let  $\beta_{U_0} = (1, -1, 1, -1, 0, \dots, 0, 0.5, -0.5, 2, -2, 0.5, 0.5, 0, \dots, 0)^\top$  and  $U_0 = SG(X^\top \beta_{U_0}) + 0.2\epsilon_{U_0}$ , where  $\epsilon_{U_0}$  follows a standard normal distribution. The function  $G(\cdot)$  is the cumulative distribution function of a standard normal distribution. Set  $\tilde{U} = SG(X^\top \beta_{\tilde{U}}) + 0.2\epsilon_{\tilde{U}}$ , where  $\epsilon_{\tilde{U}}$  follows a standard normal distribution, and the  $q$ -th coordinate of  $\beta_{\tilde{U}}$  satisfies that  $\beta_{\tilde{U},q} = \beta_{U_0,q}$  for  $q \neq 2, 4$  and  $\beta_{\tilde{U},2} = \beta_{\tilde{U},4} = 1$ . The target outcome is generated by setting  $Y_0 = \text{sgn}(U_0 - u_{0,1/4})$ , where  $u_{0,1/4}$  is the first quartile of  $U_0$ . We further introduce a weighting parameter  $\omega$  and generate the auxiliary outcome  $Y_1$  by setting  $Y_1 = \text{sgn}(U_1 - u_{1,3/4})$ , where  $U_1 = (1 - \omega)U_0 + \omega\tilde{U}$ , and  $u_{1,3/4}$  is the third quartile of  $U_1$ .
- (ii) We set  $n = 0.2N$  and  $\beta_{U_0} = (1, -1, 1, -1, 0, \dots, 0)^\top$ . We generate  $U_0$  based on a binomial distribution  $B\{8, G(X^\top \beta_{U_0})\}$ , where the number of trials equals 8 and the success probability equals  $G(X^\top \beta_{U_0})$ . Then, we corrupt this  $U_0$ : when  $U_0 \leq 3$ , we set  $U_1 = U_0 + B(3, \omega)$ ; when  $U_0 > 4$ , we set  $U_1 = U_0 - B(3, \omega)$ . The target outcome is set as  $Y_0 = 1\{U_0 > 0\}$ ; the auxiliary outcomes are set as  $Y_j = 1\{U_1 - (2j - 1)\}$ , where  $j = 1, 2, 3, 4$ .

Scenarios (i) and (ii) both involve a parameter  $\omega$  that controls the relevancy between auxiliary and target outcomes. The relevance between auxiliary and target outcomes gradually decreases with the increase of  $\omega$ . Specifically, with the increase of  $\omega$ ,  $\mathbf{w}_{\mathcal{J}}^*$  involves more against-subspace bias. When  $\omega = 0$ , because  $U_0 = U_1$  in both settings, we can show that  $\beta_0^* = \gamma \mathbf{w}_{\mathcal{J}}^*$  for some  $\gamma$  under a Gaussian design.

For covariate vector  $\mathbf{X}$ , we have the following two designs. In Design I, the covariate vector  $\mathbf{X}$  follows Gaussian distribution  $N(\mathbf{0}, \mathbf{I}_p)$ . In Design II, we first generate a  $p$ -dimensional vector following  $N(\mathbf{0}, \Sigma_p)$ , where the  $(l, k)$ th coordinate of  $\Sigma_p$  is

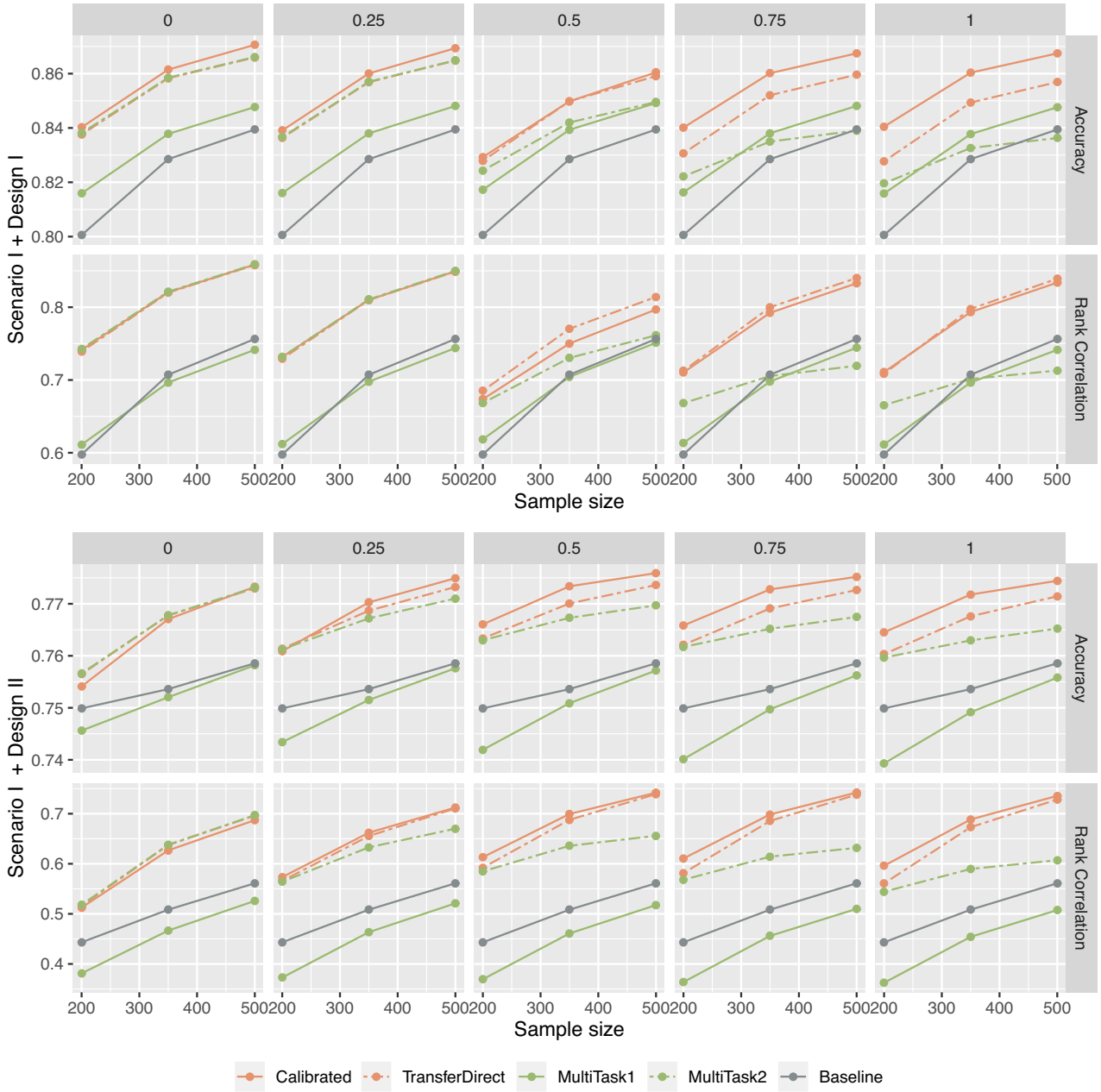


FIGURE 1 Simulation results for Scenario I with the change of sample sizes and  $\omega$ .

$0.5^{|l-k|}$ ; then, for  $l = 1, \dots, \lfloor p/4 \rfloor$ , we replace the  $4l$ -th coordinates in the generated vector with a binary variable. This binary variable is 1 if and only if the generated coordinate is greater than 0. Compared with Design I, Design II has correlated and discrete covariates. We test our methods using both designs for Scenarios (i) and (ii).

To compare the performance of different approaches, we generate a testing dataset with sample size  $n = 10^4$  and calculate two scores. Let  $\hat{\mathbb{E}}_{\text{test}}[\cdot]$  be the empirical expectation calculated using the testing dataset. The first score is the accuracy. Given an estimated decision rule  $\hat{a}_0(\mathbf{X}) = \text{sgn}(\mathbf{X}^\top \hat{\boldsymbol{\beta}}_0 + \hat{c}_0)$ , the accuracy is defined as  $\hat{\mathbb{E}}_{\text{test}}[1\{Y_0 = \hat{a}_0(\mathbf{X})\}]$ . The other score is the rank correlation. We calculate the rank correlation between

$\mathbf{X}^\top \boldsymbol{\beta}_{U_0}$  and  $\mathbf{X}^\top \hat{\boldsymbol{\beta}}_0$  and use it as a proxy of the estimation error. In these simulations, we vary the sample size of the training dataset from  $N = 200, 350$ , to 500 and fix  $p = 1000$ . In Scenario (i), we change  $\omega$  from 0 to 1 with an increment of 0.25. In Scenario (ii), we change  $\omega$  from 0 to 0.3 with an increment of 0.1. We repeat each simulation setting for 500 times.

Figures 1 and 2 illustrate how the performance metrics vary with the increase of sample sizes and  $\omega$ , for simulation Scenarios (i) and (ii), respectively. In Scenario (i), in terms of the accuracy and the rank correlation, the proposed method outperforms the baseline approach regardless of the change of sample sizes and  $\omega$ . Compared with MultiTask1 and MultiTask2, our proposed method and TransferDirect are more robust w.r.t. the change of

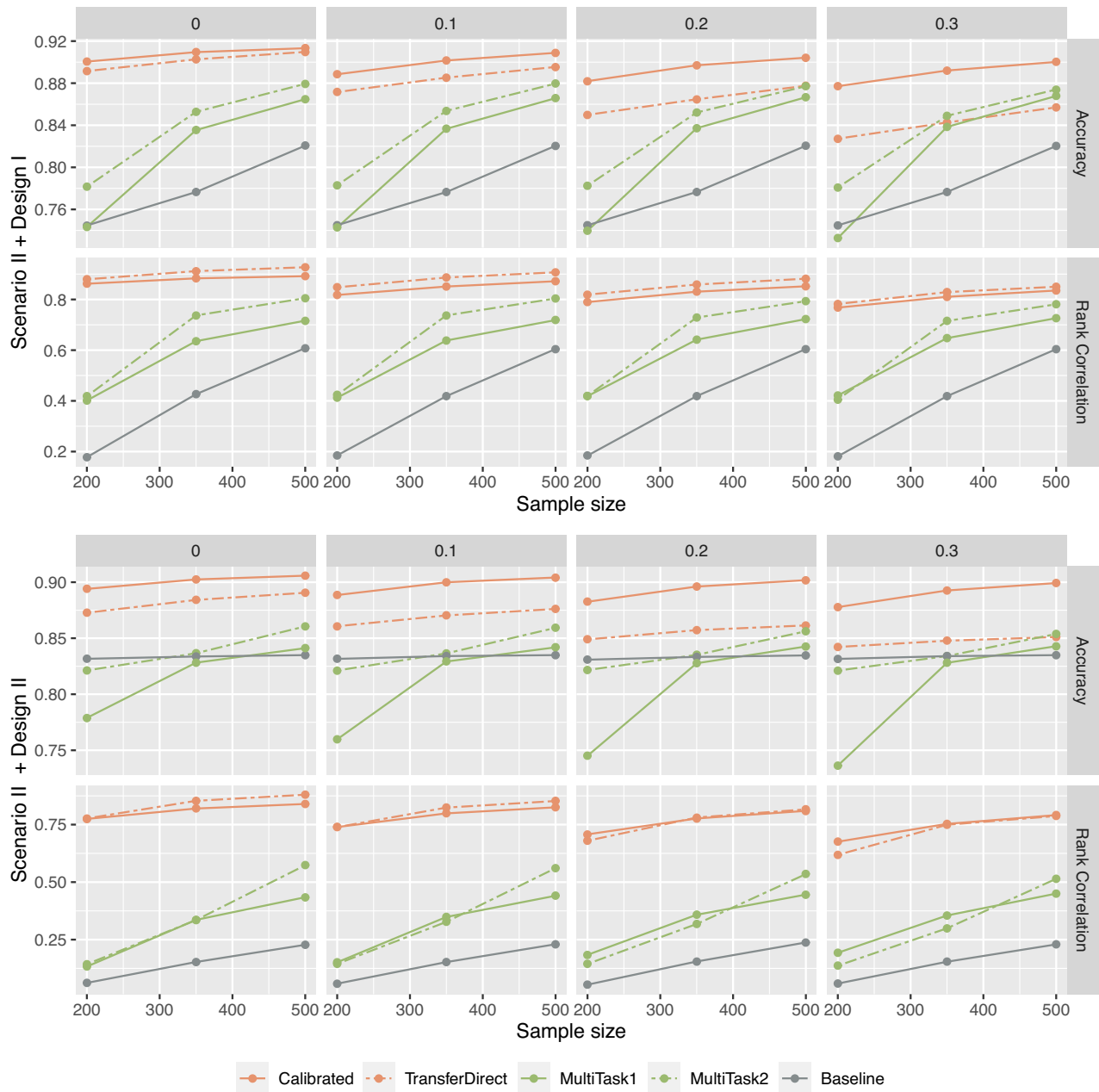


FIGURE 2 Simulation results for Scenario II with the change of sample sizes and  $\omega$ .

$\omega$ ; compared with TransferDirect, our proposed method shows great advantages in terms of prediction accuracy. In Scenario (ii), our proposed method also performs better than other methods regardless of the change of sample sizes and  $\omega$ .

## 5 APPLICATION TO PREDICTING OUTCOMES AFTER THA

In this section, we apply our proposed method to predict the event of not achieving MCID in terms of HOOS JR scores for THA patients. In this dataset, we have 202 patients who underwent an index THA hospitalization, and we consider 13 variables, including race, Risk Assessment and Prediction Tool, and

preoperative HOOS JR survey responses, as covariates. The target outcome is the event of not achieving the anchor-based MCID in their (overall) improvement (Fontana et al., 2019);  $Y_0 = 1$ , if the patient did not achieve the MCID and  $Y_0 = -1$ , otherwise. For the choice of auxiliary outcomes, several other survey questionnaires (eg, improvement in pain) are considered. If the survey outcome is continuous or ordinal, we calculate the quartiles and then define a binary outcome based on whether the original outcome surpasses each quartile or not, and use the transformed outcomes as auxiliary outcomes.

To compare different methods, we randomly split the dataset into a training dataset (70% of the entire dataset) and a testing dataset (30% of the entire dataset). We fit the proposed and other

**TABLE 1** Comparison of the mean (standard error) of accuracy and the Area Under Receiver Operating Characteristic Curve (AUC) estimated from five methods by repeated sample-splittings of the real data.

Method	Accuracy	AUC
Calibrated	0.746 (0.003)	0.712 (0.004)
TransferDirect	0.740 (0.003)	0.701 (0.004)
MultiTask1	0.712 (0.003)	0.660 (0.004)
MultiTask2	0.739 (0.003)	0.713 (0.004)
Baseline	0.731 (0.003)	0.663 (0.004)

comparison methods on the training dataset and calculate the accuracy and the Area Under Receiver Operating Characteristic Curve (AUC) on the testing dataset. The accuracy reflects the estimation error attributed to both coefficients and intercept estimation, allowing for a scalar multiplier difference, ie, if two sets of coefficients and intercepts are proportional, their accuracy will be the same. The AUC calculated by using the coefficients with varying intercepts primarily reflects the estimation error of coefficients (up to a scalar difference). The entire procedure is repeated 500 times. The mean and standard error of the accuracy and AUC are reported in Table 1.

The results therein show that the proposed method achieves the highest accuracy compared with all other methods in terms of prediction accuracy; the proposed method performs comparable to MultiTask2 in terms of the AUC. This implies that the coefficients derived from MultiTask2 and the proposed method are similar (up to a scalar difference), but the intercept estimation from MultiTask2 is more biased. This also implies that the estimation error of MultiTask2 is mainly attributed to the within-subspace bias rather than the against-subspace bias. The online [Supplementary Materials](#) provide an additional application to MCID prediction where our proposed method achieves the highest AUCs among all methods.

## 6 DISCUSSION

In this work, we develop a robust and flexible learning approach to improving decision rule estimation using auxiliary outcomes. Our approach involves a two-step procedure that takes advantage of the information provided by auxiliary outcomes and retains robustness against the bias introduced by auxiliary outcomes. Our novel bias decomposition allows for weaker required conditions and achieves superior performance against existing approaches.

One possible extension is to propose a transfer learning approach under a more relaxed condition. The proposed estimator achieves a fast rate when  $\inf_{\gamma} \sup_{j \in \mathcal{J}} \|\gamma \mathbf{w}_j^* - \beta_0^*\|_1 \ll s^* \sqrt{\log p/n}$ , which is less restrictive than the requirement in Li et al. (2022), Tian and Feng (2023). A more mild condition is, for example, that  $\sup_{j \in \mathcal{J}} \inf_{\gamma} \|\gamma \mathbf{w}_j^* - \beta_0^*\|_1 \ll s^* \sqrt{\log p/n}$  or  $\sup_{j \in \mathcal{J}} \inf_{\gamma} \|\gamma \mathbf{w}_j^* - \beta_0^*\|_2 \ll \sqrt{s^* \log p/n}$ . It could be interesting to see a transfer learning approach that achieves faster convergence rates under these milder conditions. In addition, the proposed method can also be modified to suit various practi-

cal needs. For example, we may use a distributed learner (Duan et al., 2022) to overcome the communication barrier in the first step. This communication barrier comes from the fact that the datasets from different owners (eg, hospitals) cannot be pooled on a single machine due to privacy regulations (eg, HIPAA on sharing medical records).

## ACKNOWLEDGMENTS

The authors thank Dr. Jiwei Zhao, the AE, and two referees for their helpful comments. Muxuan Liang and Jaeyoung Park are the co-first authors of this paper.

## SUPPLEMENTARY MATERIALS

Supplementary material is available at *Biometrics* online.

Web Appendices, referenced in Sections 3–5, are available with this paper at the Biometrics website on Oxford Academic. The R codes for implementing the proposed methods are also posted online with this article.

## FUNDING

None declared.

## CONFLICT OF INTEREST

None declared.

## DATA AVAILABILITY

The data that support the findings in this paper are available from the corresponding author upon reasonable request.

## REFERENCES

- Ando, R. K. and Zhang, T. (2005). A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6, 1817–1853.
- Argyriou, A., Pontil, M., Ying, Y. and Micchelli, C. (2007). A spectral regularization framework for multi-task structure learning. In *Advances in Neural Information Processing Systems*, 20, New York, USA: Curran Associates, Inc.
- Bakker, B. and Heskes, T. (2003). Task clustering and gating for Bayesian multitask learning. *Journal of Machine Learning Research*, 4, 83–99.
- Bastani, H. (2021). Predicting with proxies: transfer learning in high dimension. *Management Science*, 67, 2964–2984.
- Caruana, R. (1997). Multitask learning. *Machine Learning*, 28, 41–75.
- Duan, R., Ning, Y. and Chen, Y. (2022). Heterogeneity-aware and communication-efficient distributed statistical inference. *Biometrika*, 109, 67–83.
- Fontana, M. A., Lyman, S., Sarker, G. K., Padgett, D. E. and MacLean, C. H. (2019). Can machine learning algorithms predict which patients will achieve minimally clinically important differences from total joint arthroplasty? *Clinical Orthopaedics and Related Research*, 477, 1267.
- Gong, P., Ye, J. and Zhang, C. (2013). Multi-stage multi-task feature learning. *Journal of Machine Learning Research*, 14, 2979–3010.
- Gong, P., Zhou, J., Fan, W. and Ye, J. (2014). Efficient multi-task feature learning with calibration. In *Proceedings of the 20th ACM SIGKDD*

- International Conference on Knowledge Discovery and Data Mining, pp. 761–770.
- Hernández-Lobato, D. and Hernández-Lobato, J. M. (2013). Learning feature selection dependencies in multi-task learning. In *Advances in Neural Information Processing Systems*, 26, New York, USA: Curran Associates, Inc.
- Katakam, A., Karhade, A. V., Collins, A., Shin, D., Bragdon, C., Chen, A. F. et al. (2022). Development of machine learning algorithms to predict achievement of minimal clinically important difference for the KOOS-PS following total knee arthroplasty. *Journal of Orthopaedic Research*, 40, 808–815.
- Kunze, K. N., Karhade, A. V., Sadauskas, A. J., Schwab, J. H. and Levine, B. R. (2020). Development of machine learning algorithms to predict clinically meaningful improvement for the patient-reported health state after total hip arthroplasty. *Journal of Arthroplasty*, 35, 2119–2123.
- Li, S., Cai, T. T. and Li, H. (2022). Transfer learning for high-dimensional linear regression: prediction, estimation and minimax optimality. *Journal of the Royal Statistical Society Series B*, 84, 149–173.
- Li, S., Liu, Z.-Q. and Chan, A. B. (2014). Heterogeneous multi-task learning for human pose estimation with deep convolutional neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 482–489. IEEE Computer Society: DC, USA.
- Liu, P., Qiu, X. and Huang, X. (2017). Adversarial multi-task learning for text classification. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, 1, 1–10.
- Liu, W., Mei, T., Zhang, Y., Che, C. and Luo, J. (2015). Multi-task deep visual-semantic embedding for video thumbnail selection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3707–3715.
- Lounici, K., Pontil, M., Tsybakov, A. B. and Van De Geer, S. (2009). Taking advantage of sparsity in multi-task learning, arXiv:0903.1468, preprint: not peer reviewed.
- Maurer, A., Pontil, M. and Romera-Paredes, B. (2013). Sparse coding for multitask and transfer learning. In *Proceedings of the 30th International Conference on Machine Learning*, 343–351, New York, USA: The Association for Computing Machinery.
- Mrkšić, N., Ó Séaghdha, D., Thomson, B., Gašić, M., Su, P.-H., Vandyke, D. et al. (2015). Multi-domain dialog state tracking using recurrent neural networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*, 2, 794–799.
- Obozinski, G., Wainwright, M. J. and Jordan, M. I. (2008). High-dimensional union support recovery in multivariate regression. In *Advances in Neural Information Processing Systems*, 21, New York, USA: Curran Associates.
- Olivas, E. S., Guerrero, J. D. M., Martínez-Sober, M., Magdalena-Benedito, J. R., Serrano, L. et al. (2009). *Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods, and Techniques: Algorithms, Methods, and Techniques*. Hershey, PA, USA: IGI Global.
- Rao, N., Cox, C., Nowak, R. and Rogers, T. (2013). Sparse overlapping sets lasso for multitask learning and its application to fmri analysis. In *Advances in Neural Information Processing Systems*, 26, New York, USA: Curran Associates, Inc.
- Shinohara, Y. (2016). Adversarial multi-task learning of deep neural networks for robust speech recognition. In *Interspeech*, 2369–2372.
- Tian, Y. and Feng, Y. (2023). Transfer learning under high-dimensional generalized linear models. *Journal of the American Statistical Association*, 118, 2684–2697.
- Titsias, M. and Lázaro-Gredilla, M. (2011). Spike and slab variational inference for multi-task and multiple kernel learning. In *Advances in Neural Information Processing Systems*, 24, 2339–2347, Curran Associates, Inc, New York, USA.
- Wang, X., Bi, J., Yu, S., Sun, J. and Song, M. (2016). Multiplicative multi-task feature learning. *Journal of Machine Learning Research*, 17, 2820–2852.
- Yang, X., Kim, S. and Xing, E. P. (2009). Heterogeneous multitask learning with joint sparsity constraints. In *Advances in Neural Information Processing Systems*, 22, 2151–2159. New York, USA: Curran Associates, Inc.
- Yu, K., Tresp, V. and Schwaighofer, A. (2005). Learning Gaussian processes from multiple tasks. In *Proceedings of the 22nd International Conference on Machine Learning*, 1012–1019. New York, USA: The Association for Computing Machinery.
- Zhang, J., Ghahramani, Z. and Yang, Y. (2008). Flexible latent variable models for multi-task learning. *Machine Learning*, 73, 221–242.
- Zhang, W., Li, R., Zeng, T., Sun, Q., Kumar, S. et al. (2016). Deep model based transfer and multi-task learning for biological image analysis. *IEEE Transactions on Big Data*, 6, 322–333.
- Zhang, Z., Luo, P., Loy, C. C. and Tang, X. (2014). Facial landmark detection by deep multi-task learning. In *European Conference on Computer Vision*, 94–108. Switzerland: Springer Cham.
- Zhu, J., Chen, N. and Xing, E. (2011). Infinite latent SVM for classification and multi-task learning. In *Advances in Neural Information Processing Systems*, 24, New York, USA: Curran Associates, Inc.
- Zhuang, F., Qi, Z., Duan, K., Xi, D., Zhu, Y., Zhu, H. et al. (2020). A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109, 43–76.