

# Introduction to Methodologies of SIHR

Zhenyu(Zach) Wang

April 11, 2024

The package **SIHR** aims to perform statistical inference in high-dimensional generalized linear models with continuous and binary outcomes. It provides tools for constructing confidence intervals and performing hypothesis tests for low-dimensional objectives in both one-sample and two-sample regression settings.

## 1 Introduction

We consider the high-dimensional GLMs: for  $1 \leq i \leq n$ ,

$$\mathbb{E}(y_i | X_{i\cdot}) = f(X_{i\cdot}^\top \beta), \quad \text{with } f(z) = \begin{cases} z & \text{for linear model;} \\ \exp(z) / [1 + \exp(z)] & \text{for logistic model;} \end{cases} \quad (1)$$

where  $\beta \in \mathbb{R}^p$  denotes the high-dimensional regression vector,  $y_i \in \mathbb{R}$  and  $X_{i\cdot} \in \mathbb{R}^p$  denote respectively the outcome and the measured covariates of the  $i$ -th observation. Throughout the paper, define  $\Sigma = \mathbb{E}X_{i\cdot}X_{i\cdot}^\top$  and assume  $\beta$  to be a sparse vector with its sparsity level denoted as  $\|\beta\|_0$ . In addition to the one-sample setting, we examine the statistical inference methods for the two-sample regression models. Particularly, we generalize the regression model in (1) and consider:

$$\mathbb{E}(y_i^{(k)} | X_{i\cdot}^{(k)}) = f(X_{i\cdot}^{(k)\top} \beta^{(k)}) \quad \text{with } k = 1, 2 \text{ and } 1 \leq i \leq n_k, \quad (2)$$

where  $f(\cdot)$  is the pre-specified link function defined as (1),  $\beta^{(k)} \in \mathbb{R}^p$  denotes the high-dimensional regression vector in  $k$ -th sample,  $y_i^{(k)} \in \mathbb{R}$  and  $X_{i\cdot}^{(k)} \in \mathbb{R}^p$  denote respectively the outcome and the measured covariates in the  $k$ -th sample.

### 1.1 Package Components

This package consists of five main functions **LF**, **QF**, **CATE**, **InnProd**, and **Dist** implementing the statistical inferences for five different quantities, under the one-sample model (1) or two-sample model (2).

1. **LF**, abbreviated for linear functional, implements the inference approach for  $x_{\text{new}}^\top \beta$ , with  $x_{\text{new}} \in \mathbb{R}^p$  denoting a loading vector. With  $x_{\text{new}} = e_j$  as a special case, **LF** infers the regression coefficient  $\beta_j$ .
2. **QF**, abbreviated for quadratic functional, makes inferences for  $\beta^\top A \beta$ .  $A$  is either a pre-specified sub-matrix or the unknown covariance matrix  $\Sigma$ .
3. **CATE**, abbreviated for conditional average treatment effect, is to make inference for  $f(x_{\text{new}}^\top \beta^{(2)}) - f(x_{\text{new}}^\top \beta^{(1)})$ . This difference measures the discrepancy between conditional means, closely related to the conditional average treatment effect for the new observation with covariates  $x_{\text{new}}$ .
4. **InnProd**, abbreviated for inner products, implements the statistical inference for  $\beta^{(1)\top} A \beta^{(2)}$ . The inner products measure the similarity between the high-dimensional vectors  $\beta^{(1)}$  and  $\beta^{(2)}$ , which is useful in capturing the genetic relatedness in the GWAS applications.
5. **Dist**, short-handed for distance, makes inferences for the weighted distances  $\gamma^\top A \gamma$  with  $\gamma = \beta^{(2)} - \beta^{(1)}$ . The distance measure is useful in comparing different high-dimensional regression vectors.

## 1.2 Outlines

In section 2.2, we propose a unified inference method for  $x_{\text{new}}^\top \beta$  under linear and logistic outcome models. We also discuss inferences for quadratic functionals  $\beta_G^\top A \beta_G$  and  $\beta_G^\top \Sigma_{G,G} \beta_G$  in section 2.3. In the case of the two-sample high-dimensional regression model (2), we develop the inference method for conditional treatment effect  $\Delta(x_{\text{new}}) = f(x_{\text{new}}^\top \beta^{(2)}) - f(x_{\text{new}}^\top \beta^{(1)})$  in section 2.4; we consider inference for  $\beta_G^{(1)\top} A \beta_G^{(2)}$  and  $\beta_G^{(1)\top} \Sigma_{G,G} \beta_G^{(2)}$  in section 2.5 and  $\gamma_G^\top A \gamma_G$  and  $\gamma_G^\top \Sigma_{G,G} \gamma_G$  with  $\gamma = \beta^{(2)} - \beta^{(1)}$  in section 2.6.

## 2 Methodologies

We briefly review the penalized maximum likelihood estimator of  $\beta$  in the high-dimensional GLM (1), defined as:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \ell(\beta) + \lambda_0 \sum_{j=2}^p \frac{\|X_{\cdot j}\|_2}{\sqrt{n}} |\beta_j|, \quad (3)$$

with  $X_{\cdot j}$  denoting the  $j$ -th column of  $X$ , the first column of  $X$  set as the constant 1, and

$$\ell(\beta) = \begin{cases} \frac{1}{n} \sum_{i=1}^n (y_i - X_i^\top \beta)^2 & \text{for linear model} \\ -\frac{1}{n} \sum_{i=1}^n y_i \log \left[ \frac{f(X_i^\top \beta)}{1 - f(X_i^\top \beta)} \right] - \frac{1}{n} \sum_{i=1}^n \log(1 - f(X_i^\top \beta)) & \text{for GLM with binary outcome.} \end{cases} \quad (4)$$

The tuning parameter  $\lambda_0 \asymp \sqrt{\log p/n}$  is chosen by cross-validation. In the penalized regression (3), we do not penalize the intercept coefficient  $\beta_1$ . The penalized estimators have been shown to achieve the optimal convergence rates and satisfy desirable variable selection properties [10, 1, 14, 12]. However, these estimators are not ready for statistical inference due to the non-negligible estimation bias induced by the penalty term [11, 8, 13].

### 2.1 Linear functional for linear model

To illustrate the idea of constructing the inference method, we start with the linear functional for the linear model, which will be extended to a unified version in the section 2.2. For the linear model in (1), we define  $\epsilon_i = y_i - X_i^\top \beta$  and rewrite the model as  $y_i = X_i^\top \beta + \epsilon_i$  for  $1 \leq i \leq n$ . Given the vector  $x_{\text{new}} \in \mathbb{R}^p$ , we construct the point estimator and the CI for  $x_{\text{new}}^\top \beta$ .

A natural idea for the point estimator is to use the plug-in estimator  $x_{\text{new}}^\top \hat{\beta}$  with the penalized estimator  $\hat{\beta}$  defined in (3). However, the bias  $x_{\text{new}}^\top (\hat{\beta} - \beta)$  is not negligible. The work Cai et al. [3] proposed the bias-corrected estimator as,

$$\widehat{x_{\text{new}}^\top \beta} = x_{\text{new}}^\top \hat{\beta} + \hat{u}^\top \frac{1}{n} \sum_{i=1}^n X_i \cdot (y_i - X_i^\top \hat{\beta}), \quad (5)$$

where the second term on the right hand side in (5) is the estimate of negative bias  $-x_{\text{new}}^\top (\hat{\beta} - \beta)$ , and the projection direction  $\hat{u}$  is defined as

$$\hat{u} = \arg \min_{u \in \mathbb{R}^p} u^\top \hat{\Sigma} u \quad \text{subject to: } \|\hat{\Sigma} u - x_{\text{new}}\|_\infty \leq \|x_{\text{new}}\|_2 \lambda \quad (6)$$

$$\left| x_{\text{new}}^\top \hat{\Sigma} u - \|x_{\text{new}}\|_2^2 \right| \leq \|x_{\text{new}}\|_2^2 \lambda, \quad (7)$$

where  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i \cdot X_i^\top$  and  $\lambda \asymp \sqrt{\log p/n}$ . The bias-corrected estimator  $\widehat{x_{\text{new}}^\top \beta}$  satisfies the following error decomposition,

$$\widehat{x_{\text{new}}^\top \beta} - x_{\text{new}}^\top \beta = \underbrace{\hat{u}^\top \frac{1}{n} \sum_{i=1}^n X_i^\top \epsilon_i}_{\text{asyp. normal}} + \underbrace{\left( \hat{\Sigma} \hat{u} - x_{\text{new}} \right)^\top (\beta - \hat{\beta})}_{\text{remaining bias}}.$$

The first constraint in (6) controls the remaining bias term in the above equation while the second constraint in (7) is crucial to ensuring the asymptotic normality of  $\widehat{x_{\text{new}}^\top \beta} - x_{\text{new}}^\top \beta$  for any vector  $x_{\text{new}}$  such that the

variance of the ‘‘asympt. normal’’ term always dominates the ‘‘remaining bias’’ term. Based on the asymptotic normality, we construct the CI for  $x_{\text{new}}^\top \beta$  as

$$\text{CI} = \left( \widehat{x_{\text{new}}^\top \beta} - z_{\alpha/2} \sqrt{\widehat{V}}, \widehat{x_{\text{new}}^\top \beta} + z_{\alpha/2} \sqrt{\widehat{V}} \right) \quad \text{with } \widehat{V} = \frac{\widehat{\sigma}^2}{n} \widehat{u}^\top \widehat{\Sigma} \widehat{u},$$

where  $\widehat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - X_i^\top \widehat{\beta})^2$  and  $z_{\alpha/2}$  denotes the upper  $\alpha/2$  quantile for the standard normal distribution.

## 2.2 Linear functional for GLM

In this subsection, we generalize the inference method specifically for the linear model in Section 2.1 to GLM in (1). Given the initial estimator  $\widehat{\beta}$ , the key step is to estimate the bias  $x_{\text{new}}^\top (\widehat{\beta} - \beta)$ . We can propose a unified version of the bias-corrected estimator for  $x_{\text{new}}^\top \beta$  as

$$\widehat{x_{\text{new}}^\top \beta} = x_{\text{new}}^\top \widehat{\beta} + \widehat{u}^\top \frac{1}{n} \sum_{i=1}^n \omega(X_i^\top \widehat{\beta}) \left( y_i - f(X_i^\top \widehat{\beta}) \right) X_i, \quad (8)$$

with the second term on the right hand side of (8) being the estimate of  $-x_{\text{new}}^\top (\widehat{\beta} - \beta)$ . In consideration of different link functions  $f(\cdot)$  in (1), we shall specify in the following how to construct the projection direction  $\widehat{u}$  and the weight function  $\omega : \mathbb{R} \mapsto \mathbb{R}$  in (8). In Table 1, we consider different GLM models and present

Model	Outcome Type	$f(z)$	$f'(z)$	$\omega(z)$	Weighting
linear	Continuous	$z$	$1$	$1$	
logistic	Binary	$\frac{e^z}{1+e^z}$	$\frac{e^z}{(1+e^z)^2}$	$\frac{(1+e^z)^2}{e^z}$	Linearization
logistic_alter	Binary	$\frac{e^z}{1+e^z}$	$\frac{e^z}{(1+e^z)^2}$	$1$	Link-specific

Table 1: Definitions of the functions  $\omega$  and  $f$  for different GLMs.

the corresponding functions  $f(\cdot)$  and  $\omega(\cdot)$ , together with the derivative  $f'(\cdot)$ . Note that there are two ways of specifying the weights  $w(z)$  for logistic regression. The linearization weighting is proposed in Guo et al. [7] specifically for logistic regression; while Cai et al. [4] constructed the link-specific weighting method for general link function  $f(\cdot)$ . The projection direction  $\widehat{u} \in \mathbb{R}^p$  in (8) is constructed as follows:

$$\begin{aligned} \widehat{u} = \arg \min_{u \in \mathbb{R}^p} u^\top & \left[ \frac{1}{n} \sum_{i=1}^n \omega(X_i^\top \widehat{\beta}) f'(X_i^\top \widehat{\beta}) X_i X_i^\top \right] u \quad \text{subject to:} \\ & \left\| \frac{1}{n} \sum_{i=1}^n \omega(X_i^\top \widehat{\beta}) f'(X_i^\top \widehat{\beta}) X_i X_i^\top u - x_{\text{new}} \right\|_\infty \leq \|x_{\text{new}}\|_2 \lambda \\ & \left| x_{\text{new}}^\top \frac{1}{n} \sum_{i=1}^n \omega(X_i^\top \widehat{\beta}) f'(X_i^\top \widehat{\beta}) X_i X_i^\top u - \|x_{\text{new}}\|_2^2 \right| \leq \|x_{\text{new}}\|_2^2 \lambda. \end{aligned} \quad (9)$$

It has been established that  $\widehat{x_{\text{new}}^\top \beta}$  in (8) is asymptotically unbiased and normal for the linear model [3], the logistic model [6, 4], and the probit model [4]. The variance of  $\widehat{x_{\text{new}}^\top \beta}$  can be estimated by  $\widehat{V}$ , defined as

$$\widehat{V} = \widehat{u}^\top \left[ \frac{1}{n^2} \sum_{i=1}^n \left( \omega(X_i^\top \widehat{\beta}) \right)^2 \widehat{\sigma}_i^2 X_i X_i^\top \right] \widehat{u} \quad \text{with :} \quad (10)$$

$$\widehat{\sigma}_i^2 = \begin{cases} \frac{1}{n} \sum_{j=1}^n \left( y_j - X_j^\top \widehat{\beta} \right)^2, & \text{for linear model} \\ f(X_i^\top \widehat{\beta}) (1 - f(X_i^\top \widehat{\beta})), & \text{for GLM with binary outcome.} \end{cases} \quad (11)$$

Based on the asymptotic normality, the CI for  $x_{\text{new}}^\top \beta$  is:

$$\text{CI} = \left( \widehat{x_{\text{new}}^\top \beta} - z_{\alpha/2} \sqrt{\widehat{V}}, \widehat{x_{\text{new}}^\top \beta} + z_{\alpha/2} \sqrt{\widehat{V}} \right).$$

Subsequently, for the binary outcome case, we estimate the case probability  $\mathbb{P}(y_i = 1 \mid X_i = x_{\text{new}})$  by  $f(\widehat{x_{\text{new}}^\top \beta})$  and construct the CI for  $f(x_{\text{new}}^\top \beta)$  as:

$$\text{CI} = \left( f\left(\widehat{x_{\text{new}}^\top \beta} - z_{\alpha/2} \sqrt{\widehat{V}}\right), f\left(\widehat{x_{\text{new}}^\top \beta} + z_{\alpha/2} \sqrt{\widehat{V}}\right) \right).$$

### 2.3 Quadratic functional for GLM

We now move our focus to inference for the quadratic functional  $Q_A = \beta_G^\top A \beta_G$ , where  $G \subset \{1, \dots, p\}$  and  $A \in \mathbb{R}^{|G| \times |G|}$  denotes a pre-specified matrix of interest. Without loss of generality, we set  $G = \{1, 2, \dots, |G|\}$ . In the following, we propose a unified version of the point estimator and CI under the GLM (1). With the initial estimator  $\widehat{\beta}$  defined in (3), the plug-in estimator  $\widehat{\beta}_G^\top A \widehat{\beta}_G$  suffers from the following error,

$$\widehat{\beta}_G^\top A \widehat{\beta}_G - \beta_G^\top A \beta_G = 2\widehat{\beta}_G^\top A (\widehat{\beta}_G - \beta_G) - (\widehat{\beta}_G - \beta_G)^\top A (\widehat{\beta}_G - \beta_G).$$

The last term in the above decomposition  $(\widehat{\beta}_G - \beta_G)^\top A (\widehat{\beta}_G - \beta_G)$  is the higher-order approximation error under regular conditions; thus the bias mainly comes from the term  $2\widehat{\beta}_G^\top A (\widehat{\beta}_G - \beta_G)$ , which can be expressed as  $2x_{\text{new}}^\top (\widehat{\beta} - \beta)$  with  $x_{\text{new}} = (\widehat{\beta}_G^\top A, \mathbf{0})^\top$ . Hence the term can be estimated directly by applying the linear functional approach in section 2.2. Utilizing this idea, Guo et al. [7, 5] proposed the following estimator of  $Q_A$ ,

$$\widehat{Q}_A = \widehat{\beta}_G^\top A \widehat{\beta}_G + 2\widehat{u}_A^\top \left[ \frac{1}{n} \sum_{i=1}^n \omega(X_i^\top \widehat{\beta}) (y_i - f(X_i^\top \widehat{\beta})) X_i \right],$$

with the second term being the estimate of  $-2\widehat{\beta}_G^\top A (\widehat{\beta}_G - \beta_G)$ , where  $\widehat{u}_A$  is the projection direction defined in (9) with  $x_{\text{new}} = (\widehat{\beta}_G^\top A, \mathbf{0})^\top$ . Since  $Q_A$  is non-negative if  $A$  is positive semi-definite, we truncate  $\widehat{Q}_A$  at 0 and define  $\widehat{Q}_A = \max(\widehat{Q}_A, 0)$ . We further estimate the variance of the  $\widehat{Q}_A$  by

$$\widehat{V}_A(\tau) = 4\widehat{u}_A^\top \left[ \frac{1}{n^2} \sum_{i=1}^n \omega^2(X_i^\top \widehat{\beta}) \widehat{\sigma}_i^2 X_i X_i^\top \right] \widehat{u}_A + \frac{\tau}{n}, \quad (12)$$

where the term  $\tau/n$  with  $\tau > 0$  (default value  $\tau = 1$ ) is introduced as an upper bound for the term  $(\widehat{\beta}_G - \beta_G)^\top A (\widehat{\beta}_G - \beta_G)$ , and  $\widehat{\sigma}_i^2$  is defined in (11). Then given a fixed value of  $\tau$ , we construct the CI as  $\text{CI}(\tau) = \left( \max\left(\widehat{Q}_A - z_{\alpha/2} \sqrt{\widehat{V}_A(\tau)}, 0\right), \widehat{Q}_A + z_{\alpha/2} \sqrt{\widehat{V}_A(\tau)} \right)$ .

Now we turn to the estimation of  $Q_\Sigma = \beta_G^\top \Sigma_{G,G} \beta_G$  where the matrix  $\Sigma_{G,G}$  is unknown and estimated by  $\widehat{\Sigma}_{G,G} = \frac{1}{n} \sum_{i=1}^n X_{iG} X_{iG}^\top$ . Decompose the error of the plug-in estimator  $\widehat{\beta}_G^\top \widehat{\Sigma}_{G,G} \widehat{\beta}_G$ :

$$\widehat{\beta}_G^\top \widehat{\Sigma}_{G,G} \widehat{\beta}_G - \beta_G^\top \Sigma_{G,G} \beta_G = 2\widehat{\beta}_G^\top \widehat{\Sigma}_{G,G} (\widehat{\beta}_G - \beta_G) + \beta_G^\top (\widehat{\Sigma}_{G,G} - \Sigma_{G,G}) \beta_G - (\widehat{\beta}_G - \beta_G)^\top \widehat{\Sigma}_{G,G} (\widehat{\beta}_G - \beta_G).$$

The first term  $\widehat{\beta}_G^\top \widehat{\Sigma}_{G,G} (\widehat{\beta}_G - \beta_G)$  is estimated by applying linear functional approach in Section 2.2 with  $x_{\text{new}} = (\widehat{\beta}_G^\top \widehat{\Sigma}_{G,G}, \mathbf{0})^\top$ ; the second term  $\beta_G^\top (\widehat{\Sigma}_{G,G} - \Sigma_{G,G}) \beta_G$  can be controlled asymptotically by central limit theorem; and the last term  $(\widehat{\beta}_G - \beta_G)^\top \widehat{\Sigma}_{G,G} (\widehat{\beta}_G - \beta_G)$  is negligible due to high-order bias. Guo et al. [7] proposed the following estimator of  $Q_\Sigma$

$$\widehat{Q}_\Sigma = \widehat{\beta}_G^\top \widehat{\Sigma}_{G,G} \widehat{\beta}_G + 2\widehat{u}_\Sigma^\top \left[ \frac{1}{n} \sum_{i=1}^n \omega(X_i^\top \widehat{\beta}) (y_i - f(X_i^\top \widehat{\beta})) X_i \right],$$

where  $\widehat{u}_\Sigma$  is the projection direction constructed in (9) with  $x_{\text{new}} = (\widehat{\beta}_G^\top \widehat{\Sigma}_{G,G}, \mathbf{0})^\top$ . We introduce the estimator  $\widehat{Q}_\Sigma = \max(\widehat{Q}_\Sigma, 0)$  and estimate its variance as

$$\widehat{V}_\Sigma(\tau) = 4\widehat{u}_\Sigma^\top \left[ \frac{1}{n^2} \sum_{i=1}^n \omega^2(X_i^\top \widehat{\beta}) \widehat{\sigma}_i^2 X_i X_i^\top \right] \widehat{u}_\Sigma + \frac{1}{n^2} \sum_{i=1}^n \left( \widehat{\beta}_G^\top X_{i,G} X_{i,G}^\top \widehat{\beta}_G - \widehat{\beta}_G^\top \widehat{\Sigma}_{G,G} \widehat{\beta}_G \right)^2 + \frac{\tau}{n}, \quad (13)$$

where  $\tau > 0$ , the term  $\tau/n$  is introduced as an upper bound for the term  $(\widehat{\beta}_G - \beta_G)^\top \widehat{\Sigma}_{G,G} (\widehat{\beta}_G - \beta_G)$ , and  $\widehat{\sigma}_i^2$  is defined in (11). Then, thanks to the asymptotic normality, for a fixed value of  $\tau$ , we can construct the CI as

$$\text{CI}(\tau) = \left( \max \left( \widehat{Q}_\Sigma - z_{\alpha/2} \sqrt{\widehat{V}_\Sigma(\tau)}, 0 \right), \widehat{Q}_\Sigma + z_{\alpha/2} \sqrt{\widehat{V}_\Sigma(\tau)} \right).$$

## 2.4 Conditional average treatment effects

The inference methods proposed for one sample can be generalized to make inferences for conditional average treatment effects, which can be expressed as the difference between two linear functionals. Let  $A_i \in \{1, 2\}$  denote the treatment assignment for  $i$ -th observation. Consider the two-sample GLMs as

$$\mathbb{E}(y_i | X_i, A_i = 1) = f(X_i^\top \beta^{(1)}) \quad \text{and} \quad \mathbb{E}(y_i | X_i, A_i = 2) = f(X_i^\top \beta^{(2)}),$$

where  $f$  is the link function listed in table 1. Then, for a future individual  $X_i = x_{\text{new}}$ , we define  $\Delta(x_{\text{new}}) = \mathbb{E}(y_i | X_i, A_i = 2) - \mathbb{E}(y_i | X_i, A_i = 1)$ , that measures the difference of the conditional mean of assignment of treatment for the individual with covariates  $x_{\text{new}}$ .

Following (8), we construct the bias-corrected point estimators of  $x_{\text{new}}^\top \widehat{\beta}^{(1)}$  and  $x_{\text{new}}^\top \widehat{\beta}^{(2)}$ , together with their corresponding variance  $\widehat{V}_{(1)}$  and  $\widehat{V}_{(2)}$  as (10). The paper Cai et al. [3] proposed to estimate  $\Delta(x_{\text{new}})$  by  $\widehat{\Delta}(x_{\text{new}})$  as:

$$\widehat{\Delta}(x_{\text{new}}) = f(x_{\text{new}}^\top \widehat{\beta}^{(2)}) - f(x_{\text{new}}^\top \widehat{\beta}^{(1)}).$$

Its variance can be estimated with delta method by:

$$\widehat{V}_\Delta = \left( f'(x_{\text{new}}^\top \widehat{\beta}^{(1)}) \right)^2 \widehat{V}_{(1)} + \left( f'(x_{\text{new}}^\top \widehat{\beta}^{(2)}) \right)^2 \widehat{V}_{(2)}.$$

Then we construct the CI as  $\text{CI} = \left( \widehat{\Delta}(x_{\text{new}}) - z_{\alpha/2} \sqrt{\widehat{V}_\Delta}, \widehat{\Delta}(x_{\text{new}}) + z_{\alpha/2} \sqrt{\widehat{V}_\Delta} \right)$ .

## 2.5 Inner product of regression vectors

The paper Guo et al. [5], Ma et al. [9] have carefully investigated the CI construction for  $\beta_G^{(1)\top} A \beta_G^{(2)}$ , provided with a pre-specified submatrix  $A \in \mathbb{R}^{|\mathcal{G}| \times |\mathcal{G}|}$  and the set of indices  $\mathcal{G} \in \{1, \dots, p\}$ . Let  $\widehat{\beta}^{(1)}$  and  $\widehat{\beta}^{(2)}$  respectively be the initial estimators for their corresponding sample in (2), the plug-in but biased estimator is  $\widehat{\beta}_G^{(1)\top} A \widehat{\beta}_G^{(2)}$ . Its bias can be decomposed as:

$$\begin{aligned} \widehat{\beta}_G^{(1)\top} A \widehat{\beta}_G^{(2)} - \beta_G^{(1)\top} A \beta_G^{(2)} &= \widehat{\beta}_G^{(2)\top} A \left( \widehat{\beta}_G^{(1)} - \beta_G^{(1)} \right) + \widehat{\beta}_G^{(1)\top} A \left( \widehat{\beta}_G^{(2)} - \beta_G^{(2)} \right) \\ &\quad - \left( \widehat{\beta}_G^{(1)} - \beta_G^{(1)} \right)^\top A \left( \widehat{\beta}_G^{(2)} - \beta_G^{(2)} \right). \end{aligned}$$

The key step is to estimate the error components  $\widehat{\beta}_G^{(2)\top} A \left( \widehat{\beta}_G^{(1)} - \beta_G^{(1)} \right)$  and  $\widehat{\beta}_G^{(1)\top} A \left( \widehat{\beta}_G^{(2)} - \beta_G^{(2)} \right)$ . Then the following procedures can be interpreted as applying Linear Functional twice on two independent samples. To be specific, we propose the following bias-corrected estimator for  $\beta_G^{(1)\top} A \beta_G^{(2)}$

$$\begin{aligned} \widehat{\beta}_G^{(1)\top} A \widehat{\beta}_G^{(2)} &= \widehat{\beta}_G^{(1)\top} A \widehat{\beta}_G^{(2)} + \widehat{u}_1^\top \frac{1}{n_1} \sum_{i=1}^{n_1} \omega(X_i^{(1)\top} \widehat{\beta}^{(1)}) \left( y_i^{(1)} - f(X_i^{(1)\top} \widehat{\beta}^{(1)}) \right) X_i^{(1)} \\ &\quad + \widehat{u}_2^\top \frac{1}{n_2} \sum_{i=1}^{n_2} \omega(X_i^{(2)\top} \widehat{\beta}^{(2)}) \left( y_i^{(2)} - f(X_i^{(2)\top} \widehat{\beta}^{(2)}) \right) X_i^{(2)}, \end{aligned} \tag{14}$$

with the second term and the third term in right-hand-side of (14) estimating  $-\widehat{\beta}_G^{(2)\top} A \left( \widehat{\beta}_G^{(1)} - \beta_G^{(1)} \right)$  and  $-\widehat{\beta}_G^{(1)\top} A \left( \widehat{\beta}_G^{(2)} - \beta_G^{(2)} \right)$  respectively, where  $\widehat{u}_1$  is the projection direction defined in (9) with  $x_{\text{new}} = (\widehat{\beta}_G^{(2)\top} A, \mathbf{0})^\top$

and  $\widehat{u}_2$  is the projection direction defined in (9) with  $x_{\text{new}} = (\widehat{\beta}_G^{(1)\top} A, \mathbf{0})^\top$ . The corresponding variance of  $\widehat{\beta}_G^{(1)\top} A \widehat{\beta}_G^{(2)}$ , when  $A$  is a known positive definite matrix, is estimated as

$$\widehat{V}_A(\tau) = \widehat{V}^{(1)} + \widehat{V}^{(2)} + \frac{\tau}{\min(n_1, n_2)},$$

where  $\widehat{V}^{(k)}$  is computed as (10) for the  $k$ -th regression model ( $k = 1, 2$ ) in (2) and  $\tau > 0$ , the term  $\tau / \min(n_1, n_2)$  is introduced as an upper bound for the term  $(\widehat{\beta}_G^{(1)} - \beta_G^{(1)})^\top A (\widehat{\beta}_G^{(2)} - \beta_G^{(2)})$ .

When  $A$  is not specified, we treat  $A = \Sigma_{G,G}$ , which is unknown. As a natural generalization, the quantity  $\beta_G^{(1)\top} \Sigma_{G,G} \beta_G^{(2)}$  is well defined if the two regression models in (2) share the design covariance matrix  $\Sigma = \mathbb{E} X_i^{(1)} X_i^{(1)\top} = \mathbb{E} X_i^{(2)} X_i^{(2)\top}$ . We follow the above procedures replacing  $A$  by  $\widehat{\Sigma}_{G,G} = \frac{1}{n_1+n_2} \sum_{i=1}^{n_1+n_2} X_{i,G} X_{i,G}^\top$  where  $X$  is the row-combined matrix of  $X^{(1)}$  and  $X^{(2)}$ . The variance of  $\beta_G^{(1)\top} \Sigma_{G,G} \beta_G^{(2)}$  is now estimated as

$$\widehat{V}_\Sigma(\tau) = \widehat{V}^{(1)} + \widehat{V}^{(2)} + \frac{1}{(n_1 + n_2)^2} \sum_{i=1}^{n_1+n_2} \left( \widehat{\beta}_G^{(1)\top} X_{i,G} X_{i,G}^\top \widehat{\beta}_G^{(2)} - \widehat{\beta}_G^{(1)\top} \widehat{\Sigma}_{G,G} \widehat{\beta}_G^{(2)} \right)^2 + \frac{\tau}{\min(n_1, n_2)}.$$

Depending on whether the submatrix  $A$  is specified or not, the CI is

$$\text{CI}(\tau) = \begin{cases} \left( \widehat{\beta}_G^{(1)\top} A \widehat{\beta}_G^{(2)} - z_{\alpha/2} \widehat{V}_A(\tau), \widehat{\beta}_G^{(1)\top} A \widehat{\beta}_G^{(2)} + z_{\alpha/2} \widehat{V}_A(\tau) \right) & \text{if } A \text{ is specified} \\ \left( \widehat{\beta}_G^{(1)\top} \widehat{\Sigma}_{G,G} \widehat{\beta}_G^{(2)} - z_{\alpha/2} \widehat{V}_\Sigma(\tau), \widehat{\beta}_G^{(1)\top} \widehat{\Sigma}_{G,G} \widehat{\beta}_G^{(2)} + z_{\alpha/2} \widehat{V}_\Sigma(\tau) \right) & \text{otherwise.} \end{cases}$$

## 2.6 Distance of regression vectors

We denote  $\gamma = \beta^{(2)} - \beta^{(1)}$  and its initial estimator  $\widehat{\gamma} = \widehat{\beta}^{(2)} - \widehat{\beta}^{(1)}$ . The quantity of interest is the distance between two regression vectors  $\gamma_G^\top A \gamma_G$ , given a pre-specified submatrix  $A \in \mathbb{R}^{|\mathcal{G}| \times |\mathcal{G}|}$  and the set of indices  $\mathcal{G} \in \{1, \dots, p\}$ . The bias of the plug-in estimator  $\widehat{\gamma}_G^\top A \widehat{\gamma}_G$  is:

$$\widehat{\gamma}_G^\top A \widehat{\gamma}_G - \gamma_G^\top A \gamma_G = 2 \widehat{\gamma}_G^\top A \left( \widehat{\beta}_G^{(2)} - \beta_G^{(2)} \right) - 2 \widehat{\gamma}_G^\top A \left( \widehat{\beta}_G^{(1)} - \beta_G^{(1)} \right) - (\widehat{\gamma}_G - \gamma_G)^\top A (\widehat{\gamma}_G - \gamma_G).$$

The key step is to estimate the error components  $\widehat{\gamma}_G^\top A \left( \widehat{\beta}_G^{(1)} - \beta_G^{(1)} \right)$  and  $\widehat{\gamma}_G^\top A \left( \widehat{\beta}_G^{(2)} - \beta_G^{(2)} \right)$  in the above decomposition. We apply linear functional techniques twice here, and propose the bias-corrected estimator:

$$\begin{aligned} \widehat{\gamma}_G^\top A \gamma_G &= \widehat{\gamma}_G^\top A \widehat{\gamma}_G - 2 \widehat{u}_1^\top \frac{1}{n_1} \sum_{i=1}^{n_1} \omega(X_i^{(1)\top} \widehat{\beta}^{(1)}) \left( y_i^{(1)} - f(X_i^{(1)\top} \widehat{\beta}^{(1)}) \right) X_i^{(1)} \\ &\quad + 2 \widehat{u}_2^\top \frac{1}{n_2} \sum_{i=1}^{n_2} \omega(X_i^{(2)\top} \widehat{\beta}^{(2)}) \left( y_i^{(2)} - f(X_i^{(2)\top} \widehat{\beta}^{(2)}) \right) X_i^{(2)}. \end{aligned} \tag{15}$$

Then by non-negative distance, we define  $\widehat{\gamma}_G^\top A \gamma_G = \max \left\{ \widehat{\gamma}_G^\top A \widehat{\gamma}_G, 0 \right\}$ . The second term on right-hand-side of (15) is to estimate  $-2 x_{\text{new}}^\top (\widehat{\beta}_G^{(1)} - \beta_G^{(1)})$  with  $x_{\text{new}} = (\widehat{\gamma}_G^\top A, \mathbf{0})^\top$ ; and the third term on right-hand-side of (15) is to estimate  $-2 x_{\text{new}}^\top (\widehat{\beta}_G^{(2)} - \beta_G^{(2)})$  with  $x_{\text{new}} = (\widehat{\gamma}_G^\top A, \mathbf{0})^\top$  as well. The corresponding asymptotic variance for the bias-corrected estimator is

$$\widehat{V}_A(\tau) = 4 \widehat{V}^{(1)} + 4 \widehat{V}^{(2)} + \frac{\tau}{\min(n_1, n_2)},$$

where  $\widehat{V}^{(k)}$  is computed as (10) for the  $k$ -th regression model ( $k = 1, 2$ ) and  $\tau > 0$ , the term  $\tau / \min(n_1, n_2)$  is introduced as an upper bound for the term  $(\widehat{\gamma}_G - \gamma_G)^\top A (\widehat{\gamma}_G - \gamma_G)$ . With asymptotic normality, we construct the CI

$$\text{CI}(\tau) = \left( \max \left( \widehat{\gamma}_G^\top A \gamma_G - z_{\alpha/2} \sqrt{\widehat{V}_A(\tau)}, 0 \right), \widehat{\gamma}_G^\top A \gamma_G + z_{\alpha/2} \sqrt{\widehat{V}_A(\tau)} \right).$$

When the submatrix  $A$  is not specified, we treat  $A = \Sigma_{G,G}$ , which is unknown. The point estimator  $\widehat{\gamma_G^\top \Sigma_{G,G} \gamma_G}$  can be computed similarly as outlined in (15). In this case, the submatrix  $A$  is substituted with  $\widehat{\Sigma}_{G,G}$  and the resulting value is truncated at 0, where  $\widehat{\Sigma}_{G,G} = \frac{1}{n_1+n_2} \sum_{i=1}^{n_1+n_2} X_{i,G} X_{i,G}^\top$  with  $X$  as the row-combined matrix of  $X^{(1)}$  and  $X^{(2)}$ . Its corresponding asymptotic variance is

$$\widehat{V}_\Sigma = 4\widehat{V}^{(1)} + 4\widehat{V}^{(2)} + \frac{1}{(n_1+n_2)^2} \sum_{i=1}^{n_1+n_2} \left( \widehat{\gamma}_G^\top X_{i,G} X_{i,G}^\top \widehat{\gamma}_G - \widehat{\gamma}_G^\top \widehat{\Sigma}_{G,G} \widehat{\gamma}_G \right)^2 + \frac{\tau}{\min(n_1, n_2)}.$$

Next we present its CI

$$\text{CI}(\tau) = \left( \max \left( \widehat{\gamma}_G^\top \widehat{\Sigma} \widehat{\gamma}_G - z_{\alpha/2} \sqrt{\widehat{V}_\Sigma(\tau)}, 0 \right), \widehat{\gamma}_G^\top \widehat{\Sigma} \widehat{\gamma}_G + z_{\alpha/2} \sqrt{\widehat{V}_\Sigma(\tau)} \right).$$

## 3 Others

### 3.1 Construction of Projection Direction

The construction of projection directions are key to the bias correction step, see (8). In the following, we introduce the equivalent dual problem of constructing the projection direction. The constrained optimizer  $\widehat{u} \in \mathbb{R}^p$  can be computed in the form of  $\widehat{u} = -\frac{1}{2} \left[ \widehat{\mathbf{v}}_{-1} + \frac{x_*}{\|x_*\|_2} \widehat{\mathbf{v}}_1 \right]$ , where,  $\widehat{\mathbf{v}} \in \mathbb{R}^{p+1}$  is defined as

$$\widehat{\mathbf{v}} = \arg \min_{\mathbf{v} \in \mathbb{R}^{p+1}} \left\{ \frac{1}{4n} \mathbf{v}^\top \mathbf{H}^\top X^\top \text{Diag}(\mathbf{w}) \text{Diag}(\mathbf{f}') X \mathbf{H} \mathbf{v} + x_{\text{new}}^\top \mathbf{H} \mathbf{v} + \lambda_n \|x_{\text{new}}\|_2 \cdot \|\mathbf{v}\|_1 \right\}, \quad (16)$$

with  $\mathbf{H} = \left[ \frac{x_{\text{new}}}{\|x_{\text{new}}\|_2}, \mathbf{I}_{p \times p} \right] \in \mathbb{R}^{p \times (p+1)}$ ,  $\mathbf{w} = \left( \omega(X_1^\top \widehat{\beta}), \dots, \omega(X_n^\top \widehat{\beta}) \right)^\top$  and  $\mathbf{f}' = \left( f'(X_1^\top \widehat{\beta}), \dots, f'(X_n^\top \widehat{\beta}) \right)^\top$ . We refer to Proposition 2 in Cai et al. [2] for the detailed derivation of the dual problem (16). In this dual problem, when  $\widehat{\Sigma}$  is singular and the tuning parameter  $\lambda_n > 0$  gets sufficiently close to 0, the dual problem cannot be solved as the minimum value converges to negative infinity. Hence we choose the smallest  $\lambda_n > 0$  such that the dual problem has a finite minimum value. Such selection of the tuning parameter dated at least back to Javanmard and Montanari [8].

## References

- [1] P. J. Bickel, Y. Ritov, and A. B. Tsybakov. Simultaneous analysis of lasso and dantzig selector. *The Annals of statistics*, 37(4):1705–1732, 2009.
- [2] T. Cai, T. Cai, and Z. Guo. Optimal statistical inference for individualized treatment effects in high-dimensional models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2019.
- [3] T. Cai, T. Tony Cai, and Z. Guo. Optimal statistical inference for individualized treatment effects in high-dimensional models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 83(4):669–719, 2021.
- [4] T. T. Cai, Z. Guo, and R. Ma. Statistical inference for high-dimensional generalized linear models with binary outcomes. *Journal of the American Statistical Association*, pages 1–14, 2021.
- [5] Z. Guo, W. Wang, T. T. Cai, and H. Li. Optimal estimation of genetic relatedness in high-dimensional linear models. *Journal of the American Statistical Association*, 114:358–369, 2019.
- [6] Z. Guo, P. Rakshit, D. S. Herman, and J. Chen. Inference for the case probability in high-dimensional logistic regression. *The Journal of Machine Learning Research*, 22(1):11480–11533, 2021.
- [7] Z. Guo, C. Renaux, P. Bühlmann, and T. Cai. Group inference in high dimensions with applications to hierarchical testing. *Electronic Journal of Statistics*, 15(2):6633–6676, 2021.

- [8] A. Javanmard and A. Montanari. Confidence intervals and hypothesis testing for high-dimensional regression. *The Journal of Machine Learning Research*, 15(1):2869–2909, 2014.
- [9] R. Ma, Z. Guo, T. T. Cai, and H. Li. Statistical inference for genetic relatedness based on high-dimensional logistic regression. *arXiv preprint arXiv:2202.10007*, 2022.
- [10] N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the lasso. *The annals of statistics*, 34(3):1436–1462, 2006.
- [11] S. van de Geer, P. Bühlmann, Y. Ritov, and R. Dezeure. On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42:1166–1202, 2014.
- [12] M. J. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using  $\ell_1$ -constrained quadratic programming (lasso). *IEEE transactions on information theory*, 55(5):2183–2202, 2009.
- [13] C.-H. Zhang and S. S. Zhang. Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):217–242, 2014.
- [14] P. Zhao and B. Yu. On model selection consistency of lasso. *The Journal of Machine Learning Research*, 7:2541–2563, 2006.