

# Test of Significance for High-dimensional Thresholds with Application to Individualized Minimal Clinically Important Difference



Huijie Feng<sup>\*a</sup>, Jingyi Duan<sup>\*a</sup>, Yang Ning<sup>a</sup>, Jiwei Zhao<sup>b</sup>

Department of Statistics and Data Science, Cornell University<sup>a</sup>; Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison<sup>b</sup>

## Introduction

**Motivation** A common and natural approach to assess clinical significance is to determine minimal clinically important difference(MCID). Since the magnitude of MCID would depend on patients' demographic variables, it is of scientific interest to estimate the individualized MCID (iMCID) based on each individual patient's clinical profile, which is typically complex and high-dimensional data.

**Goal** We aim at developing statistical methods that incorporate the high-dimensional data into both **magnitude determination and uncertainty quantification of iMCID**.

### Problem Setup

$X \in \mathbb{R}$ : a **continuous variable** representing the score change collected from the PRO

$Y = \pm 1$ : a **binary variable** derived from the patient's response to the anchor question

$Z \in \mathbb{R}^d$ : patient's clinical profile including demographic variables  $n$  i.i.d samples,  $\{(x_i, y_i, z_i)\}_{i=1}^n$  of  $(X, Y, Z)$ ,  $d \gg n$

The objective function is formulated as

$$\beta^* = \underset{\beta}{\operatorname{argmin}} R(\beta), \text{ where } R(\beta) = \mathbb{E} \left[ w(Y) L_{01} \left\{ Y (X - \beta^T Z) \right\} \right],$$

where  $L_{01}(u) = \frac{1}{2} \{1 - \operatorname{sign}(u)\}$  is the 0-1 loss,  $w(1) = 1/\pi$ ,

$w(-1) = 1/(1 - \pi)$  and  $\pi = \mathbb{P}(Y = 1)$ .

Denote  $\beta^* = (\theta^*, \gamma^{*T})^T$ , where  $\theta^*$  is an arbitrary one-dimensional component of  $\beta^*$ . We start from considering the hypothesis testing problem  $H_0: \theta^* = 0$  versus  $H_1: \theta^* \neq 0$  where we treat  $\gamma^*$  as a high-dimensional nuisance parameter.

## Problem Overview

**Estimation** Under  $d \gg n$ , estimating  $\beta^*$  is **hard**:

- The non-smoothness of  $L_{01}(u)$  would cause the estimator to have a **nonstandard convergence rate**, which happens even in the fixed low dimensional case (Kim et al., 1990).
- Minimizing the empirical risk function based on the 0-1 loss is computationally **NP-hard** and is often very difficult to implement.
- There's no estimators of  $\beta^*$  with root-n convergence rate (Feng et al. 2022).

### Inference

- Non-regular model: **non-Gaussian limiting distribution** (Kim et al., 1990).
- Under high dimensionality, there exists **bias** induced by the penalty term.

## Methodology

**Smoothed Surrogate loss** To tackle the non-smoothness of  $L_{01}(u)$ , Feng et al. 2022 considered the following **smoothed surrogate loss**:

$$R_\delta(\beta) = \mathbb{E} \left[ w(Y) L_{\delta,K} \left\{ Y (X - \beta^T Z) \right\} \right],$$

where  $L_{\delta,K}(u) = \int_{u/\delta}^\infty K(t) dt$  is a smoothed approximation of  $L_{01}(u)$ ,  $K$  is a kernel function and  $\delta > 0$  is a bandwidth parameter.

**Penalized Smoothed Surrogate Estimator**

$$\hat{\beta} := \underset{\beta}{\operatorname{argmin}} R_\delta^n(\beta) + P_\lambda(\beta), \star$$

$P_\lambda(\beta)$ : sparsity induced penalty;  $R_\delta^n(\beta)$ : empirical version of  $R_\delta(\beta)$ ,

$$R_\delta^n(\beta) = \frac{1}{n} \sum_{i=1}^n w(y_i) L_{\delta,K} \left( y_i (x_i - \beta^T z_i) \right).$$

**Bias Corrected Smoothed Decorrelated Score**

the classical score test breaks down:

- When plugging in the estimate of the nuisance parameter  $\gamma(\hat{\gamma})$ , the estimation error from  $\hat{\gamma}$  may become the leading term

**Our solution: bias corrected score**

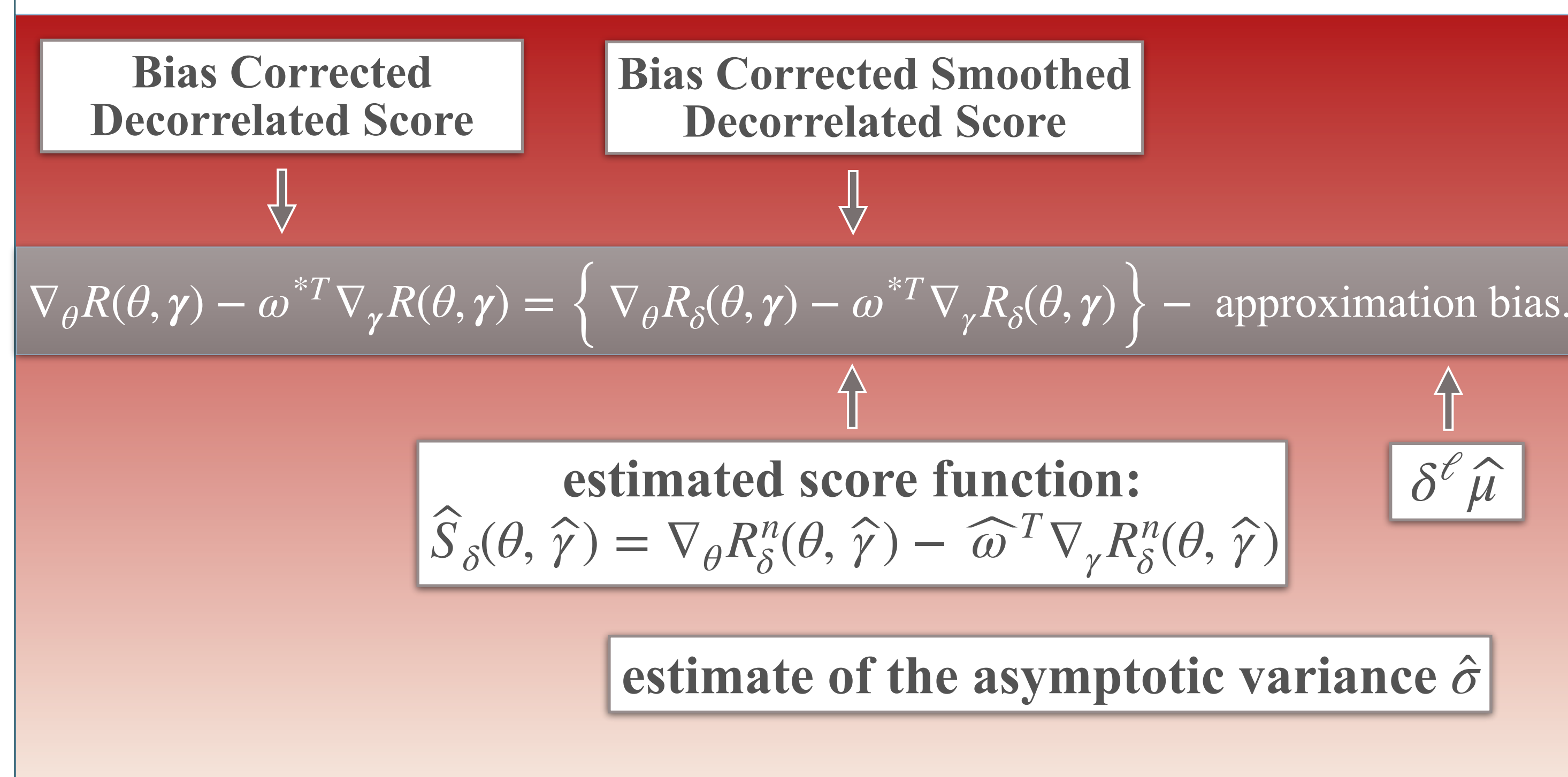
$$\nabla_\theta R(\theta, \gamma) - \omega^{*T} \nabla_\gamma R(\theta, \gamma),$$

where the decorrelation vector is  $\omega^* = \left( \nabla_{\gamma, \gamma}^2 R(\beta^*) \right)^{-1} \nabla_{\gamma, \theta}^2 R(\beta^*)$ .

- Sample version of  $R(\theta, \gamma)$  is non-differentiable

**Our solution: smoothed bias corrected score**

We approximate  $R(\theta, \gamma)$  by  $R_\delta(\theta, \gamma)$ ,  $\gamma$  is estimated by  $\hat{\gamma}$  which is given by  $\star$ .



**The test statistic is defined as**

$$\hat{U}_n = \sqrt{n\delta} \left( \frac{\hat{S}_\delta(0, \hat{\gamma}) - \delta^\ell \hat{\mu}}{\hat{\sigma}} \right)$$

## Theory

**Theorem 1 Asymptotic Normality of The Decorrelated Score Under The Null Hypothesis**

Under  $H_0: \theta^* = 0$ , it holds that

$$\frac{\sqrt{n\delta} \hat{S}_\delta(0, \hat{\gamma}) - \sqrt{n\delta^{2\ell+1}} \mu^*}{\sigma^*} \rightarrow N(0, 1).$$

$\hat{U}_n = \sqrt{n\delta} \left( \frac{\hat{S}_\delta(0, \hat{\gamma}) - \delta^\ell \hat{\mu}}{\hat{\sigma}} \right)$ : plug in the estimators  $\hat{\mu}$  and  $\hat{\sigma}$

$|\hat{\mu} - \mu^*| = o_p(1)$  and  $|\hat{\sigma} - \sigma^*| = o_p(1) \rightarrow \hat{U}_n \rightarrow N(0, 1)$

Test function:  $T_{DS} = I \left( \left| \hat{U}_n \right| > \Phi^{-1}(1 - \alpha/2) \right)$

Type I error  $\mathbb{P}(T_{DS} = 1 | H_0) \rightarrow \alpha$

**Theorem 2 Local Asymptotic Normality**

- When the contiguous alternatives approach the null hypothesis at  $n^{-\frac{\ell}{2\ell+1}}$ , the proposed test  $\rightarrow$  to a normal distribution
- If the alternatives deviate from the null hypothesis faster than  $n^{-\frac{\ell}{2\ell+1}}$ , the asymptotic power of our test is 1

## Simulation

**Model:**  $Y = \operatorname{sign}(X - \beta^{*T} Z + \epsilon)$

- SDS:** the proposed smoothed decorrelated score test
  - DS:** the decorrelated score test method (Ning and Liu, 2017)
  - Honest:** Honest confidence region method (Belloni et al., 2016)
- DS and Honest methods are devised for the logistic regression only

Table 6.1: The empirical Type I error rate of the tests under the Heteroskedastic Gaussian scenario from SDS, DS and Honest methods.

d	method	$s = 3$			$s = 10$		
		$\rho = 0.2$	$\rho = 0.5$	$\rho = 0.7$	$\rho = 0.2$	$\rho = 0.5$	$\rho = 0.7$
100	SDS	5.6%	5.0%	6.4%	4.8%	4.8%	5.2%
	DS	1.2%	2.0%	2.0%	2.0%	1.8%	1.8%
	Honest	5.2%	5.6%	7.6%	5.4%	5.2%	6.8%
500	SDS	4.8%	4.4%	5.6%	5.6%	5.0%	4.8%
	DS	0.2%	0.4%	0.4%	0.2%	0.0%	0.4%
	Honest	7.0%	10.8%	7.6%	8.2%	6.8%	7.2%
1000	SDS	4.4%	6.0%	5.6%	5.0%	5.4%	5.0%
	DS	0.0%	0.4%	0.2%	0.0%	0.0%	0.4%
	Honest	10.0%	10.4%	12.6%	12.4%	6.4%	15.2%

## References

- Feng, H., Ning, Y., and Zhao, J. (2022), "Nonregular and minimax estimation of individualized thresholds in high dimension with binary responses," The Annals of Statistics, 50, 2284–2305.
- Kim, J., Pollard, D., et al. (1990), "Cube root asymptotics," The Annals of Statistics, 18, 191–219.
- Ning, Y. and Liu, H. (2017), "A general theory of hypothesis tests and confidence regions for sparse high dimensional models," The Annals of Statistics, 45, 158–195.
- Belloni, A., Chernozhukov, V., and Wei, Y. (2016), "Post-selection inference for generalized linear models with many controls," Journal of Business & Economic Statistics, 34, 606–619.