

Correction to the paper “Optimal False Discovery Rate Control for Dependent Data”

JICHUN XIE, T. TONY CAI, AND HONGZHE LI*

We have found a mistake in the proof Theorem 6 in our published paper “Optimal False Discovery Rate Control for Dependent Data” [4]. We apologize to the readers and thank Professor Jens Ledet Jensen at Aarhus University for his question which led to identification of this mistake. We provide here a corrected proof of Theorem 6 with further clarifications of the assumptions.

In the GWAS setting that we consider the paper, X_i 's are often the Z -scores with $\text{Var}(X_i) = 1$ for very large sample sizes. We assume that $\sigma_{ii} = 1$. We define the true latent parameter $\theta_{0,i}$: if $\theta_{0,i} = 0$, $X_i \sim N(0, 1)$; and if $\theta_{0,i} = 1$, $X_i \sim N(\mu_i, 1)$. Also, we denote the working latent parameter as θ_i , which is used to define the likelihood ratio $f(X_i | \theta_i = 1)/f(X_i | \theta_i = 0)$ and $f(\mathbf{X} | \theta_i = 1)/f(\mathbf{X} | \theta_i = 0)$.

Assumption (A) can be weakened as following:

Assumption (A'). The non-null proportion p satisfies $m^{-\tau_1} \leq p \leq 1 - m^{-\tau_1}$ for some constant $0 < \tau_1 < 1$.

Let the symbol “ \circ ” be the operator of Hadamard product. Assumption (B) can be clarified as following:

Assumption (B') The data $\mathbf{x}^{(m)} = (x_1, \dots, x_m)$ is an observation of the random variable $\mathbf{X}^{(m)} = (X_1, \dots, X_m)$, which follows a multivariate normal distribution given the mean $\boldsymbol{\mu}^{(m)} \circ \boldsymbol{\theta}^{(m)} = (\mu_1\theta_1, \dots, \mu_m\theta_m)$, *i.e.*

$$\mathbf{X}^{(m)} | \boldsymbol{\mu}^{(m)}, \boldsymbol{\theta}^{(m)} \sim N(\boldsymbol{\mu}^{(m)} \circ \boldsymbol{\theta}^{(m)}, \boldsymbol{\Sigma}^{(m)}).$$

Here $\boldsymbol{\Sigma}^{(m)}$ is the covariance matrix with diagonal elements equal to 1, and $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)^T$ is a random vector, with each μ_i independently following a distribution with CDF $G(\mu)$. Assume for some constant $\tau_2 > \tau_1$,

$$G\{(2\tau_2 \log m)^{1/2}\} - G\{-(2\tau_2 \log m)^{-1/2}\} = 0.$$

It guarantees that

$$\text{pr}\{|\mu_i| \geq (2\tau_2 \log m)^{1/2}\} = 1, \quad i = 1, \dots, m.$$

Remark: Compared to the original Assumption (A), Assumption (A') allows a larger range of p . The condition imposed on the CDF function $G(\mu)$ in Assumption (B') is very weak. For example, consider the case where the non-null proportion is small as $p = m^{-\tau_1}$, for some $1/2 < \tau_1 < 1$

(also known as the sparse case). By [2] and [1], if $|\mu_i| < (2\tau_1 \log m)^{1/2}$, it is not possible to test a single signal with diminishing type I and type II errors. Further, by [3], in order to almost recover all the signals, τ_2 has to be no smaller than $1 + (1 - \tau_1)^2$. Note that $1/2 < \tau_1 < 1$. Therefore, Assumption (B') imposes a weaker condition on the signal strength than what is needed for signal recovery.

Assumption (C) can be weakened as follows:

Assumption (C') The covariance matrix $\boldsymbol{\Sigma}^{(m)}$ is positive definite.

Theorem 6. Under the assumptions (A'), (B') and (C'), define $T_{OR,i}$ and $T_{MG,i}$ as in equation (6) and equation (12). Then for all $\epsilon > 0$ and for all $i = 1, \dots, m$,

$$\lim_{m \rightarrow \infty} \text{pr}(|T_{MG,i} - T_{OR,i}| > \epsilon) = 0.$$

Proof of Theorem 6. We prove the results for $i = 1$. Let $\mathbf{X}_2 = (X_2, \dots, X_m)^T$ be the subvector of the random vector \mathbf{X} without the first variable X_1 . Correspondingly, let $\boldsymbol{\theta}_2 = (\theta_2, \dots, \theta_m)^T$ and $\boldsymbol{\mu}_2 = (\mu_2, \dots, \mu_m)^T$. Define $2\epsilon = \tau_2 - \tau_1 > 0$. Then $\tau_2 - \tau_1 - \epsilon = \epsilon > 0$.

The proof has several steps.

1). We temporarily fix $\boldsymbol{\mu}$ and $\boldsymbol{\theta}_2$. WLOG, assume $\mu_1 \geq (2\tau_2 \log m)^{1/2} > 0$.

1.1) We first consider the case that the true latent variable $\theta_{0,1} = 0$. We show that with probability greater than $1 - O\{(\log m)^{-1/2}\}$,

- (1) $f(X_1 | \theta_1 = 1, \mu_1) < m^{-\tau_2 + \epsilon} \cdot f(X_1 | \theta_1 = 0, \mu_1)$,
- (2) $f(\mathbf{X} | \theta_1 = 1, \boldsymbol{\theta}_2, \boldsymbol{\mu}) < m^{-\tau_2 + \epsilon} \cdot f(\mathbf{X} | \theta_1 = 0, \boldsymbol{\theta}_2, \boldsymbol{\mu})$.

Note that

$$\begin{aligned} \frac{f(X_1 | \theta_1 = 1, \mu_1)}{f(X_1 | \theta_1 = 0, \mu_1)} &= \exp\left\{-\frac{1}{2}(X_1 - \mu_1)^2 + \frac{1}{2}X_1^2\right\} \\ (3) \qquad \qquad \qquad &= \exp\left(\mu_1 X_1 - \frac{1}{2}\mu_1^2\right). \end{aligned}$$

We assume $\theta_{0,1} = 0$, so $X_1 \sim N(0, 1)$ and for sufficiently large m

$$\begin{aligned} \text{pr}(X_1 > (\log \log m)^{1/2}) &= \Phi(-(\log \log m)^{1/2}) \\ &\leq \frac{\varphi((\log \log m)^{1/2})}{(\log \log m)^{1/2}} \end{aligned}$$

*Corresponding author.

$$(4) \quad < (\log m)^{-1/2}.$$

Here $\Phi(\cdot)$ is the cdf of $N(0, 1)$ and $\varphi(\cdot)$ is the corresponding pdf.

Note that $\mu_1 = (2\tau_2 \log m)^{1/2}$. Then for all sufficiently large m , with probability greater than $1 - (\log m)^{-1/2}$,

$$\frac{f(X_1 | \theta_1 = 1, \mu_1)}{f(X_1 | \theta_1 = 0, \mu_1)} < m^{-\tau_2 + \varepsilon}.$$

Let $\Omega = \Sigma^{-1}$ be the precision matrix of \mathbf{X} . Corresponding to the partition $\mathbf{X} = (X_1, \mathbf{X}_2)$, we can write

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \quad \text{and} \quad \Omega = \begin{pmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{21} & \Omega_{22} \end{pmatrix}.$$

where Σ_{11} and Ω_{11} are scalars.

Based on the model, given θ_2 and μ ,

$$\mathbf{X}_2 \sim N(\mu_2 \circ \theta_2, \Sigma_{22}).$$

Conditioning on θ_2 and μ , we have

$$\begin{aligned} & \frac{f(\mathbf{X} | \theta_1 = 1, \theta_2, \mu)}{f(\mathbf{X} | \theta_1 = 0, \theta_2, \mu)} \\ &= \exp \left\{ -\frac{1}{2} (X_1 - \mu_1)^2 \Omega_{11} + (X_1 - \mu_1) \Omega_{12} (\mathbf{X}_2 - \mu_2 \circ \theta_2) \right. \\ & \quad \left. - \frac{1}{2} (\mathbf{X}_2 - \mu_2 \circ \theta_2)^T \Omega_{22} (\mathbf{X}_2 - \mu_2 \circ \theta_2) \right. \end{aligned}$$

$$(5) \quad \left. + \frac{1}{2} X_1^2 \Omega_{11} - X_1 \Omega_{12} (\mathbf{X}_2 - \mu_2 \circ \theta_2) + \frac{1}{2} (\mathbf{X}_2 - \mu_2 \circ \theta_2)^T \Omega_{22} (\mathbf{X}_2 - \mu_2 \circ \theta_2) \right\}$$

$$(6) \quad = \exp \left\{ X_1 \mu_1 \Omega_{11} - \frac{1}{2} \mu_1^2 \Omega_{11} - \mu_1 \Omega_{11} \frac{\Omega_{12} (\mathbf{X}_2 - \mu_2 \circ \theta_2)}{\Omega_{11}} \right\}$$

Let $Z_2 = \Omega_{12} (\mathbf{X}_2 - \mu_2 \circ \theta_2) / \Omega_{11}$. Then

$$Z_2 | (\theta_2, \mu_2) \sim N(0, \Omega_{12} \Sigma_{22} \Omega_{21} / \Omega_{11}^2).$$

We now show that the variance of Z_2 is upper bounded by 1.

By the equality

$$\Sigma \Omega = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{21} & \Omega_{22} \end{pmatrix} = \begin{pmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{I} \end{pmatrix},$$

we have

$$(7) \quad \Sigma_{12} \Omega_{21} = 1 - \Sigma_{11} \Omega_{11}.$$

By another equality $\Omega \Sigma \Omega = \Omega$ with the same partition, we have

$$(8) \quad \Omega_{11}^2 \Sigma_{11} + 2\Omega_{11} \Sigma_{12} \Omega_{21} + \Omega_{12} \Sigma_{22} \Omega_{21} = \Omega_{11}.$$

By (7) and (8),

$$(9) \quad 0 \leq \Omega_{12} \Sigma_{22} \Omega_{21} / \Omega_{11}^2 = (\Sigma_{11} \Omega_{11} - 1) / \Omega_{11} = (\Omega_{11} - 1) / \Omega_{11} < 1.$$

Also by (9), we have $\Omega_{11} \geq 1$, and

$$(10) \quad \text{pr}(Z_2 < -(\log \log m)^{1/2} | \theta_2, \mu) < \Phi(-(\log \log m)^{1/2}) < (\log m)^{-1/2}.$$

By (4) and (10), given θ_2 and μ , with probability greater than $1 - 2(\log m)^{-1/2}$,

$$\frac{f(\mathbf{X} | \theta_1 = 1, \theta_2, \mu)}{f(\mathbf{X} | \theta_1 = 0, \theta_2, \mu)} < \exp\{2(\log \log m)^{1/2} \mu_1 \Omega_{11} - \mu_1^2 \Omega_{11} / 2\} < m^{-\tau_2 + \varepsilon},$$

where the second inequality is due to $\Omega_{11} \geq 1$ and $\mu_1 = (2\tau_2 \log m)^{1/2}$. This implies (2).

1. ii) We now turn to the case where the true latent variable $\theta_{0,1} = 1$. We show that with probability greater than $1 - 2\{(\log m)^{-1/2}\}$,

$$(11) \quad f(X_1 | \theta_1 = 1, \mu_1) > c_2 m^{\tau_2 - \varepsilon} \cdot f(X_1 | \theta_1 = 0, \mu_1),$$

$$(12) \quad f(\mathbf{X} | \theta_1 = 1, \theta_2, \mu) > c_2 m^{\tau_2 - \varepsilon} \cdot f(\mathbf{X} | \theta_1 = 0, \theta_2, \mu).$$

Since now $\theta_{0,1} = 1$, $X_1 \sim N(\mu_1, 1)$,

$$(13) \quad \text{pr}(X_1 - \mu_1 < -(\log \log m)^{1/2}) = \Phi(-(\log \log m)^{1/2}) < (\log m)^{-1/2}.$$

By (3) and (13), with probability greater than $1 - (\log m)^{-1/2}$,

$$\frac{f(X_1 | \theta_1 = 1, \mu_1)}{f(X_1 | \theta_1 = 0, \mu_1)} = \exp\{\mu_1(x_1 - \mu_1) + \mu_1^2 / 2\} > m^{\tau_2 - \varepsilon}.$$

In addition, given θ_2 and μ_2 ,

$$(14) \quad \text{pr}(Z_2 > (\log \log m)^{1/2} | \theta_2, \mu_2) = (\log m)^{-1/2}.$$

By (6), (13) and (14) and Assumption (C'), we can follow the proof of Step (1.i) and show that with probability greater than $1 - 2(\log m)^{-1/2}$,

$$\frac{f(\mathbf{X} | \theta_1 = 1, \theta_2, \mu)}{f(\mathbf{X} | \theta_1 = 0, \theta_2, \mu)} > m^{\tau_2 - \varepsilon}.$$

2) Now consider θ_2 and μ as random vectors.

When $\theta_{0,1} = 0$, for each given θ_2 and μ ,

$$f(\mathbf{X} | \theta_1 = 1, \theta_2, \mu) < m^{-\tau_2 + \varepsilon} \cdot f(\mathbf{X} | \theta_1 = 0, \theta_2, \mu)$$

holds with probability greater than $1 - 2(\log m)^{-1/2}$. Therefore, with probability greater than $1 - 2(\log m)^{-1/2}$, we have

$$\begin{aligned} & \sum_{\theta_2} f(\mathbf{X} \mid \theta_1 = 1, \theta_2, \boldsymbol{\mu}) \text{pr}(\theta_2) \\ & < m^{-\tau_2 + \varepsilon} \cdot \sum_{\theta_2} f(\mathbf{X} \mid \theta_1 = 0, \theta_2, \boldsymbol{\mu}) \text{pr}(\theta_2). \end{aligned}$$

We conclude that

$$(15) \quad \frac{f(\mathbf{X} \mid \theta_1 = 1, \boldsymbol{\mu})}{f(\mathbf{X} \mid \theta_1 = 0, \boldsymbol{\mu})} = \frac{\sum_{\theta_2} f(\mathbf{X} \mid \theta_1 = 1, \theta_2, \boldsymbol{\mu}) \text{pr}(\theta_2)}{\sum_{\theta_2} f(\mathbf{X} \mid \theta_1 = 0, \theta_2, \boldsymbol{\mu}) \text{pr}(\theta_2)} < m^{-\tau_2 + \varepsilon}$$

holds with probability greater than $1 - 2(\log m)^{-1/2}$.

By (1), (15) and the following equality

$$\begin{aligned} \frac{f(X_1 \mid \theta_1 = 1)}{f(X_1 \mid \theta_1 = 0)} &= \frac{\int f(X_1 \mid \theta_1 = 1, \boldsymbol{\mu}) dG(\boldsymbol{\mu})}{\int f(X_1 \mid \theta_1 = 0, \boldsymbol{\mu}) dG(\boldsymbol{\mu})}, \\ \frac{f(\mathbf{X} \mid \theta_1 = 1)}{f(\mathbf{X} \mid \theta_1 = 0)} &= \frac{\int f(\mathbf{X} \mid \theta_1 = 1, \boldsymbol{\mu}) dG(\boldsymbol{\mu})}{\int f(\mathbf{X} \mid \theta_1 = 0, \boldsymbol{\mu}) dG(\boldsymbol{\mu})}, \end{aligned}$$

we have when $\theta_{0,1} = 0$, with probability greater than $1 - 2(\log m)^{-1/2}$,

$$(16) \quad \frac{f(X_1 \mid \theta_1 = 1)}{f(X_1 \mid \theta_1 = 0)} < m^{-\tau_2 + \varepsilon} \quad \text{and} \quad \frac{f(\mathbf{X} \mid \theta_1 = 1)}{f(\mathbf{X} \mid \theta_1 = 0)} < m^{-\tau_2 + \varepsilon}$$

Similarly, when $\theta_{0,i} = 1$, we can show that with probability greater than $1 - 2(\log m)^{-1/2}$,

$$(17) \quad \frac{f(X_1 \mid \theta_1 = 1)}{f(X_1 \mid \theta_1 = 0)} > m^{\tau_2 - \varepsilon} \quad \text{and} \quad \frac{f(\mathbf{X} \mid \theta_1 = 1)}{f(\mathbf{X} \mid \theta_1 = 0)} > m^{\tau_2 - \varepsilon}$$

3) We are now ready to prove Theorem 6.

It is easy to show

$$\begin{aligned} T_{MG,1} &= \frac{1-p}{(1-p) + pf(X_1 \mid \theta_1 = 0)/f(X_1 \mid \theta_1 = 0)} \\ T_{OR,1} &= \frac{1-p}{(1-p) + pf(\mathbf{X} \mid \theta_1 = 0)/f(\mathbf{X} \mid \theta_1 = 0)} \end{aligned}$$

By Assumption (A'), $cm^{-\tau_1} \leq p/(1-p) \leq cm^{\tau_1}$.

When $\theta_{0,1} = 0$, with probability greater than $1 - O((\log m)^{-1/2})$,

$$\begin{aligned} T_{MG,1} &\geq \frac{(1-p)}{1-p+pm^{-\tau_2+\varepsilon}} \\ &\geq \frac{1}{1+pm^{-\tau_2+\varepsilon}/(1-p)} \geq \frac{1}{1+cm^{\tau_1-\tau_2+\varepsilon}}, \end{aligned}$$

which yields

$$1 - T_{MG,1} \leq \frac{cm^{-(\tau_2-\tau_1-\varepsilon)}}{1+cm^{-(\tau_2-\tau_1-\varepsilon)}}.$$

Similarly, it can be shown that this result holds for $T_{OR,1}$. By Assumption (A'), with probability greater than

$$1 - O((\log m)^{-1/2}),$$

$$\begin{aligned} |T_{OR,1} - T_{MG,1}| &\leq |1 - T_{OR,1}| + |1 - T_{MG,1}| \\ &= \frac{2cm^{-(\tau_2-\tau_1-\varepsilon)}}{1+cm^{-(\tau_2-\tau_1-\varepsilon)}} = O(m^{-(\tau_2-\tau_1-\varepsilon)}) \rightarrow 0. \end{aligned}$$

When $\theta_{0,1} = 1$, with probability greater than $1 - O((\log m)^{-1/2})$,

$$\begin{aligned} T_{MG,1} &\leq \frac{(1-p)}{1-p+pc_2m^{\tau_2-\varepsilon}} \\ &\leq \frac{1}{1+pm^{\tau_2-\varepsilon}/(1-p)} \leq \frac{1}{1+cm^{\tau_2-\tau_1-\varepsilon}}. \end{aligned}$$

Same result holds for $T_{OR,1}$. By Assumption (A'), with probability greater than $1 - O((\log m)^{-1/2})$,

$$\begin{aligned} |T_{OR,1} - T_{MG,1}| &\leq \frac{2}{1+cm^{\tau_2-\tau_1-\varepsilon}} \\ &= O(m^{-(\tau_2-\tau_1-\varepsilon)}) \rightarrow 0. \quad \square \end{aligned}$$

Received 21 August 2014

REFERENCES

- [1] JI, P. and JIN, J. (2012). UPS delivers optimal phase diagram in high-dimensional variable selection. *The Annals of Statistics* **40-3** 73–103. [MR3013180](#)
- [2] JIN, J. and DONOHO, D. (2004). Higher criticism for detecting sparse heterogeneous mixtures. *The Annals of Statistics* **32-3** 962–994. [MR2065195](#)
- [3] XIE, J., CAI, T. and LI, H. (2011). Sample size and power analysis for sparse signal recovery in genome-wide association studies. *Biometrika* **98-2** 273–290. [MR2806428](#)
- [4] XIE, J., CAI, T. T., MARIS, J. and LI, H. (2011). Optimal false discovery rate control for dependent data. *Statistics and Its Interface* **4** 417–430. [MR2868825](#)

Jichun Xie

Department of Biostatistics and Bioinformatics
School of Medicine
Duke University
USA

E-mail address: jichun.xie@duke.edu

T. Tony Cai

Department of Statistics
The Wharton School
University of Pennsylvania
USA

E-mail address: tc@wharton.upenn.edu

Hongzhe Li

Department of Biostatistics and Epidemiology
School of Medicine
University of Pennsylvania
USA

E-mail address: hongzhe@upenn.edu