

# Supplementary Note: Controlling bias and inflation in epigenome- and transcriptome-wide association studies using the empirical null distribution

Maarten van Iterson<sup>1</sup>, Erik van Zwet<sup>2</sup>, the BIOS Consortium<sup>3</sup>, P. Eline Slagboom<sup>1</sup>, and Bastiaan T. Heijmans<sup>1</sup>

<sup>1</sup>Leiden University Medical Center, Department of Molecular Epidemiology, Leiden, The Netherlands.

<sup>2</sup>Leiden University Medical Center, Department of Medical Statistics and Bioinformatics, Leiden, The Netherlands.

<sup>3</sup>the BIOS consortium

May 25, 2016

## Contents

<b>1</b>	<b>The BIOS consortium</b>	<b>1</b>
<b>2</b>	<b>The genomic inflation factor is affected by true associations</b>	<b>2</b>
<b>3</b>	<b>Applying genomic control is the same as using an inflated or overdispersed empirical null</b>	<b>3</b>
<b>4</b>	<b>Unobserved covariates introduce inflation and bias</b>	<b>3</b>
<b>5</b>	<b>Partial least squares for imputation of white blood cell composition</b>	<b>4</b>
<b>6</b>	<b>Different ways to calculate genomic inflation</b>	<b>4</b>

## 1 The BIOS consortium

The mission of the BIOS Consortium is to create a large-scale data infrastructure and to bring together BBMRI researchers focusing on integrative omics studies in Dutch Biobanks (<https://www.bbmri.nl/?p=259>).

**Management Team:** Bastiaan T. Heijmans (chair)<sup>1</sup>, Peter A.C. 't Hoen<sup>2</sup>, Joyce van Meurs<sup>3</sup>, Rick Jansen<sup>5</sup>, Lude Franke<sup>6</sup>.

**Cohort collection:** Dorret I. Boomsma<sup>7</sup>, René Pool<sup>7</sup>, Jenny van Dongen<sup>7</sup>, Jouke J. Hottenga<sup>7</sup> (Netherlands Twin Register); Marleen MJ van Greevenbroek<sup>8</sup>, Coen D.A. Stehouwer<sup>8</sup>, Carla J.H. van der Kallen<sup>8</sup>, Casper G. Schalkwijk<sup>8</sup> (Cohort study on Diabetes and Atherosclerosis Maastricht); Cisca Wijmenga<sup>6</sup>, Lude Franke<sup>6</sup>, Sasha Zhernakova<sup>6</sup>, Ettje F. Tigchelaar<sup>6</sup> (LifeLines Deep); P. Eline Slagboom<sup>1</sup>, Marian Beekman<sup>1</sup>, Joris Deelen<sup>1</sup>, Diana van Heemst<sup>9</sup> (Leiden Longevity Study); Jan H. Veldink<sup>10</sup>, Leonard H. van den Berg<sup>10</sup> (Prospective ALS

Study Netherlands); Cornelia M. van Duijn<sup>4</sup>, Bert A. Hofman<sup>11</sup>, Aaron Isaacs<sup>4</sup>, André G. Uitterlinden<sup>3</sup> (Rotterdam Study).

**Data Generation:** Joyce van Meurs (Chair)<sup>3</sup>, P. Mila Jhamai<sup>3</sup>, Michael Verbiest<sup>3</sup>, H. Eka D. Suchiman<sup>1</sup>, Marlijn Verkerk<sup>3</sup>, Ruud van der Breggen<sup>1</sup>, Jeroen van Rooij<sup>3</sup>, Nico Lakenberg<sup>1</sup>.

**Data management and computational infrastructure:** Hailiang Mei (Chair)<sup>12</sup>, Maarten van Iterson<sup>1</sup>, Michiel van Galen<sup>2</sup>, Jan Bot<sup>13</sup>, Dasha V. Zhernakova<sup>5</sup>, Rick Jansen<sup>4</sup>, Peter van 't Hof<sup>12</sup>, Patrick Deelen<sup>5</sup>, Irene Nooren<sup>13</sup>, Peter A.C. 't Hoen<sup>2</sup>, Bastiaan T. Heijmans<sup>1</sup>, Matthijs Moed<sup>1</sup>.

**Data Analysis Group:** Lude Franke (Co-Chair)<sup>6</sup>, Martijn Vermaat<sup>2</sup>, Dasha V. Zhernakova<sup>6</sup>, René Luijk<sup>1</sup>, Marc Jan Bonder<sup>6</sup>, Maarten van Iterson<sup>1</sup>, Patrick Deelen<sup>6</sup>, Freerk van Dijk<sup>14</sup>, Michiel van Galen<sup>2</sup>, Wibowo Arindrarto<sup>12</sup>, Szymon M. Kielbasa<sup>15</sup>, Morris A. Swertz<sup>14</sup>, Erik. W van Zwet<sup>15</sup>, Rick Jansen<sup>5</sup>, Peter-Bram 't Hoen (Co-Chair)<sup>2</sup>, Bastiaan T. Heijmans (Co-Chair)<sup>1</sup>.

1. Molecular Epidemiology Section, Department of Medical Statistics and Bioinformatics, Leiden University Medical Center, Leiden, The Netherlands
2. Department of Human Genetics, Leiden University Medical Center, Leiden, The Netherlands
3. Department of Internal Medicine, ErasmusMC, Rotterdam, The Netherlands
4. Department of Genetic Epidemiology, ErasmusMC, Rotterdam, The Netherlands
5. Department of Psychiatry, VU University Medical Center, Neuroscience Campus Amsterdam, Amsterdam, The Netherlands
6. Department of Genetics, University of Groningen, University Medical Centre Groningen, Groningen, The Netherlands
7. Department of Biological Psychology, VU University Amsterdam, Neuroscience Campus Amsterdam, Amsterdam, The Netherlands
8. Department of Internal Medicine and School for Cardiovascular Diseases (CARIM), Maastricht University Medical Center, Maastricht, The Netherlands
9. Department of Gerontology and Geriatrics, Leiden University Medical Center, Leiden, The Netherlands
10. Department of Neurology, Brain Center Rudolf Magnus, University Medical Center Utrecht, Utrecht, The Netherlands
11. Department of Epidemiology, ErasmusMC, Rotterdam, The Netherlands
12. Sequence Analysis Support Core, Leiden University Medical Center, Leiden, The Netherlands
13. SURFsara, Amsterdam, the Netherlands
14. Genomics Coordination Center, University Medical Center Groningen, University of Groningen, Groningen, the Netherlands
15. Medical Statistics Section, Department of Medical Statistics and Bioinformatics, Leiden University Medical Center, Leiden, The Netherlands

## 2 The genomic inflation factor is affect by true associations

Given a set of test-statistics  $z_2, \dots, z_p$  the (squared) genomic inflation factor is given by[1]:

$$\lambda^2 = \frac{\text{med}\{z_2^2, \dots, z_p^2\}}{0.457}. \quad (1)$$

The median of the squared test-statistics will be the ordered test-statistic at position  $p/2$  or  $(p+2)/2$ , if  $p$  is odd or even, respectively. Since, the set of test-statistics represents  $p_1$  test-statistics following the null distribution and  $p - p_1$  the alternative. The set ordered test-statistics will be given by  $\{z_2^2, \dots, z_{p_1}^2, z_{p_1+2}^2, \dots, z_p^2\}$ . Furthermore, it is known in advance that  $p_1 > p/2$  or  $p_1 > (p+1)/2$  it follows that  $\text{med}\{z_2^2, \dots, z_p^2\} > \text{med}\{z_2^2, \dots, z_{p_1}^2\} > 0.457$  and thus  $\lambda^2 > 1$ .

### 3 Applying genomic control is the same as using an inflated or overdispersed empirical null

Consider a two-sided test for normally distributed test-statistics,  $T_i$  for  $i = 1, \dots, p$ . Genomic control divides test-statistics by the inflation factor,  $\lambda$ , before the calculation of  $P$  values,  $U_i$ .

$$\begin{aligned}
 U_i &= 2 \left[ 1 - \Phi \left( \left| \frac{T_i}{\lambda} \right| \right) \right] \\
 &= 2 \left[ 1 - \Pr \left\{ Z \leq \left| \frac{T_i}{\lambda} \right| \right\} \right] \\
 &= 2 [1 - \Pr \{ \lambda Z \leq |T_i| \}] \\
 &= 2 [1 - \Pr \{ X \leq |T_i| \}]
 \end{aligned} \tag{2}$$

Here,  $Z \sim \mathcal{N}(0, 1)$  with CDF given by  $\Phi$ . Furthermore,  $X \sim \mathcal{N}(0, \lambda^2)$ , represents the overdispersed or inflation normal distribution.

And in general if both inflation and bias are present.

$$\begin{aligned}
 U_i &= 2 \left[ 1 - \Phi \left( \left| \frac{T_i - \mu}{\sigma} \right| \right) \right] \\
 &= 2 [1 - \Pr \{ \sigma Z + \mu \leq |T_i| \}] \\
 &= 2 [1 - \Pr \{ X \leq |T_i| \}]
 \end{aligned} \tag{3}$$

Here,  $\mu$ , represents the bias and  $\sigma$  the inflation. Furthermore, now  $X \sim \mathcal{N}(\mu, \sigma^2)$ .

### 4 Unobserved covariates introduce inflation and bias

In a note from P. Rao[2] it was shown that the omission of a variable introduces bias and decreases the variance of all least squares estimates, i.e. introduces both bias and inflation. Here is a sketch of the proof for the introduction of bias, within the framework of the main text.

Considered the omission of the unobserved technical or biological covariates,  $\mathbf{W}$ . For the sake of simplicity, we assume there are no known covariates,  $\mathbf{Z}$ , without loss of generality.

$$\begin{aligned}
 \mathbf{y}_j &= \mathbf{x}\tilde{\beta}_j + \tilde{\epsilon}_j \\
 \mathbf{y}_j &= \mathbf{x}\beta_j + \mathbf{W}\gamma_j + \epsilon_j
 \end{aligned} \tag{4}$$

The latter model is true but we are unaware of this and continue estimating the regression coefficient of interest,  $\tilde{\beta}_j$  of the former, misspecified, model.

$$\begin{aligned}
 \tilde{b}_j &= \frac{\mathbf{x}^T \mathbf{y}_j}{\mathbf{x}^T \mathbf{x}} \\
 &= \frac{\mathbf{x}^T (\mathbf{x}\beta_j + \mathbf{W}\gamma_j + \epsilon_j)}{\mathbf{x}^T \mathbf{x}} \\
 &= \beta_j + \frac{\mathbf{x}^T \mathbf{W}\gamma_j + \mathbf{x}^T \epsilon_j}{\mathbf{x}^T \mathbf{x}}
 \end{aligned} \tag{5}$$

where, we substituted the true model for  $\mathbf{y}_j$ . Now, since  $E[\epsilon_j] = 0$ , the expected regression coefficient is given by:

$$E[\tilde{b}_j] = \beta_j + \frac{\mathbf{x}^T \mathbf{W}}{\mathbf{x}^T \mathbf{x}} \gamma_j \tag{6}$$

and the bias is given by the last fraction, which can be interpreted as a weighted sum of correlations between the covariate of interest,  $\mathbf{x}$  and  $q$  omitted variables.

$$\sum_k^q \gamma_{jk} \text{cor}(\mathbf{x}, \mathbf{w}_k) \quad (7)$$

If all weights are zero or all correlations than the bias will be zero to. The bias is not equal to zero if an omitted variable is confounded with the outcome, i.e. both  $\gamma_{jk}$  and  $\mathbf{x}^T \mathbf{W}$  are not zero.

## 5 Partial least squares for imputation of white blood cell composition

White blood cell counts (WBC), i.e., neutrophils, lymphocytes, monocytes, eosinophils and basophils, were measured by the standard WBC differential as part of the CBC (Complete Blood Count). However, a minority of samples lack CBC measurements. Since DNA methylation levels are informative of the white blood cell composition[3] we build a linear predictor to infer the white blood cell composition of those samples lacking WBC measurements.

Obviously, this model can not be fitted using ordinary least-squares, since  $p \gg n$ , we need some kind of regularization. Furthermore, the multivariate responses, white blood counts on five cell types, represents compositional data, i.e., the data are percentages that sum up to 100%.

We have chosen to use partial least-squares for fitting a model with cell counts as a multivariate response and the  $> 400.000$  CpGs age and sex as covariates. It is known that the WBCC is dependent on age and gender. The advantage of partial least-squares is that it both can handle multivariate responses and high-dimensional ( $p \gg n$ ) covariates. We used the R-package pls[4] to fit the model and optimize the number of pls-components using five-fold cross-validation. The fitted model was used to predict WBCC using the 450K data, age and sex of those samples lacking WBCC.

The pls-approach has been validated by splitting the data with WBCC available in a train and test set. Fit the pls-model on the train set and predict WBCC on the test set. Correlation (Pearson) between predicted and measured WBCC range from 0.86 – 0.37 for lymphocytes and basophils respectively, the intraclass correlation was 0.84 – 0.25.

This approach has been implemented in a R package <https://github.com/mvaniterson/wbccPredictor>.

## 6 Different ways to calculate genomic inflation

Devlin and Roeder [1] propose the genomic inflation factor is as the ratio of the median of  $\chi_1^2$ -distributed test-statistics divided by the median of a  $\chi_1^2$  of 0.456.

$$\lambda = \frac{\text{median}(t_1, t_2, \dots, t_p)}{0.456}, \quad (8)$$

where  $t_i \sim \chi_2^2$  at least under the null hypothesis.

Since, in EWAS/TWAS usually test-statistics,  $z_1, z_2, \dots, z_p$  are derived that follow (approximately) a normal distribution the following formula can be used to estimate the amount of inflation:

$$\lambda^2 = \frac{\text{median}(z_1^2, z_2^2, \dots, z_p^2)}{0.456}, \quad (9)$$

since,  $z_i^2$  will approximately follow a  $\chi_1^2$ -distribution. Furthermore, the inflation of the z-scale is the square-root of the  $\chi_2^2$ -scale to indicate this we introduced the  $\lambda^2$ .

Another way to estimate the amount of inflation is using the absolute value of z-score and divide them by  $0.456^2$ .

Occasionally,  $P$  values are used to estimate the amount of inflation by comparing median of the minus  $\log_{10}$ -transformed  $P$  values with the median of minus  $\log_{10}$ -transformed random uniformly distributed statistics. However, a simple calculation shows these follow an exponential distribution with rate parameter  $\log_e 10$  and thus the theoretical median is  $\log_{10} 2$ .

## References

- [1] B. Devlin and K. Roeder. Genomic control for association studies. *Biometrics*, 55(4), 1999.
- [2] P. Rao. Some notes on misspecification in multiple regression. *The American Statistician*, 25(5), 1971.
- [3] E.A. Houseman et al. DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC Bioinformatics*, 8(13), 2012.
- [4] Wehrens R. Mevik, B.-H. The pls Package: Principal Component and Partial Least Squares Regression in R. *Journal of Statistical Software*, 18(2), 2007.