

# Methods for variance estimation of high dimensional data

*Xuelong Wang*

*2020-01-29*

## Contents

<b>1</b>	<b>Motivation</b>	<b>1</b>
<b>2</b>	<b>Linear regression</b>	<b>1</b>
<b>3</b>	<b>GCTA method</b>	<b>1</b>
3.1	GCTA approach . . . . .	1
3.2	The proposed method . . . . .	3
3.3	Interactive effect . . . . .	5
3.4	available software . . . . .	6
<b>4</b>	<b>EigenPrism</b>	<b>6</b>
4.1	Model assumption . . . . .	6
4.2	Simulation results . . . . .	6
<b>5</b>	<b>Variance estimation in high-dimensional linear models</b>	<b>7</b>
5.1	Model assumption . . . . .	7
5.2	Signal Estimation for $\Sigma = I$ . . . . .	7
5.3	Simulation results . . . . .	8
	<b>References</b>	<b>8</b>

## 1 Motivation

## 2 Linear regression

## 3 GCTA method

### 3.1 GCTA approach

The GCTA approach estimates variances of weak effects. . .

#### 3.1.1 Model assumption

GCTA approach is built on a linear mix model(LMM):

$$Y_i = \mu_i + \sum_{j=1}^p X_{ij}\beta_j + \epsilon_i. \quad (1)$$

where  $Y_i$  denotes a outcome (quantitative measurement) and  $x_{ij}, j = 1, \dots, p$  are the standardized covariates measurements for subject  $i$ . Besides we also assume the independence between the covariates and error terms,  $\epsilon_i \perp x_{jk}$ . The equation 1 may be re-expressed as

$$Y_i = \mu + X_i^T \beta + \epsilon_i. \quad (2)$$

where  $X_i$  is a  $p \times 1$  vector,  $Y = (y_1, \dots, y_n)^T$ ,  $\mu = (\mu_1, \dots, \mu_n)^T$ ,  $\beta = (\beta_1, \dots, \beta_p)^T$ , and  $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T$ .

The goal here is to estimate how much variation of the outcome is related to the covariates when  $n \ll p$ . Based on the assumptions of model 1, the variance of  $y$  could be composed into two parts,

$$Var(Y_i) = Var(X_i^T \beta) + Var(\epsilon).$$

Based on the different assumptions of the randomness of  $X_i$  and  $\beta$ , we will have different format for the  $Var(X_i^T \beta)$ . But no matter in what situation, we will have

$$\begin{aligned} Var(X_i^T \beta) &= E(X_i^T \beta (X_i^T \beta)^T) - E(X_i^T \beta) E(X_i^T \beta)^T \\ &= E(X_i^T \beta (X_i^T \beta)^T) - E(X_i^T \beta) E(X_i^T \beta)^T \text{ b/c } E(X_i) = 0 \\ &= E(X_i^T \beta (X_i^T \beta)^T) - E(X_i^T) E(\beta) (E(X_i^T) E(\beta))^T \\ &= E(X_i^T \beta (X_i^T \beta)^T) - 0 \text{ b/c } E(X_i) = 0 \text{ or } E(\beta) = 0 \\ &= E(X_i^T \beta (X_i^T \beta)^T) \end{aligned}$$

One particular situation we are interested in is that  $X$  is random with  $E(X_i) = 0$ ,  $Var(X_i) = \Sigma$  and  $\beta$  could be fixed effect or random effects. For example, if  $\beta$  is fixed then

$$Var(X_i^T \beta) = E(X_i^T \beta (X_i^T \beta)^T) = E(\text{trace}(X_i^T \beta (X_i^T \beta)^T)) = \beta^T \Sigma \beta,$$

If we further assume that  $\Sigma = I_p$ , then

$$\sigma_g^2 = \beta^T \Sigma \beta = \sum_i \beta_i^2.$$

Then if we also assume  $X_i$ 's are uncorrelated random vectors, then we will have

$$Var(X\beta) = E(X\beta\beta^T X^T),$$

$$Var(X\beta)_{ii} = E(X_i^T B X_i) = \text{trac}(B I_p) = \sum_i \beta_i^2 \quad Var(X\beta)_{ij} = 0$$

So we have

$$Var(X\beta) = \sum_i \beta_i^2 I_n.$$

The next question is how to estimate  $\sum_i \beta_i^2$ . In original paper of GCTA, it assumes that  $X$ 's are fixed and  $\beta \sim N(0_p, \sigma_g^2/p I_p)$ , and we have

$$Y = \mu + X\beta + \epsilon.$$

Then the

$$Var(Y) = Var(X\beta) + \sigma_\epsilon^2 I_n = X X^T / p \sigma_g^2 I_p + \sigma_\epsilon^2 I_n = K \sigma_g^2 + \sigma_\epsilon^2 I_n$$

, where  $K = X X^T / p$ . Let the narrow-sense of the heritability ratio as

$$h^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma_\epsilon^2}.$$

Besides, (Jiang et al. 2016) provided more rigorous decription and proof of the LMM and MLMM. In the paper, it assume that  $X$  is also a random matrix and it is independnt with  $\beta$  and  $\epsilon$ . It also mentioned that LMM is a condntional model, on  $\mu$  and  $X$ , so no matter  $X$  is nonradom.

A restricted maximum likelihood method (REML) method is used to estimate the  $\hat{\sigma}_g^2$  and  $\hat{\sigma}_e^2$ . It turns out to be that GCTA method can also estimate the  $Var(X_i\beta)$ . But GCTA will still assume the  $K\sigma_g^2$  and try to estimate  $\sigma_g^2$  under our setting up, and we have

$$K = 1/p(XX^T), \quad K_{ii} = \frac{\sum_k X_{ik}^2}{p} \rightarrow 1 \text{ and } K_{ij} = \frac{\sum_j X_{iK}X_{jK}}{p} \rightarrow 0 \text{ as } p \rightarrow \infty??$$

So we could use the  $\hat{\sigma}_g^2$  to estimate  $\sum_j \hat{\beta}_j^2??$

### 3.1.2 Advantages of the GCTA apporach

For the environmental study mentioned before, the GCTA method demonstrates some advantages than other variance estimation approaches.

- a working random effects model to estimate  $Var(X\beta)$
- Don't need to select the casual covariates, so that could work with weak signal problem
- relatively little bias compared to other methods

### 3.1.3 Two more Obstacles

Although we discussed the GCTA approach could be a good tool for the environmental health analysis, there are still issues we need to tackle.

- Theoretical analysis of The GCTA approach suggests the Independence of causal covariates, but most of the environmental data are high correlated.
- In SNP studies, the number of covariates is large and the number of interactive terms is also going to be very large, which makes the interactive effect even harder to be estimated. Therefore, interaction effect usually is not considered in SNPs studies. Although in environmental studies the number of predictors is not large (within 40), directly applying the GCTA method to estimate the interactive effect still hardly guarantee good performance.

## 3.2 The proposed method

With those two problems in mind, we develop a new method by modifying the GCTA method for correlated covariates. The main idea is to transform the correlated covariates into uncorrelated ones. The transformation process is also called decorrelation. We consider a linear transformation so that the transformation does not change the variance structure.

### 3.2.1 Transformation for correlated covariates

The linear transformation is

$$Z = A^{-1}X,$$

where  $X$  are the covariates vector,  $A$  is a linear transformation operator which is a full rank square matrix. After transformation, the covariance of the new covariates  $Z$  will be

$$Var(Z) = I_p.$$

Moreover, based on the model assumed by GCTA (model 2), we have

$$Y = \mu + X^T \beta + \epsilon = Z^T A^T \beta + \epsilon = Z^T \alpha + \epsilon,$$

where  $\alpha = A^T \beta$ . Let's look the total effect of \*X and \*Z:

$$Var(X^T \beta) = Var(Z^T A^T \beta) = Var(Z^T \alpha).$$

Therefore, the  $Z$  will be the uncorrelated predictors and  $Z^T \alpha$  should keep the same total cumulative effect as  $X^T \beta$ . If  $X$  follows a normal distribution, i.e.  $X \sim N(0, \Sigma)$ , then the  $Z \sim N(0, I_p)$ . Therefore,  $Z$ 's elements are independent to each other, which is the exact condition we want for the GCTA approach. Although for non-normal covariates the decorrelation procedure only reduces linear association with no guarantee of independence, it still can improve the performance of GCTA method.

### 3.2.2 Decorrelation procedure

There are many methods and algorithms for data decorrelation. One of commonly used methods is to apply the eigenvalue decomposition to the covariance matrix. Let  $\Sigma_X$  be the covariance matrix of  $X$ , so  $\Sigma_X$  is a symmetric and positive-definite. Then eigenvalue decomposition of  $\Sigma_X$  will be

$$Var(X) = \Sigma_X = U \Lambda U^T,$$

where  $X$  is the random vector with dim as  $p \times 1$ ,  $\Sigma_X$  is  $p \times p$  symmetry and p.d. matrix,  $\Lambda$  is a diagonal matrix with each diagonal element as the eigenvalue. If the  $\Sigma_X$  is full rank, then we could just take the reciprocal of each square root of eigenvalue as following.

$$\Sigma_X^{-\frac{1}{2}} = U \Lambda^{-\frac{1}{2}} U^T, \text{ and } \Lambda^{-\frac{1}{2}} = \begin{bmatrix} \lambda_1^{-\frac{1}{2}} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_p^{-\frac{1}{2}} \end{bmatrix}.$$

So that after transformation, the  $\Sigma_X^{-\frac{1}{2}} X$  will have an identity covariance matrix as following:

$$Var(\Sigma_X^{-\frac{1}{2}} X) = \Sigma_X^{-\frac{1}{2}} \Sigma_X \Sigma_X^{-\frac{1}{2}} = U \Lambda^{-\frac{1}{2}} U^T U \Lambda^{-1} U^T U \Lambda^{-\frac{1}{2}} U^T = I_p.$$

If the  $\Sigma_X$  is not full rank, then we can still use the eigenvalue decomposition. But the procedure cannot guarantee the identity covariance matrix anymore. The reason is that some eigenvalues will be zero, so we can not take the reciprocal. One straightforward solution is just leave them there:

$$Var(X) = \Sigma_X = U \Lambda U^T = [U_1 \quad U_2] \begin{bmatrix} \Lambda_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} = U_1 \Lambda_1 U_1^T,$$

where  $U_1$  is a  $p \times r$  matrix with  $r < p$  and in most of case  $r = n$  the sample size. Then after applying the same procedure we get following,

$$\tilde{\Sigma}_X^{-\frac{1}{2}} = U_1 \Lambda_1^{-\frac{1}{2}} U_1^T,$$

Where  $\tilde{\Sigma}_X^{-1}$  is the Moore Penrose inverse. After transformation the  $X$  we have,

$$Var(\tilde{\Sigma}_X^{-\frac{1}{2}} X) = \tilde{\Sigma}_X^{-\frac{1}{2}} \Sigma_X \tilde{\Sigma}_X^{-\frac{1}{2}} = U_1 \Lambda_1^{-\frac{1}{2}} U_1^T = U_1 U_1^T,$$

Where  $U_1 U_1^T + U_2 U_2^T = I_p$  and  $(U_1 U_1^T)^T U_1 U_1^T = U_1 U_1^T$ . Besides,  $U_1 U_1^T$  and  $U_2 U_2^T$  are idempotent and  $rank(U_2 U_2^T) + rank(U_1 U_1^T) = p$ .

### 3.3 Interactive effect

For analysing the interactive effect, we need to consider interactive terms in our model. Let's just consider a 2-way interaction model

$$Y_i = \mu_i + \sum_{j=1}^p x_{ij}\beta_j + \sum_{l \neq k} \gamma_{lk}x_{il}x_{ik} + \epsilon_i, \quad (3)$$

where  $\gamma_{jk}$  denotes interactive coefficients. Anything else will be same as GCTA model 2. This model also can be expressed in the matrix form

$$Y = \mu + X\beta + X\Gamma X^T + \epsilon. \quad (4)$$

Where  $\Gamma$  is a  $p \times p$  matrix with element as  $\gamma_{jk}$ . Let's also assume that  $X_i \perp \epsilon_i$ , then the variance of  $Y_i$  can be decompose as following

$$Var(Y_i) = Var(X_i^T \beta) + Var(X_i^T \Gamma X_i) + 2Cov(X_i^T \beta, X_i^T \Gamma X_i) + Var(\epsilon_i). \quad (5)$$

After adding the interactive terms, the situation becomes complicated.

1. Besides the interactive effect there is an additional covariance term of  $X_i^T \beta, X_i^T \Gamma X_i$  to deal with.
2. The main and interactive terms are bonded to be dependent, even though all elements  $X$  are independent. Same situation for the 2-way interactive terms, they are also dependent.

As we mentioned before, independence of covariates is suggested for GCTA approach to work well, so we cannot guarantee the performance of GCTA approach.

To handle the covariance terms, we now focus on the situations where  $Cov(X_i^T \beta, X_i^T \Gamma X_i) = 0$ , so that we don't have to worry about it. For the cases where we cannot ignore the covariance term, it's hard to estimate both of the effects well. The reason is that the covariance term will be somehow mixed into both main and interactive effect estimations, so it is not easy to separate covariance part from the effects' estimation. We will discuss it latter in this paper. Let's just assume that covariates are independent and centered to each other and there is no square terms in the model 3, e.i.  $\gamma_{jj} = 0$ ,

$$\begin{aligned} Cov(X^T \beta, X^T \Gamma X) &= E[(X^T \beta - E(X^T \beta))(X^T \Gamma X - E(X^T \Gamma X))] \\ &= E[X^T \beta (X^T \Gamma X - E(X^T \Gamma X))] \\ &= E[X^T \beta \cdot X^T \Gamma X] \\ &= E[(\sum_h^p (x_h \beta_h))(\sum_j^p \sum_k^p \gamma_{jk} x_j x_k)] \\ &= 0 \end{aligned}$$

For the second issues, we extend our proposed approach to handle the interactive terms. Although it's impossible to make the interactive terms independent with themselves or the main terms, we still can transform them into uncorrelated. Therefore, we could combine the main and interactive term together as a larger covariate matrix

$$X_{i,t} = \begin{bmatrix} X_i \\ X_{i,inter} \end{bmatrix},$$

where  $X_{i,inter} = (x_{i1}x_{i2}, \dots, x_{i(p-1)}x_{ip})^T$ . Then apply the decorrelation process on the combined matrix  $X_t = (X_{1t}, \dots, X_{nt})^T$ . Given the independence of the covariates, simulation studies have shown that the

proposed method could estimate both of the cumulative and interactive effect with little bias. This also suggests that the uncorrelation of covariates may be good enough to let the GCTA works appropriately. Therefore, we may release the condition from independent to uncorrelated covariates.

### 3.4 available software

## 4 EigenPrism

### 4.1 Model assumption

### 4.2 Simulation results

#### 4.2.1 $p = 100$ , $\sum \beta^2 = 8$ and $\rho = 0.1-0.9$

n	rho_e	MSE	est_mean	est_var_m	est_var	NA_main	N	relative_ratio	relative_ratio_var
50	0.1	0.11	0.16	0.35	0.10	0	1000	2.4	0.00
50	0.3	0.09	0.34	0.35	0.09	0	1000	2.9	0.00
50	0.5	0.07	0.52	0.35	0.07	0	1000	3.9	0.00
50	0.7	0.05	0.71	0.35	0.05	0	1000	5.5	0.00
50	0.9	0.03	0.90	0.35	0.03	0	1000	9.0	0.01
75	0.1	0.06	0.14	0.26	0.06	0	1000	3.3	0.00
75	0.3	0.05	0.33	0.26	0.05	0	1000	4.5	0.00
75	0.5	0.03	0.52	0.26	0.03	0	1000	6.7	0.00
75	0.7	0.02	0.71	0.26	0.02	0	1000	11.4	0.00
75	0.9	0.01	0.91	0.26	0.01	0	1000	24.0	0.02
100	0.1	0.04	0.12	0.21	0.04	0	1000	4.5	0.00
100	0.3	0.03	0.32	0.21	0.03	0	1000	6.1	0.00
100	0.5	0.02	0.51	0.21	0.02	0	1000	9.5	0.00
100	0.7	0.01	0.71	0.21	0.01	0	1000	18.1	0.00
100	0.9	0.00	0.91	0.21	0.00	0	1000	48.8	0.02
150	0.1	NaN	NaN	NaN	NA	1000	1000	NA	NA
150	0.3	NaN	NaN	NaN	NA	1000	1000	NA	NA
150	0.5	NaN	NaN	NaN	NA	1000	1000	NA	NA
150	0.7	NaN	NaN	NaN	NA	1000	1000	NA	NA
150	0.9	NaN	NaN	NaN	NA	1000	1000	NA	NA

#### 4.2.2 $p = 500$ , $\sum \beta^2 = 8$ and $\rho = 0.1-0.9$

n	rho_e	MSE	est_mean	est_var_m	est_var	NA_main	N	relative_ratio	relative_ratio_var
50	0.1	0.30	0.26	0.60	0.27	1	1000	NA	NA
50	0.3	0.28	0.42	0.60	0.26	1	1000	NA	NA
50	0.5	0.27	0.59	0.60	0.26	1	1000	NA	NA
75	0.1	0.17	0.20	0.43	0.15	0	1000	1.8	0
75	0.3	0.16	0.38	0.43	0.15	0	1000	1.9	0
75	0.5	0.15	0.55	0.43	0.14	0	1000	2.0	0
100	0.1	0.10	0.18	0.34	0.09	0	1000	2.7	0
100	0.3	0.09	0.36	0.34	0.09	0	1000	3.0	0
100	0.5	0.08	0.54	0.34	0.08	0	1000	3.3	0
150	0.1	0.05	0.15	0.24	0.05	0	1000	3.9	0
150	0.3	0.04	0.34	0.24	0.04	0	1000	4.5	0
150	0.5	0.04	0.52	0.24	0.04	0	1000	5.4	0
250	0.1	0.02	0.12	0.16	0.02	0	1000	5.8	0
250	0.3	0.02	0.31	0.16	0.02	0	1000	7.3	0
250	0.5	0.01	0.51	0.16	0.01	0	949	9.7	0
500	0.1	0.01	0.10	0.10	0.01	0	1000	10.4	0
500	0.3	0.01	0.30	0.10	0.01	0	1000	14.3	0

## 5 Variance estimation in high-dimensional linear models

### 5.1 Model assumption

### 5.2 Signal Estimation for $\Sigma = I$

$$E\left(\frac{1}{n}\|y\|^2\right) = \tau^2 + \sigma^2, \quad E\left(\frac{1}{n^2}\|X^T y\|^2\right) = \frac{d+n+1}{n}\tau^2 + \frac{d}{n}\sigma^2$$

After some linear algebra, we have the corresponding estimator is

$$\hat{\sigma}^2 = \frac{d+n+1}{n(n+1)}\|y\|^2 - \frac{1}{n(n+1)}\|X^T y\|^2, \quad \hat{\tau}^2 = -\frac{d}{n(n+1)}\|y\|^2 + \frac{1}{n(n+1)}\|X^T y\|^2$$

Under some standard condition the estimators have asymptotic normality.

$$\begin{aligned} \psi_1^2 &= 2 \left\{ \frac{d}{n} (\sigma^2 + \tau^2)^2 + \sigma^4 + \tau^4 \right\} \\ \psi_2^2 &= 2 \left\{ \left(1 + \frac{d}{n}\right) (\sigma^2 + \tau^2)^2 - \sigma^4 + 3\tau^4 \right\} \\ \psi_0^2 &= \frac{2}{(\sigma^2 + \tau^2)^2} \left\{ \left(1 + \frac{d}{n}\right) (\sigma^2 + \tau^2)^2 - \sigma^4 \right\} \end{aligned}$$

If  $d/n \rightarrow \rho \in [0, \infty)$ , then

$$n^{1/2} \left( \frac{\hat{\sigma}^2 - \sigma^2}{\psi_1} \right), n^{1/2} \left( \frac{\hat{\tau}^2 - \tau^2}{\psi_2} \right), n^{1/2} \left( \frac{\hat{\tau}^2 - \tau^2}{\psi_0} \right) \rightarrow N(0, 1) \text{ in distribution.}$$

### 5.3 Simulation results

#### 5.3.1 $p = 500$ , $\sum \beta^2 = 1$ and $\rho = 0.5$

n	MSE	est_mean	est_var_m	est_var	NA_main
250	0.11	0.98	0.11	0.11	0
500	0.04	1.00	0.04	0.04	0

#### 5.3.2 $p = 500$ , $\sum \beta^2 = 1$ and $\rho = 0.1-0.9$

n	rho_e	MSE	est_mean	est_var_m	est_var	NA_main	N	relative_ratio	relative_ratio_var
250	0.1	0.85	0.09	0.02	0.02	0	100	-0.14	0.01
250	0.3	0.54	0.28	0.02	0.02	0	100	-0.20	0.01
250	0.5	0.30	0.48	0.02	0.03	0	100	-0.23	0.00
250	0.7	0.13	0.68	0.02	0.03	0	100	-0.25	0.00
250	0.9	0.05	0.88	0.02	0.03	0	100	-0.24	0.00
500	0.1	0.81	0.10	0.00	0.00	0	100	0.12	0.01
500	0.3	0.49	0.30	0.01	0.01	0	100	0.07	0.01
500	0.5	0.25	0.50	0.01	0.01	0	100	0.02	0.00
500	0.7	0.10	0.70	0.01	0.01	0	100	-0.01	0.00
500	0.9	0.02	0.90	0.01	0.01	0	100	-0.04	0.00

## References

Jiang, Jiming, Cong Li, Debashis Paul, Can Yang, Hongyu Zhao, and others. 2016. “On High-Dimensional Misspecified Mixed Model Analysis in Genome-Wide Association Study.” *The Annals of Statistics* 44 (5). Institute of Mathematical Statistics: 2127–60.